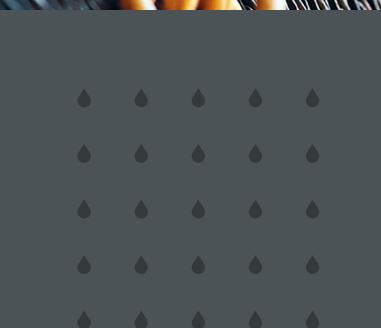


# СИСТЕМА ИНТЕГРИРОВАННОГО ПЛАНИРОВАНИЯ РАСТЕНИЕВОДСТВА

ООО «Ростовская зерновая компания «Ресурс»



## Содержание

|  |   |
|--|---|
| 1. Описание функциональных характеристик системы по модулям .....                    | 3 |
| 2. Информация об условиях использования программного обеспечения .....               | 6 |
| 3. Установка необходимого программного обеспечения (Инструкция по инсталляции) ..... | 6 |

## 1. Описание функциональных характеристик системы по модулям

- 1.1. Система интегрированного планирования Растениеводства (далее по тексту – СИПР) предназначена для автоматизированной поддержки бизнес-процессов долгосрочного и производственного планирования сезона в сегменте растениеводства. Функционал системы включает:
  - 1.1.1. Проведение оптимизационного моделирования в части:
    - модели оптимального размещения культур (структуры посевных площадей);
    - алгоритмов формирования плана посевных работ;
    - проведения моделирования балансировки сельскохозяйственной техники.
  - 1.1.2. Единый интерфейс планирования производства, основанный на ведении технологических карт культур.
  - 1.1.3. Формирование производственных задач.
  - 1.1.4. Отслеживание выполнения задач в реальном времени, с возможностью просмотров отчетов за прошлые периоды.
  - 1.1.5. Интеграцию с производственным оборудованием для получения/передачи информации.
  - 1.1.6. Инструменты оркестрации для запуска и мониторинга работы алгоритмов в среде оптимизации.
  - 1.1.7. Проведение расчетов для формирования показателей отчетности.
  - 1.1.8. Набор механизмов проверки входных данных (сценария) перед запуском моделирования:
    - полноты данных;
    - корректности данных;
    - логических проверок.
  - 1.1.9. Набор пользовательских интерфейсов для выполнения шагов планирования в соответствии с бизнес-ролью пользователя:
    - области планирования и формы ввода данных для подготовки входных данных;
    - пульт управления интегрированным планированием для запуска процесса моделирования и мониторинга процессов расчета;
    - формирование отчетности и анализ результатов;
    - факторный анализ и сравнение сценариев.
  - 1.1.10. Механизмы интеграции с внешними системами (системы оперативного контроля, корпоративные хранилища данных).

Рисунок № 1. Функциональная архитектура решения СИПР



1.2. Формирование производственной программы – единая платформа интегрированного планирования для поддержки принятия стратегических решений и согласованного с ними оперативного управления. В основе решения лежит создание оптимизационной платформы для автоматизации сквозных расчетов, учитывающих как процессы стратегического планирования, так и процессы производственного планирования сезона, с поддержкой возможности анализа различных сценариев развития Растениеводства.

Подход к автоматизации задач долгосрочного планирования и планирования производственной программы на сезон представляет собой использование единого вычислительного графа, состоящего из набора взаимосвязанных оптимизационных моделей, выполняемых по цепочке. При этом результат работы одной модели является входными данными для работы следующей по цепочке модели.

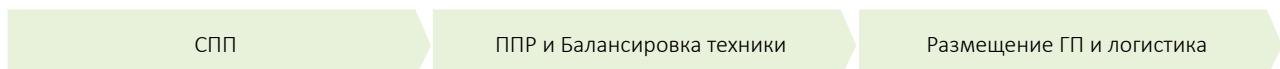
Общая схема связей и порядка выполнения оптимизационных моделей в вычислительном графе приведена на Схеме № 1, где:

**СПП** – Структура посевных площадей

**ППР** – План полевых работ

**ГП** – Готовая продукция

Схема № 1. Связь и порядок выполнения оптимизационных моделей



Автоматизированный подход поддержки процессов планирования включает:

1.2.1. Формирование сценария для моделирования:

- входные данные;
- параметры сценария.

1.2.2. Подготовку и расчет сценария.

1.2.3. Анализ результатов.

Процесс планирования производственной программы выполняется в несколько шагов. Каждый шаг – это сценарий модели конкретного вида (СПП, ППР и пр.). Иерархическая последовательность сценариев формирует цепочку оптимизационных моделей, результатом которой является вариант производственного плана.

Пользователи СИПР формируют несколько вариантов производственного плана, сравнивают их, анализируют финансовый показатели и выбирают план, который утверждается в качестве «производственной программы». Утвержденная производственная программа может быть передана в систему оперативного контроля или сезонного планирования растениеводства для реализации.

1.3. Оптимизационное моделирование – в основе цепочки моделей находится оптимизационная модель СПП, задача которой спланировать посевы на полях. Результатом моделирования СПП являются следующие данные:

- 1.3.1. структура посевных площадей;
- 1.3.2. плановая урожайность культур на полях;
- 1.3.3. оценка потребности в наемной уборочной технике с учетом имеющегося парка сельскохозяйственной техники и оборудования.

Результаты моделирования СПП используются в качестве входных данных для оптимизационной модели ППР. В результате моделирования ППР определяются сроки и ресурсы для работ по выращиванию культур и графику уборки урожая. Реализация вычислительного графа поддерживает возможность пересчета модели ППР без пересчета модели СПП.

Результаты моделирования ППР будут использоваться в качестве входных данных для следующей по цепочке группы оптимизационных моделей по размещению готовой продукции, задача которой на основе уточненного плана полевых работ сформировать план размещения продукции на местах хранения и соответствующую логистику от полей до мест размещения.

Подготовка к запуску процесса моделирования заключается в подготовке входных данных, задании параметров сценария, а также в выполнении проверки качества входных данных.

Вне зависимости от того, запускал ли пользователь проверку качества данных вручную, при запуске расчета модели (процесса оптимизации) предварительным шагом в обязательном порядке автоматически запускается проверка данных. При наличии критических ошибок процесс останавливается, уведомляя пользователя о перечне замечаний к данным. Пользователь обязан устранить критические ошибки. Некритичные замечания к данным не обязательны к устранению, пользователь может запустить расчет без их устранения.

Если проверка качества данных не выявила критичных замечаний, то осуществляется запуск оптимизационного расчета, при этом сценарий переходит в статус «В процессе расчета». Во время расчета ведется журнал расчета, в котором отражается текущая стадия расчета. По завершении расчета сценарий переходит в один из статусов: «Расчет завершен успешно. Решение найдено», «Расчет завершен успешно. Решение не найдено», «Расчет завершен неуспешно».

1.4. Аналитическая отчетность – в системе реализуется возможность формирования следующих отчетов:

- 1.4.1. производственные отчеты (выполняется на базе результатов расчета соответствующих оптимизационных моделей);
- 1.4.2. экономические отчеты (выполняется на базе результатов расчета всей цепочки моделей);
- 1.4.3. отчеты по сравнению базового сценария и сценария с инвестициями в технику и места хранения (выполняется на базе результатов расчетов всей цепочки моделей двух сценариев).

Производственные отчеты формируются по мере выполнения каждого из блоков моделей. В состав производственных отчетов системы входят следующие:

- Структура посевных площадей;
- Сводная производственная программа;
- План работ сельскохозяйственной техники в формате отчета и гант-диаграммы;
- Отчет по потребности в технике для сезонов долгосрочного планирования (с оценкой дефицита техники);
- План потребности в персонале (с оценкой дефицита);
- План размещения собственной и наемной техники по филиалам;
- План перемещения техники между филиалами;
- Сводная таблица по валовому сбору для очереди уборки;
- Схема размещения готовой продукции (с оценкой дефицита техники);
- Расчет потребности в грузовом автотранспорте для уборочной кампании.

1.5. Уровни доступа к информационным ресурсам – встроенные средства системы позволяют ограничить доступ пользователей СИПР к информационным ресурсам как на уровне доступа к компонентам, так и на уровне объектов полномочий каждого компонента системы. Объектами полномочий в компонентах системы являются:

1.5.1. Компонент Dashboards:

- аналитические отчеты;
- датасеты.

1.5.2. Компонент формы ввода:

- пользовательские формы;
- сценарии запуска оптимизации;
- датасеты.

1.5.3. Компонент Analytics:

- OLAP-кубы.

1.5.4. Компонент ML:

- оптимизационные модели.

Доступ пользователей к информационным ресурсам СИПР предполагает выполнение процедур аутентификации и авторизации. Процедуры аутентификации и авторизации пользователей выполняются в СИПР с использованием решения для управления идентификацией и доступами Keycloak с возможностью подключения к службе Active Directory.

Представляемые на данном уровне полномочия обеспечивают доступ пользователей СИПР непосредственно к информации, обрабатываемой и хранящейся в СИПР, а также к средствам автоматизированной обработки этой информации.

Полномочия на доступ к функциональности и информационным ресурсам Системы определяются ролевой моделью доступа. Ролевая модель обеспечивает разграничение доступа пользователей к информационным ресурсам.

Полномочия пользователей СИПР определяются с учетом:

- типов пользователей системы;
- уровней доступа к информационным ресурсам системы посредством прикладного программного обеспечения (далее по тексту – ПО);
- разделения функциональных обязанностей.

Управление полномочиями доступа пользователей (предоставление, изменение, аннулирование) обеспечивается средствами прикладного ПО Polymatica в интерфейсе пользователя с правами администратора.

## **2. Информация об условиях использования программного обеспечения**

- 2.1. Стоимость использования программного обеспечения на условиях неисключительной бессрочной лицензии для одного юридического лица – 15 000 000 (пятнадцать миллионов) рублей 00 копеек.
- 2.2. Для полнофункциональной работы Системы интегрированного планирования Растениеводства требуется наличие лицензий на ПО Полиматика.

## **3. Установка необходимого программного обеспечения (Инструкция по инсталляции)**

- 3.1. Для работы СИПР необходима установка СУБД PostgreSQL а также распределенной системы хранения SeaweedFS. Установку рекомендуется выполнять на один или несколько выделенных серверов.

Установка СУБД PostgreSQL выполняется согласно инструкции для соответствующего дистрибутива, расположенной в сети интернет по адресу <https://www.postgresql.org/download/>.

Последовательность установки SeaweedFS представлена и доступна на ресурсе по адресу <https://github.com/chrislusf/seaweedfs/wiki/Getting-Started>.

После установки указанного программного обеспечения, для работы системы требуется запустить сервисы **master**, **filer**, **s3**, **volume**, предварительно сконфигурировав SeaweedFS.

- 3.2. Конфигурация SeaweedFS – настроить роль для **S3 API** с полными правами, создав файл **/etc/seaweedfs/s3-config.json** со следующим содержимым:

```
{  
  "identities": [  
    {  
      "name": "anonymous",  
      "actions": [  
        "Read",  
        "List"  
      ]  
    },  
    {  
      "name": "admin",  
      "credentials": [  
        {  
          "accessKey": "qwerty",  
          "secretKey": "qwerty"  
        }  
      ],  
      "actions": [  
        "Admin",  
        "Read",  
        "List",  
        "Tagging",  
        "Write"  
      ]  
    }  
  ]  
}
```

Данную конфигурацию следует указывать при запуске сервиса S3 через параметр **config:-config=/etc/seaweedfs/s3-config.json**. Также необходимо создать необходимые бакеты хранения S3, выполнив на сервере SeaweedFS команды:

```
echo "s3.bucket.create-name mm" | weed shell  
echo "s3.bucket.create-name md" | weed shell  
echo "s3.bucket.create-name cs" | weed shell  
echo "s3.bucket.create-name cnf" | weed shell
```

- 3.3. Установка компонента Polymatica (конфигурация Polymatica ML). Перед установкой следует указать основные параметры развертывания в файле **charts/polymatica-ml/values.yaml** в следующих разделах:

#### 3.3.1. Раздел global

- RegistryUrl – адрес для доступа к docker registry в виде <ip-адрес или DNS-имя сервера **polymatica ml>:5000**
- DomainName – ip-адрес или DNS-имя для доступа к серверу
- GatewayUrl – адрес для доступа к публикуемым сервисам в виде **http://<ip-адрес или DNS-имя сервера polymatica ml>/**
- MetadataServerUrl – адрес для доступа к сервису метаданных в виде **http://<ip-адрес или**

### **DNS-имя сервера *polymatica ml*>**

- Подраздел postgres:
  - ✓ host – адрес или имя хоста PostgreSQL
  - ✓ port – порт PostgreSQL
  - ✓ user – пользователь для подключения к БД PostgreSQL
  - ✓ password – пароль для подключения к БД PostgreSQL
  - ✓ databaseMeta – имя БД сервиса метаданных в PostgreSQL
  - ✓ databaseDD – имя БД сервиса datadiscovery в PostgreSQL
  - ✓ databaseMD – имя БД сервиса modeldesigner в PostgreSQL
- Подраздел s3:
  - ✓ url – адрес сервиса S3 API в развернутой SeawedFS в виде **http://<ip-адрес** или DNS-имя сервера **straweedfs:<порт>**

### 3.3.2. Раздел modelmanager

- Подраздел postgres
  - ✓ user – пользователь для подключения к БД PostgreSQL
  - ✓ password – пароль для подключения к БД PostgreSQL
  - ✓ databaseMM – имя БД сервиса modelmanager в PostgreSQL
- Подраздел smtp
  - ✓ username – пользователь для подключения к SMTP-серверу
  - ✓ password – пароль для подключения к SMTP-серверу
  - ✓ from – адрес отправителя
  - ✓ port – порт SMTP-сервера
  - ✓ server – ip-адрес или DNS-имя SMTP-сервера
- Подраздел camunda
  - ✓ jdbcUrl – строка подключения к БД camunda в виде **jdbc:postgresql://<ip-адрес** или DNS-имя сервера **PostgreSQL:<порт>/cs\_dev1?sslmode=disable**
  - ✓ jdbcUser – пользователь для подключения к БД PostgreSQL
  - ✓ jdbcPassword – пароль для подключения к БД PostgreSQL
  - ✓ keycloakUrl – адрес для доступа к сервису keycloak в виде **http://<ip-адрес или DNS-имя сервера polymatica ml>**
- Подраздел modellauncher
  - ✓ devpi\_host – адрес для доступа к репозиторию devpi в виде **http://<ip-адрес** или DNS-имя сервера **polymatica ml:3141**

### 3.3.3. Раздел keycloak

- Подраздел postgres
  - ✓ user – пользователь для подключения к БД PostgreSQL
  - ✓ password – пароль для подключения к БД PostgreSQL
  - ✓ databaseName – имя БД сервиса keycloak в PostgreSQL

3.4. Установку Polymatica ML следует выполнять под непrivилегированным пользователем с правами sudo. Для установки запустить install.sh из каталога с дистрибутивом: DOMAINNAME=<ip-адрес или DNS-имя сервера polymatica-ml> ./install.sh В переменной DOMAINNAME указывается сетевой адрес или DNS-имя сервера, которые будут использоваться для доступа к серверу Polymatica ML. После установки рекомендуется перезапустить службу k3s: sudo systemctl restart k3s

# СИСТЕМА ИНТЕГРИРОВАННОГО ПЛАНИРОВАНИЯ РАСТЕНИЕВОДСТВА

Руководство пользователя



# ОГЛАВЛЕНИЕ

|   |    |
|---|----|
| 1. Общая часть .....  | 5  |
| 1.1. Архитектура .....                                      | 5  |
| 1.2. Интерфейс .....  | 7  |
| 1.2.1. Авторизация .....                                    | 7  |
| 1.2.2. Основное меню .....                                  | 7  |
| 1.2.3. Роли пользователей .....                             | 9  |
| 2. Компонент Исследование данных (Data Discovery, DD) ..... | 11 |
| 2.1. Главный экран .....                                    | 11 |
| 2.1.1. Интерфейс главного экрана DD.....                    | 11 |
| 2.1.2. Создание исследования.....                           | 13 |
| 2.1.3. Удаление исследования.....                           | 14 |
| 2.2. Мастер настройки .....                                 | 14 |
| 2.2.1. Интерфейс Мастера настройки .....                    | 14 |
| 2.2.2. Параметры .....                                      | 14 |
| 2.2.3. Роли атрибутов .....                                 | 15 |
| 2.2.4. Визуализации .....                                   | 16 |
| 2.2.5. Корреляция .....                                     | 29 |
| 2.2.6. Кластеризация.....                                   | 30 |
| 2.2.7. Статистические тесты .....                           | 33 |
| 2.3. Расписание .....                                       | 36 |
| 2.4. Результаты исследования.....                           | 38 |
| 2.4.1. Интерфейс экрана Результаты Исследования .....       | 38 |
| 2.4.2. Профилирование .....                                 | 41 |
| 2.4.3. Результат расчета корреляции .....                   | 45 |
| 2.4.4. Результат расчета кластеризации.....                 | 46 |
| 2.4.5. Результат расчета статистических тестов.....         | 50 |
| 2.4.6. Автоматическое построение модели в MD .....          | 51 |
| 3. Компонент Построение моделей (Model Designer, MD) .....  | 54 |
| 3.1. Главный экран MD .....                                 | 54 |
| 3.1.1. Интерфейс главного экрана MD .....                   | 54 |
| 3.1.2. Создание проекта MD .....                            | 55 |
| 3.1.3. Удаление проекта MD .....                            | 56 |
| 3.1.4. Копирование проекта MD .....                         | 56 |
| 3.1.5. Редактирование проекта MD .....                      | 56 |
| 3.1.6. Работа с проектом MD .....                           | 56 |

|  |     |
|--|-----|
| 3.2. Конструктор сценариев .....                                 | 56  |
| 3.2.1. Интерфейс конструктора сценариев .....                    | 57  |
| 3.2.2. Логирование.....  | 58  |
| 3.2.3. Создание сценария .....                                   | 59  |
| 3.2.4. Создание нескольких сценариев .....                       | 59  |
| 3.2.5. Узлы .....  | 59  |
| 3.2.6. Кросс-валидация.....                                      | 259 |
| 3.2.7. Автоподбор гиперпараметров.....                           | 259 |
| 3.2.8. Результаты моделирования.....                             | 260 |
| 3.3. Пример базового сценария.....                               | 261 |
| 4. Компонент Разработка решений (Decision Manager, DM) .....     | 264 |
| 4.1. Главный экран DM .....                                      | 264 |
| 4.1.1. Интерфейс главного экрана DM .....                        | 264 |
| 4.1.2. Создание проекта DM .....                                 | 265 |
| 4.1.3. Удаление проекта DM .....                                 | 265 |
| 4.1.4. Копирование проекта DM .....                              | 265 |
| 4.1.5. Редактирование проекта DM .....                           | 265 |
| 4.1.6. Работа с проектом DM .....                                | 266 |
| 4.2. Конструктор цепочек решений .....                           | 266 |
| 4.2.1. Интерфейс Конструктора цепочек решений .....              | 266 |
| 4.2.2. Узлы решений.....   | 268 |
| 5. Компонент Управление моделями и решениями (Model Manager, MM) | 273 |
| 5.1. Интерфейс MM.....   | 273 |
| 5.2. Репозиторий.....  | 275 |
| 5.3. Проект.....   | 276 |
| 5.3.1. Создание проекта MM.....                                  | 277 |
| 5.3.2. Удаление проекта MM.....                                  | 278 |
| 5.3.3. Редактирование проекта MM.....                            | 278 |
| 5.4. Модель .....  | 279 |
| 5.4.1. Импорт модели .....                                       | 279 |
| 5.4.2. Удаление модели .....                                     | 282 |
| 5.4.3. Работа с моделью.....                                     | 283 |
| 5.5. Версия модели .....   | 286 |
| 5.5.1. Импорт версии модели .....                                | 286 |
| 5.5.2. Удаление версии модели .....                              | 286 |
| 5.5.3. Работа с версией модели .....                             | 287 |
| 5.6. Публикация версии модели .....                              | 291 |

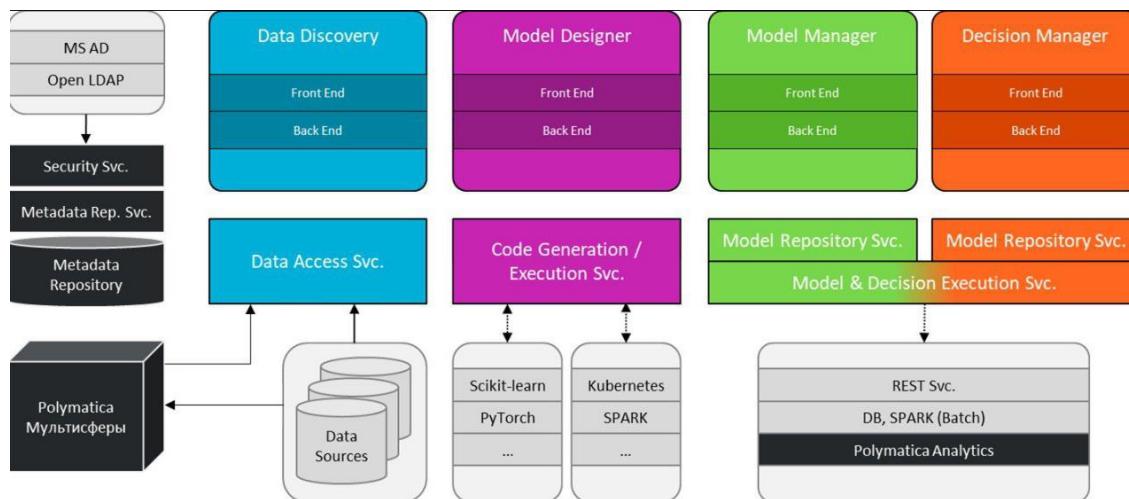
|  |     |
|--|-----|
| 5.6.1. Сервис .....  | 291 |
| 5.6.2. Пакетная публикация .....                                     | 294 |
| 5.6.3. Оценка характеристик.....                                     | 297 |
| 5.6.4. Скрипт .....  | 301 |
| 5.6.5. Публикация для скоринга мультисфер (Polymatica Analytics).... | 301 |
| 5.7. Согласование версий модели и публикаций.....                    | 302 |
| 5.7.1. Шаблон согласования.....                                      | 302 |
| 5.7.2. Процесс согласования .....                                    | 303 |
| 5.8. Библиотеки .....  | 307 |
| 5.9. Пользовательские атрибуты .....                                 | 308 |
| 6. Глобальный поиск.....   | 310 |
| 6.1. Интерфейс экрана Глобального поиска .....                       | 310 |
| 6.2. Работа с тэгами .....   | 311 |
| 6.2.1. Добавление тэга .....   | 311 |
| 6.2.2. Фильтрация при помощи тэгов .....                             | 311 |
| 6.3. Связанные объекты .....   | 312 |

# 1. Общая часть

## 1.1. Архитектура

Система представляет собой тонкий клиент (thin client) с веб-интерфейсом. В основу серверной части (Backend) Модуля заложена микросервисная архитектура.

На рисунке приведена целевая функциональная архитектура Модуля.



**Рисунок 1 Функциональная архитектура Модуля**

Основными компонентами архитектуры являются:

- Интерфейс пользователя (Frontend/Backend). Выделяются следующие подкомпоненты пользовательского интерфейса:
  - **Исследование данных (Data Discovery)** — пользовательский интерфейс для выбора исследуемых данных и их последующего интерактивного анализа. В текущей версии Модуля обеспечивается поддержка следующих операций:
    - выбор источника данных и формирования выборки;
    - визуализация выборки или ее части;
    - профилирование данных и оценка выборки по основным статистическим критериям;
    - отбор и преобразование признаков;
    - тестирование статистических гипотез о виде распределения непрерывных и дискретных наблюдений.
  - **Построение моделей машинного обучения (Model Designer)** — пользовательский интерфейс, позволяющий пользователю в графическом режиме в виде последовательности шагов (pipeline) построить модели машинного обучения, провести оценки точности моделей и сохранить модели в репозитории для дальнейшего использования.

В текущей версии Модуля обеспечивается поддержка следующих функций:

- отбор и преобразование признаков;
  - разделение выборки на обучающую, валидационную и проверочную;
  - построение моделей машинного обучения различными методами;
  - автоподбор гиперпараметров модели;
  - кросс-валидация;
  - тестирование моделей;
  - интерпретация моделей;
  - скоринг и оценка качества моделей — расчет, табличная и графическая визуализация различных характеристик модели;
  - выбор лучшей модели по заданному критерию;
  - возможность переобучения моделей при появлении новых данных;
  - сохранение и регистрация моделей в репозитории.
- **Управление жизненным циклом моделей (Model Manager)** — интерфейс пользователя для управления репозиторием моделей. Позволяет публиковать модели для применения в пакетном режиме или режиме сервиса, а также настраивать процессы согласования моделей. В текущей версии Модуля обеспечивается поддержка следующих операций:
    - ведение единого репозитория и версионности моделей;
    - публикация моделей для применения в пакетном режиме — формирование скрипта для применения модели на данных в БД;
    - публикация моделей в качестве сервиса;
    - публикация моделей из репозитория в среду Polumatica;
    - сравнение моделей и выбор лучшей;
    - снятие моделей с публикации.
  - **Управление решениями (Decision Manager)** — интерфейс пользователя для создания цепочек решений. В текущей версии Модуля обеспечивается поддержка следующих функций:
    - Построение цепочек решений в графическом интерфейсе с возможностью комбинирования моделей ИИ и экспертных правил;
    - Создание экспертных правил в графическом режиме;
    - Тестирование цепочек решений.
  - Сервис управления метаданными (Metadata Rep. Svc.) — компонент отвечает за хранение и управление внутренними метаданными модуля, такими как параметры подключения к источникам данных, проекты по построению моделей, параметры для публикации моделей и т.д. Остальные компоненты решения используют сервис управления метаданными для сохранения, чтения, изменения внутренних метаданных.

- Сервис управления безопасностью (Security Svc.) — компонент отвечает за управление правами доступа, интеграцию с каталогами пользователей (Open LDAP, MS ActiveDirectory), аутентификацию и авторизацию пользователей.
- Сервис по доступу к источникам данных (Data Access Svc.) — данный компонент отвечает за доступ к различным источникам и получение из них данных.
- Сервис по генерации и запуску вычислений (Code Generation and Execution Svc.) — данный компонент предназначен для перевода команд в выполняемый код на языке Python. В решении предполагается механизм регистрации различных библиотек Python (scikit-learn и др.), содержащих алгоритмы машинного обучения. При работе с данными в компоненте Data Discovery, а также при создании проекта и построении моделей в компоненте Model Designer — выполняется генерация кода Python в соответствии с заданными настройками.
- Сервис по управлению и публикации моделей (Model Repository Svc.) — компонент содержит процедуры, необходимые для управления жизненным циклом моделей — ведение реестра моделей, процессы согласования, процессы публикации.
- Сервис(ы) исполнения моделей (Model Execution Svc.) — это сервис(ы) исполнения моделей, т.е. каждая опубликованная модель — это исполняемый сервис.

## 1.2. Интерфейс

### 1.2.1. Авторизация

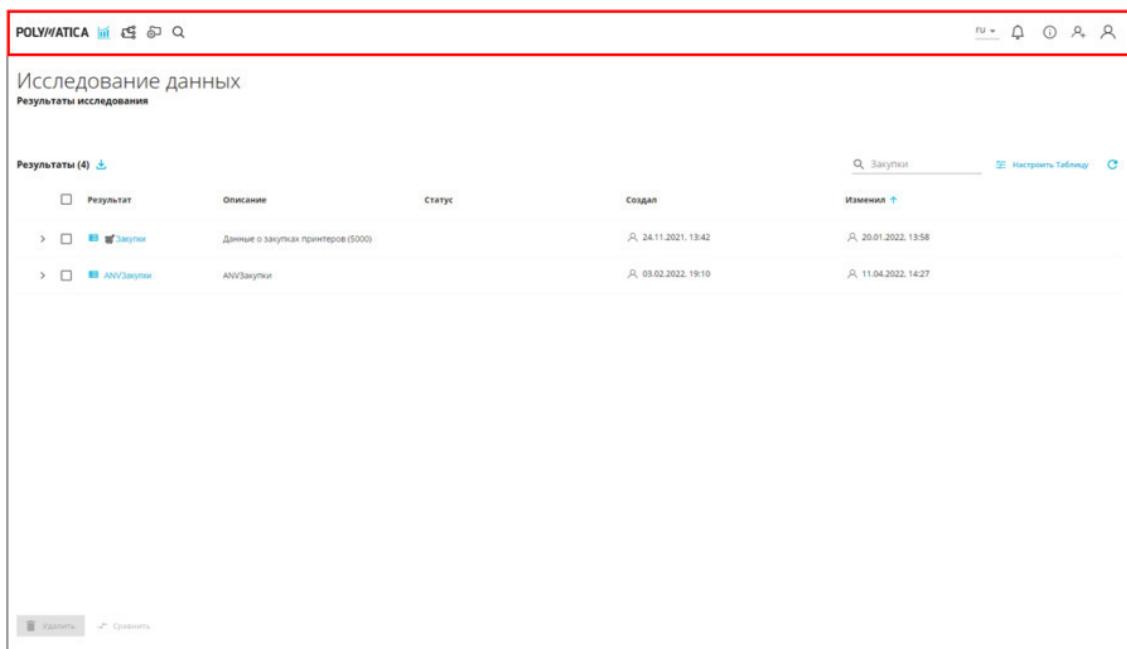
При запуске Модуля открывается экран авторизации в системе.

Логин и пароль учетной записи Пользователь должен получить у своего системного администратора.

При корректной авторизации в системе в зависимости от выданных пользователю прав откроется один из компонентов Модуля.

### 1.2.2. Основное меню

Основным элементом интерфейса является панель меню, расположенная в верхней части экрана (рисунок ниже). Она неизменна для всех компонентов Модуля.



## Рисунок 2 Интерфейс модуля. Меню

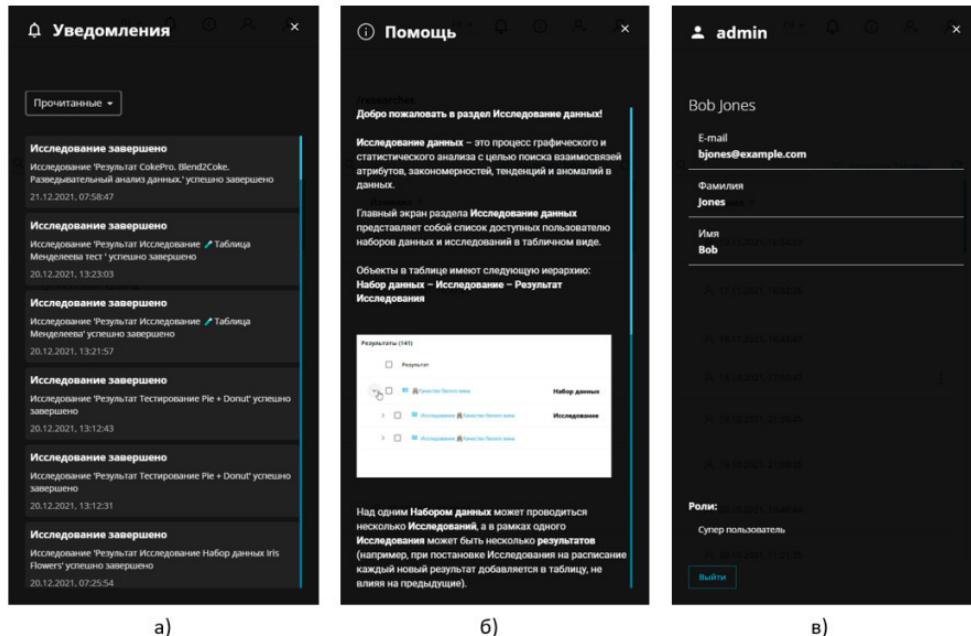
В левой части панели меню расположены объекты, которые позволяют:

- Выбор иконки открывает компонент Исследование данных (Data Discovery, DD).
- Выбор иконки открывает компонент Построение модели (Model Designer, MD).
- Выбор иконки открывает компонент Разработка решений (Decision Manager, DM).
- Выбор иконки открывает компонент Управление моделями и решениями (Model Manager, MM).
- Выбор иконки открывает интерфейс глобального поиска.

В правой части меню расположены объекты, которые позволяют:

- Для выбора языка отображения интерфейса предусмотрен список доступными языками.
- Выбор иконки открывает боковую панель с уведомлениями о завершении расчета результатов Исследований в DD (Рисунок а).
- Выбор иконки открывает боковую панель с интерактивной справкой (Рисунок б).

- Выбор иконки  открывает раздел Администрирования.
- Выбор иконки  открывает боковую панель с информацией об учетной записи пользователя, о лицензии и кнопкой выхода из Модуля (Рисунок в).



**Рисунок 3 Боковая панель Уведомления (а), боковая панель Помощь (б), боковая панель Пользователь (в)**

### 1.2.3. Роли пользователей

В зависимости от выданных пользователю прав в интерфейсе отображаются соответствующие компоненты:

- **Аналитик** имеет доступ к функционалу компонентов Исследование данных (DD) и Построение моделей (MD).
- **Менеджер моделей** имеет доступ ко всему функционалу компонента Управление моделями (MM).
- **Пользователь моделей** имеет доступ к разделу Опубликованные модели MM.
- **Согласующий** имеет доступ к разделу Согласование MM.
- **Суперпользователь** имеет полный доступ к функционалу Модуля.
- **Системный администратор** имеет доступ только к разделу Администрирование.

|                                | DD | MD | MM | Admin |
|--------------------------------|----|----|----|-------|
| <b>Аналитик</b>                | ✓  | ✓  |    |       |
| <b>Менеджер моделей</b>        |    |    | ✓  |       |
| <b>Пользователь моделей</b>    |    |    | ✓  |       |
| <b>Согласующий</b>             |    |    | ✓  |       |
| <b>Суперпользователь</b>       | ✓  | ✓  | ✓  | ✓     |
| <b>Системный администратор</b> |    |    |    | ✓     |

Таблица 1 Доступ к компонентам по ролям

## 2. Компонент Исследование данных (Data Discovery, DD)

Этап исследования данных представляет собой процесс графического и статистического анализа, необходимый для проверки общего качества данных и поиска взаимосвязей атрибутов, тенденций и аномалий. Полученная на данном этапе информация позволяет сформулировать гипотезы о том, как данные помогут решить поставленную задачу.

Для этого в компоненте Исследование данных предусмотрены следующие инструменты:

- Графическое представление данных (одномерные и многомерные графики и таблицы).
- Профилирование данных (статистические характеристики).
- Корреляционные матрицы.
- Кластерный анализ.
- Статистические тесты.
- Периодический запуск исследования данных (постановка исследований на расписание).

Компонент включает:

- Главный экран со списком доступных исследований.
- Мастер настройки исследования.
- Экран с результатами исследования.
- Окно постановки исследования на расписание.
- Вспомогательные окна настройки вида, пример данных и т.д.

### 2.1. Главный экран

#### 2.1.1. Интерфейс главного экрана DD

Главный экран раздела **Исследование данных** открывается при выборе иконки  в левом верхнем меню и представляет собой список доступных пользователю наборов данных и исследований в табличном виде (рисунок ниже).

The screenshot shows the 'Исследование данных' (Data Discovery) component interface. At the top, there's a search bar and navigation icons. Below it, a header says 'Исследование данных' and 'Результаты исследования'. A table titled 'Результаты (975)' lists items with columns: 'Результат' (Result), 'Описание' (Description), 'Статус' (Status), 'Создал' (Created by), and 'Изменил' (Changed by). The table includes entries like 'testestest', 'sqldtest33', and various 'Тестирование Iris Flowers' and 'Boston Housing Data' entries. At the bottom of the table are buttons for 'Удалить' (Delete) and 'Сохранить' (Save).

**Рисунок 4 Главный экран компонента Исследование данных (Data Discovery)**

Объекты в таблице имеют Иерархию в соответствии с рисунком ниже.

The diagram illustrates the object hierarchy. On the left, a tree view shows a 'Результат' (Result) node expanded to show a 'Набор данных' (Dataset) node labeled 'Качество белого вина' (Wine Quality). This dataset node has several 'Исследование' (Research) nodes, each associated with a specific 'Результат исследования' (Research Result) node. The results are all labeled 'Запуск по расписанию Исследование Качество белого вина' (Scheduled Research Wine Quality).

**Рисунок 5 Пример иерархии объектов в таблице Главного экрана компонента Исследование данных**

Для одного **Набора данных** (имеет символическое обозначение ) может проводиться несколько **Исследований** (обозначается ), например, с разным набором визуализаций и объектов. В рамках одного Исследования может быть несколько **Результатов** (обозначается ) — например, при постановке Исследования на расписание.

Таблица с доступными исследованиями имеет гибкие настройки отображения. Так пользователь может:

- изменить ширину любого столбца (для этого необходимо перетащить границу его заголовка  до нужной ширины);
- сортировать таблицу (для этого необходимо выбрать иконку  рядом с заголовком сортируемого столбца);
- скрывать/отображать столбцы и изменять их порядок в окне **Вид таблицы** ( **Настроить Таблицу**) в правом верхнем углу таблицы; при выборе иконки  столбец скроется, при наведении на иконку  активируется возможность перемещения столбца);
- сбросить внесенные изменения также в окне **Вид таблицы** (для этого выбрать кнопку «**Сбросить**»).

Для быстрого поиска объекта в таблице предусмотрено поле  **Поиск...** в правой верхней части таблицы.

Объекты таблицы можно выгрузить в формате Excel. Для этого нужно выбрать иконку **Экспорта в excel** .

Другие подсказки по интерфейсу Главного экрана:

- Описание набора данных подтягивается при его регистрации (подробнее в Руководстве Администратора), для Исследования задается в окне **Мастера настройки**.
- Исследование может иметь **Статус**, связанный с постановкой на расписание (иконка  — Запланировано), результат исследования может иметь **Статус**, связанный с ходом выполнения исследования (иконка  — Обрабатывается, иконка  — Готово).
- Иконка  рядом с датой создания и изменения позволяет узнать, кто является автором исследования/изменения исследования.
- Иконка  позволяет открыть выпадающее меню, в котором Пользователь может **Запустить исследование**.

## 2.1.2. Создание исследования

Для создания нового Исследования необходимо выполнить следующие шаги:

- Выбрать необходимый набор данных из списка доступных.
- В открывшемся **Мастере настройки** выбрать и задать настройки:
  - визуализаций (подробнее в разделе Визуализации),
  - корреляции (подробнее в разделе Корреляция),
  - кластеризации (подробнее в разделе Кластеризация),
  - статистических тестов (подробнее в разделе Статистические тесты).

- Запустить расчет Исследования, выбрав пункт **Запустить** в выпадающем меню : рядом с необходимым Исследованием.
- После успешного выполнения Исследования (статус «Готово» ) открыть результаты Исследования, выбрав из списка необходимый.

### **2.1.3. Удаление исследования**

Для удаления Исследования (результата Исследования) необходимо выполнить следующие шаги:

- Выбрать чекбокс рядом с удаляемым объектом (при выборе чекбокса рядом с Набором данных выбираются все объекты, нижестоящие в иерархии — это все Исследования данного набора и их результаты).
- Выбрать кнопку **«Удалить»**.

## **2.2. Мастер настройки**

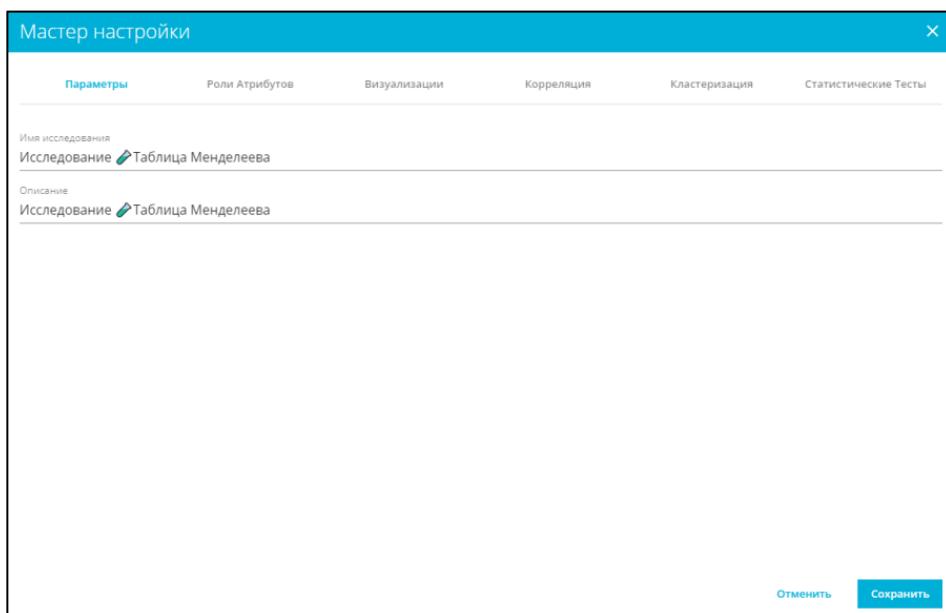
### **2.2.1. Интерфейс Мастера настройки**

При создании нового исследования открывается окно **Мастер настройки**, который включает в себя следующие вкладки:

- Параметры.
- Роли атрибутов.
- Визуализации.
- Корреляция.
- Кластеризация.
- Статистические тесты.

### **2.2.2. Параметры**

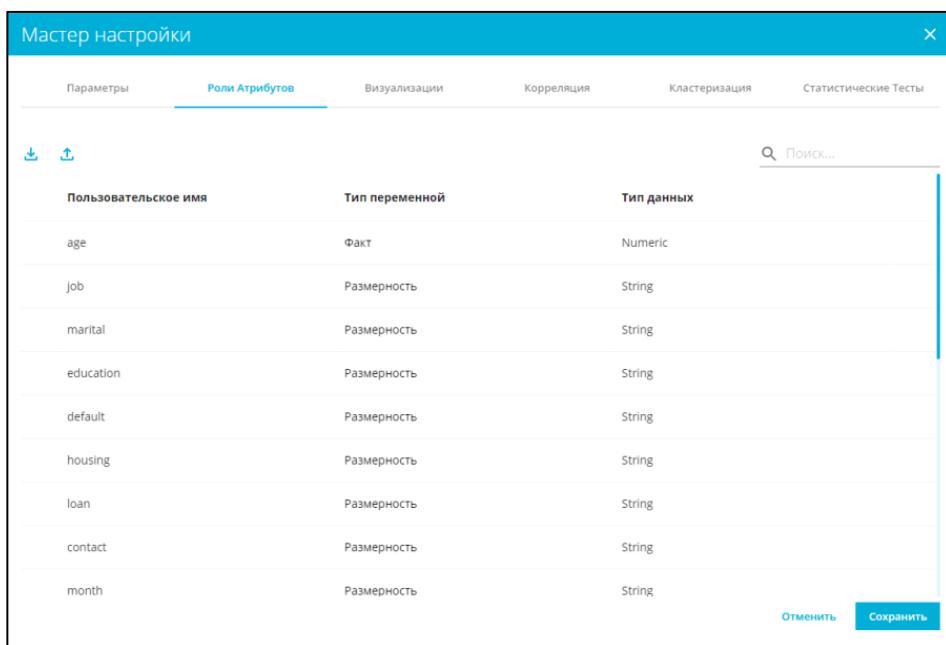
Во вкладке **Параметры** пользователь может задать Название исследования и Описание (рисунок ниже), которые отображаются в таблице на Главном экране компонента.



**Рисунок 6 Вкладка «Параметры» Мастера настройки**

### **2.2.3. Роли атрибутов**

Во вкладке **Роли атрибутов** Пользователь может ознакомиться с параметрами атрибутов набора данных и изменить тип переменной.



**Рисунок 7 Вкладка «Роли атрибутов» Мастера настройки**

Пользовательское имя задается при добавлении Набора данных в разделе Администрирования (подробнее [Руководство администратора](#)).

В компоненте предусмотрены следующие типы переменных и типы данных (задаются по умолчанию):

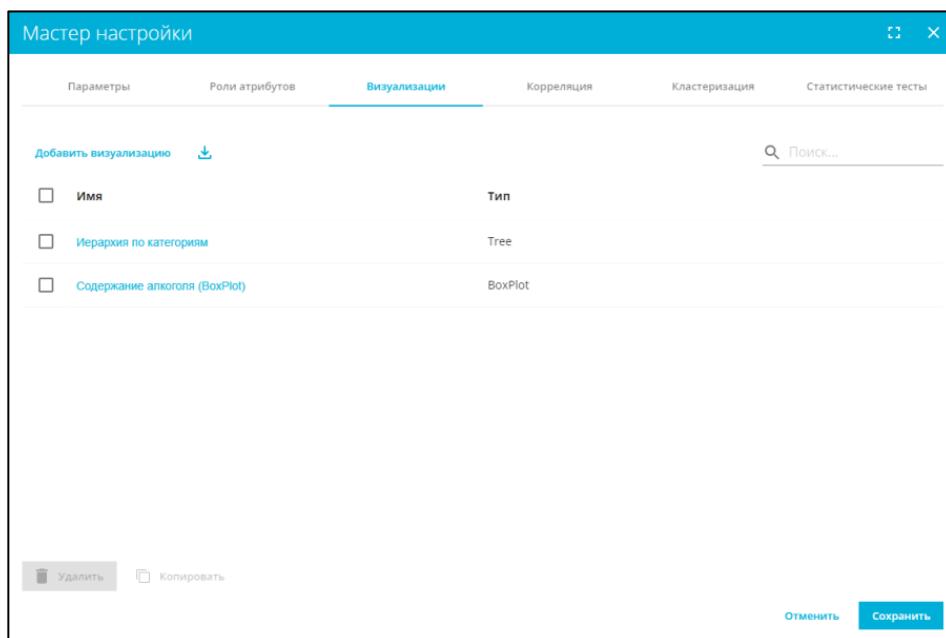
- **Факт:**
  - Numeric (включает и целые числа, и числа с плавающей точкой).
  - Date (дата).
  - DateTime (дата и время).
- **Размерность:**
  - String (строка).

Т.к. визуализации могут строится лишь с определенными типами переменных (подробнее [Визуализации](#)), предусмотрено изменение Типа переменной. Для этого необходимо щелкнуть левой кнопкой мыши по выбранному атрибуту и в открывшемся меню выбрать необходимый тип. Тип данных изменить нельзя.

Изменить Тип переменных также можно выгрузив Excel из системы, внести изменения, и загрузить обратно. Для этого нужно выбрать иконки **Экспорта в Excel**  и **Импорта из Excel** .

## 2.2.4. Визуализации

Данная вкладка представляет собой список визуализаций, которые будут рассчитаны в ходе Исследования.



**Рисунок 8 Вкладка «Визуализации» Мастера настройки**

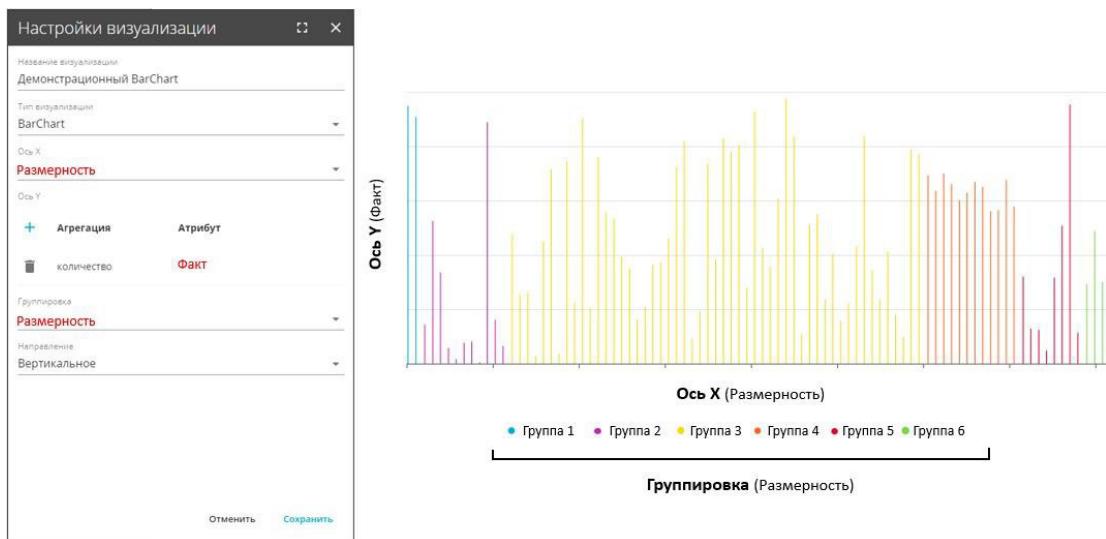
Для создания новой Визуализации необходимо выбрать **«Добавить визуализацию»**. В результате откроется окно **Настройки визуализации**.

В компоненте предусмотрены следующие типы визуализаций:

- Столбчатая диаграмма (Bar Chart).
- Диаграмма с накоплением (Stacked Bar Chart).
- Круговая диаграмма (Pie Chart).
- Кольцевая диаграмма (Donut Chart).
- Линейный график (Line Chart).
- Точечный график (Scatter Plot).
- Точечный 3D график (Scatter Plot 3D).
- Диаграмма размаха (BoxPlot).
- Древовидная диаграмма (Tree).
- Тепловая карта (HeatMap).
- Матрица сопряженности.

#### 2.2.4.1. Столбчатая диаграмма (Bar Chart)

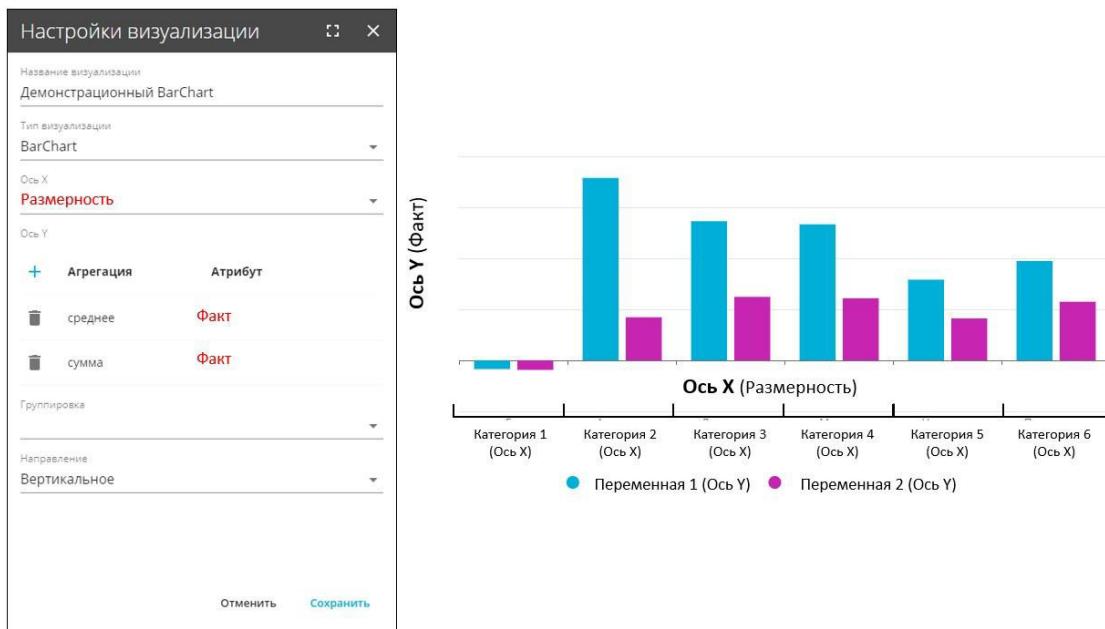
Столбчатая диаграмма отображает числовое сравнение между разными категориями. На одной оси (параметр **Ось X**) представлены конкретные сравниваемые категории, а на другой (параметр **Ось Y**) — шкала числовых значений. Диаграмма также позволяет отражать группировку по другой категории при помощи цвета (параметр **Группировка**).



**Рисунок 9 Столбчатая диаграмма с вертикальным направлением столбцов и группировкой**

Столбцы могут быть расположены вертикально или горизонтально (параметр **Направление**).

Можно выбрать несколько атрибутов (параметр **Ось Y**), но тогда будет недоступна Группировка.

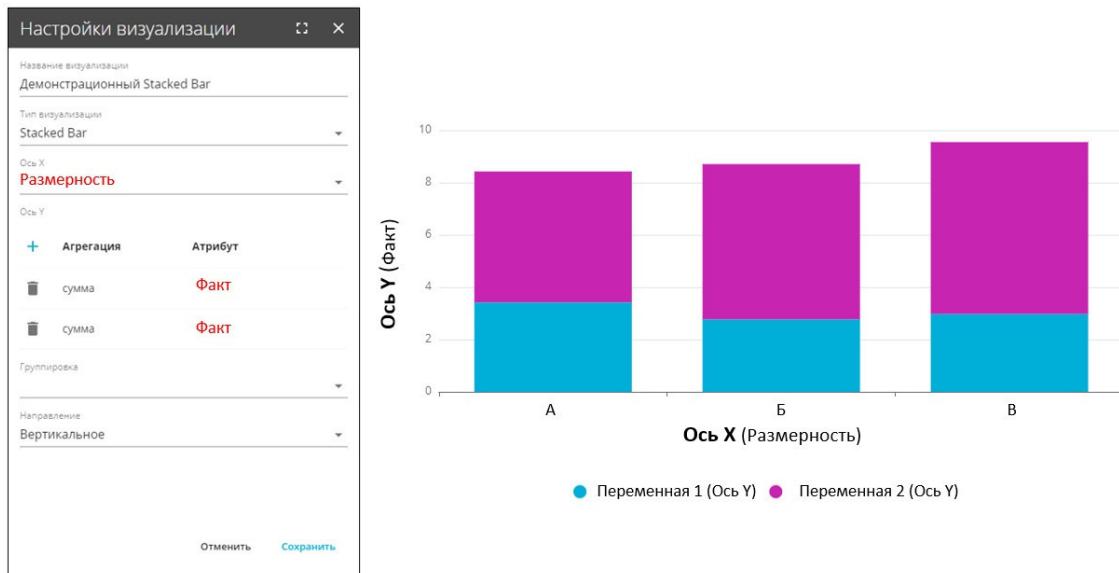


**Рисунок 10 Столбчатая диаграмма с несколькими атрибутами**

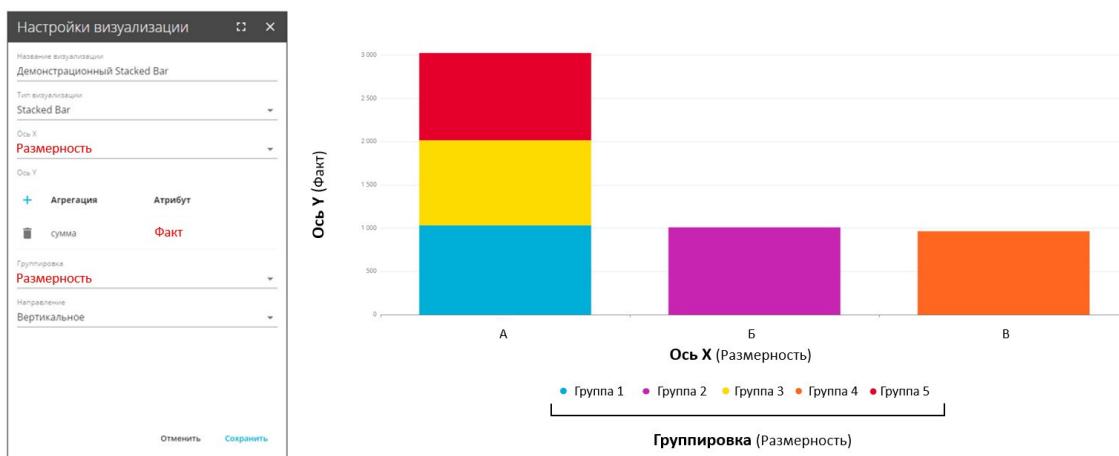
#### 2.2.4.2. Диаграмма с накоплением (Stacked Bar Chart)

Диаграмма с накоплением показывает вклад нескольких элементов данных (параметр **Ось Y**) в суммирующий результат в виде столбцов, расположенных друг над другом. Высота каждого столбца пропорциональна значению соответствующего элемента данных. Сравнение происходит между разными категориями (параметр **Ось X**). Можно выбрать несколько атрибутов (параметр **Ось Y**), но тогда будет недоступна Группировка (параметр **Группировка**).

Столбцы могут быть расположены вертикально или горизонтально (параметр **Направление**).



**Рисунок 11 Диаграмма с накоплением и несколькими атрибутами Ось Y**



**Рисунок 12 Диаграмма с накоплением и параметром Группировка**

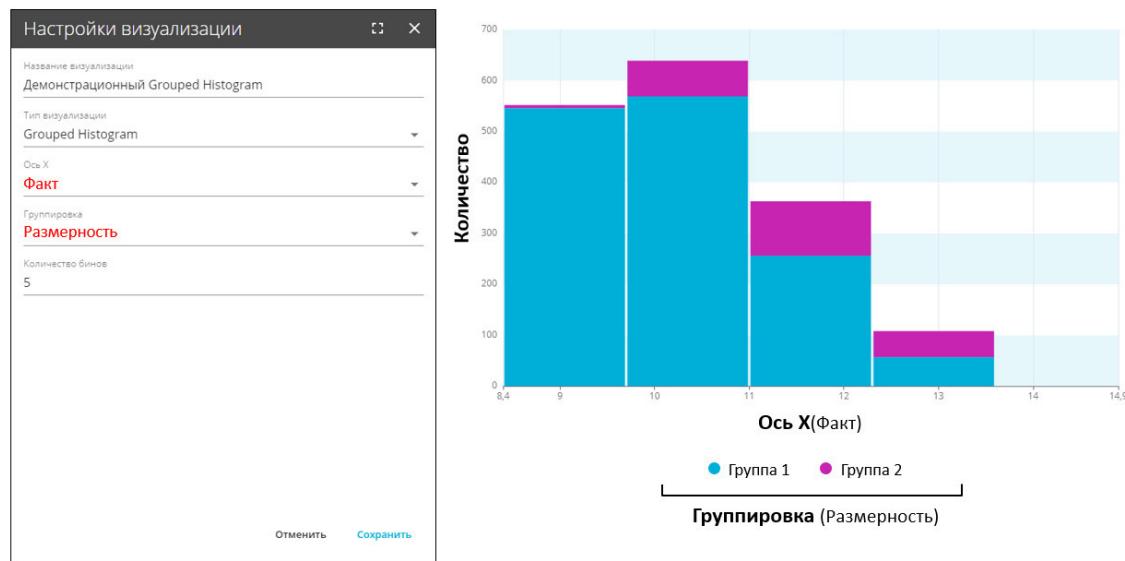
#### 2.2.4.3. Гистограмма с группировкой (Grouped Histogram)

Гистограмма агрегирует числовые данные (**параметр Ось X**) по группам с равными интервалами, которые называют **бинами**, и отображает частоту встречаемости значений в каждом из бинов (**параметр Количество бинов**).



**Рисунок 13 Гистограмма**

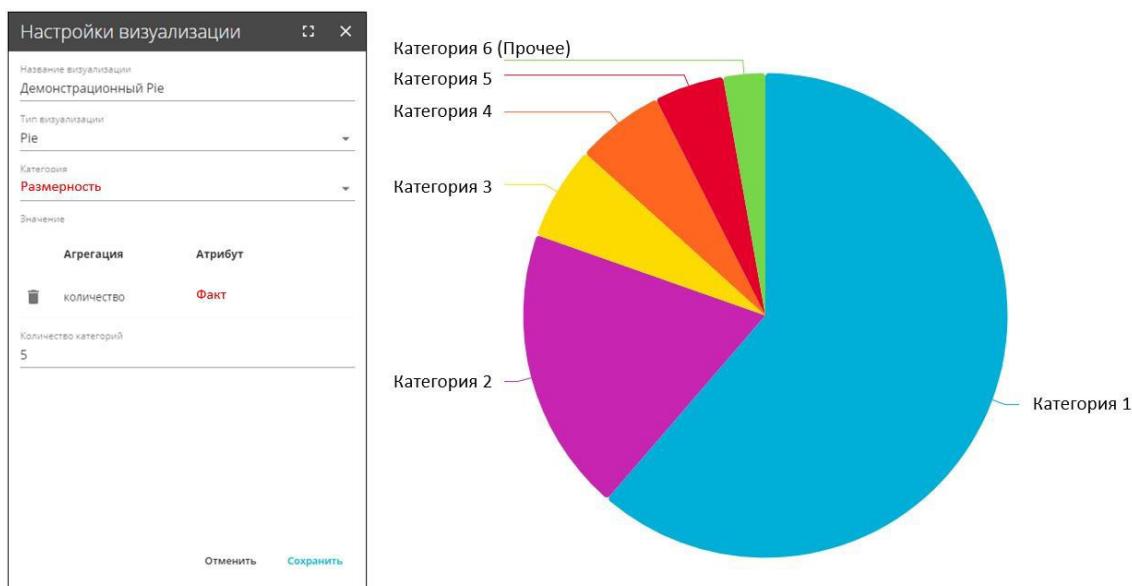
Для построения гистограммы с группировкой необходимо задать атрибут для Группировки (**параметр Группировка**).



**Рисунок 14 Сгруппированная гистограмма**

#### 2.2.4.4. Круговая диаграмма (Pie Chart)

Круговая диаграмма показывает процентное соотношение между категориями (параметр **Категория**). Длина каждой дуги представляет собой пропорциональную длину категории (параметр **Значение**) от длины всей окружности, равной 100%.

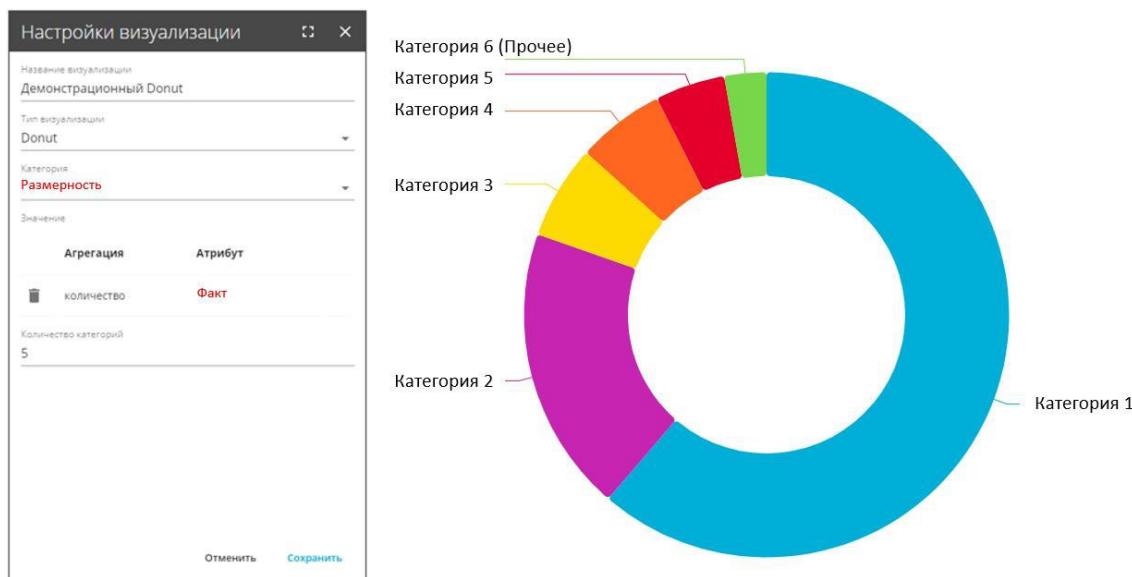


**Рисунок 15 Круговая диаграмма**

Если данные содержат множество категорий, пользователь может сократить количество отображаемых на диаграмме категорий для удобства интерпретации (параметр **Количество категорий**). Не вошедшие категории (согласно сортировки по убыванию) объединяются в единую категорию.

#### 2.2.4.5. Кольцевая диаграмма (Donut Chart)

Кольцевая диаграмма представляет собой круговую диаграмму, но с вырезанной центральной частью. Это позволяет пользователю сосредоточить внимание на длине дуг, а не на сопоставлении пропорций сегментов.

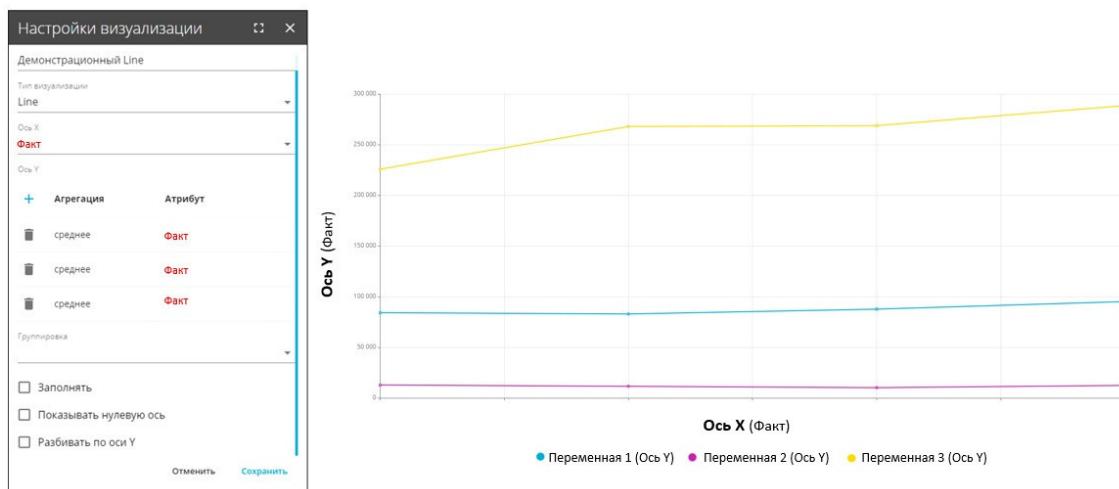


**Рисунок 16 Кольцевая диаграмма**

Аналогично круговой диаграмме пользователь может сократить количество отображаемых на диаграмме категорий (параметр **Количество категорий**). Не вошедшие категории (согласно сортировки по убыванию) объединяются в единую категорию.

#### 2.2.4.6. Линейный график (Line Chart)

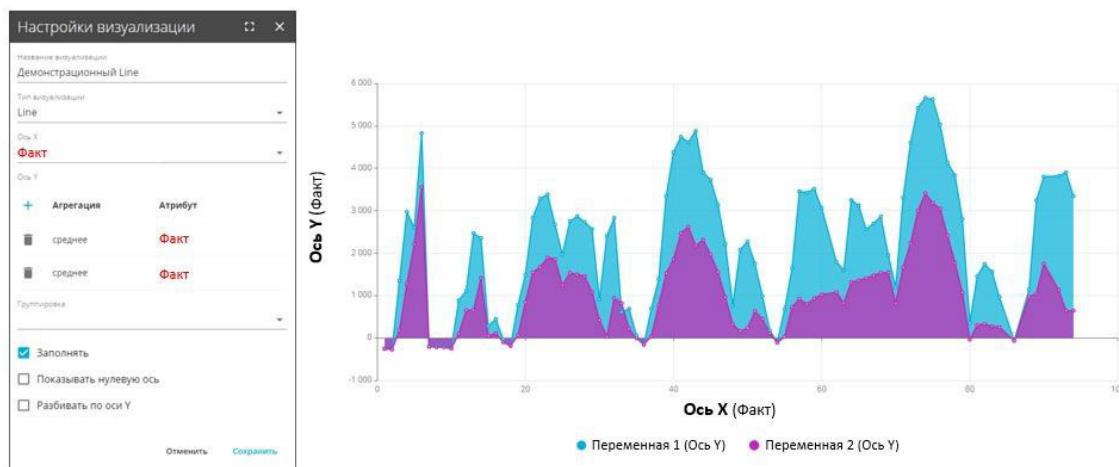
Линейный график отображает изменение показателей (параметр **Ось Y**) за некоторый интервал (параметр **Ось X**). Как правило, на оси Y отмечаются количественные значения, а на оси X шкала последовательностей.



**Рисунок 17 Линейный график**

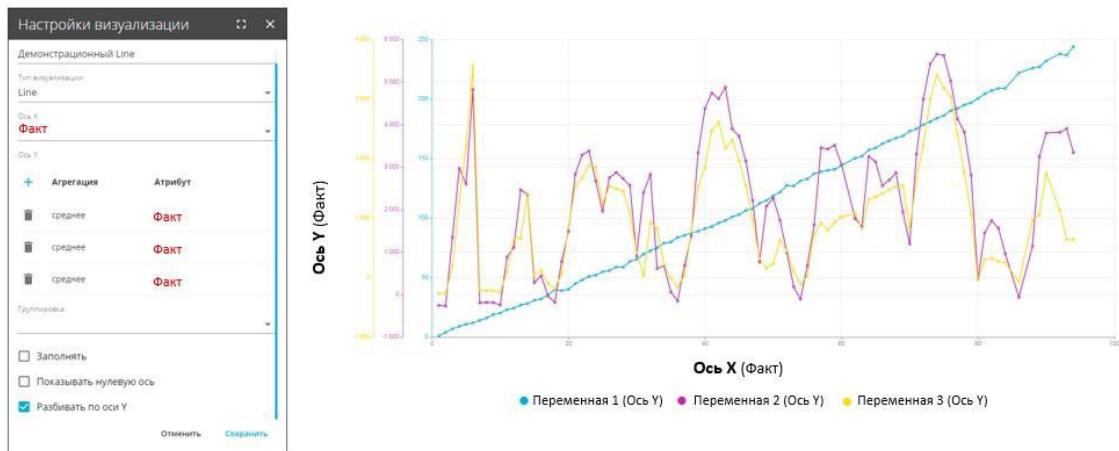
Также предусмотрена возможность группировки наблюдений по цвету (параметр **Группировка**).

Выбрав чекбокс **Заполнить**, пользователь может построить накопительную диаграмму, в которой область ниже линии заполнена определённым цветом.



**Рисунок 18 Линейный график с накоплением**

Многоосевая диаграмма позволяет сравнить несколько показателей с различными диапазонами значений и использует для этого две или более оси Y и одну общую ось X (чекбокс **Разбивать по оси Y**).



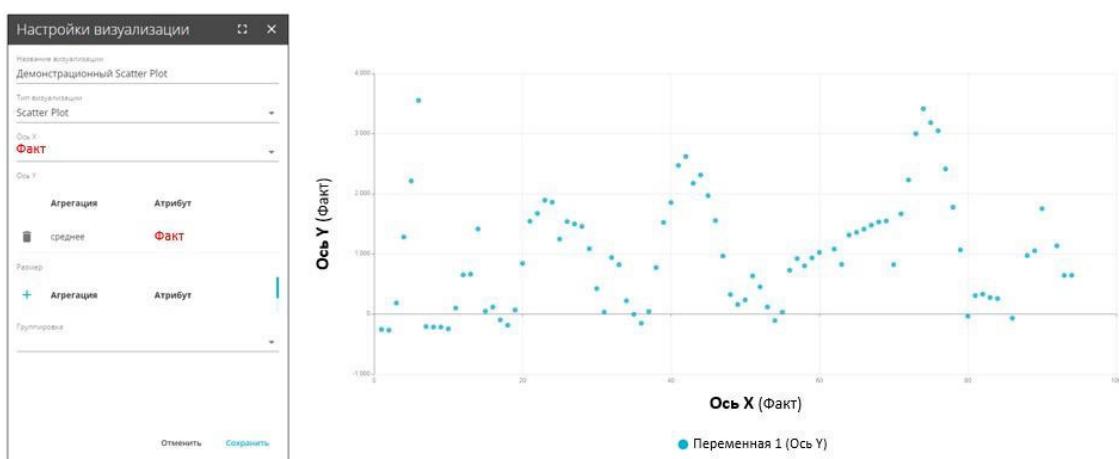
**Рисунок 19 Многоосевой линейный график**

#### 2.2.4.7. Точечный график (Scatter Plot)

Точечный график активно используется для:

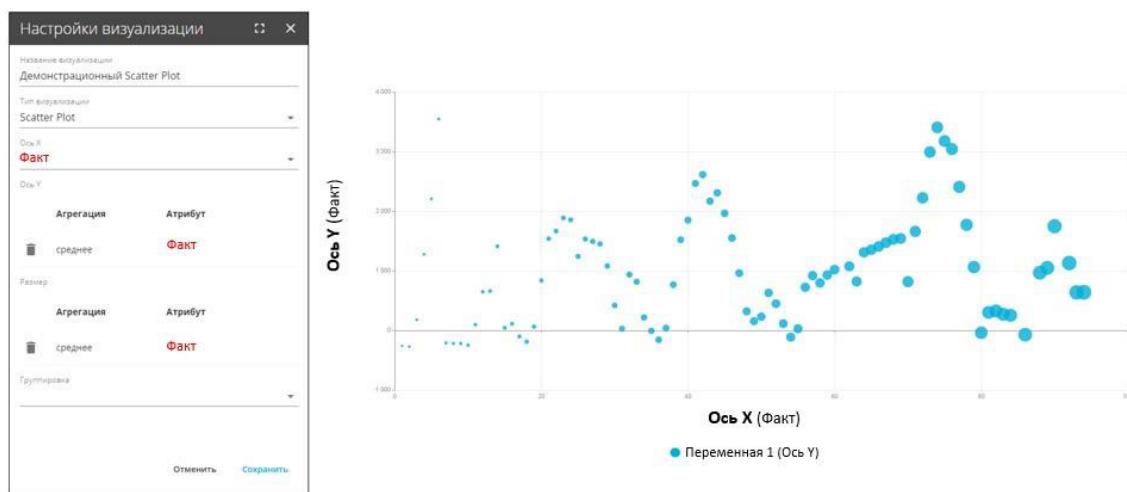
- Оценки корреляции переменных;
- Обнаружения выбросов и ошибок в данных;
- Оценки скученности наблюдений.

Для построения графика пользователю необходимо задать две количественные переменные (параметры **Ось X** и **Ось Y**).



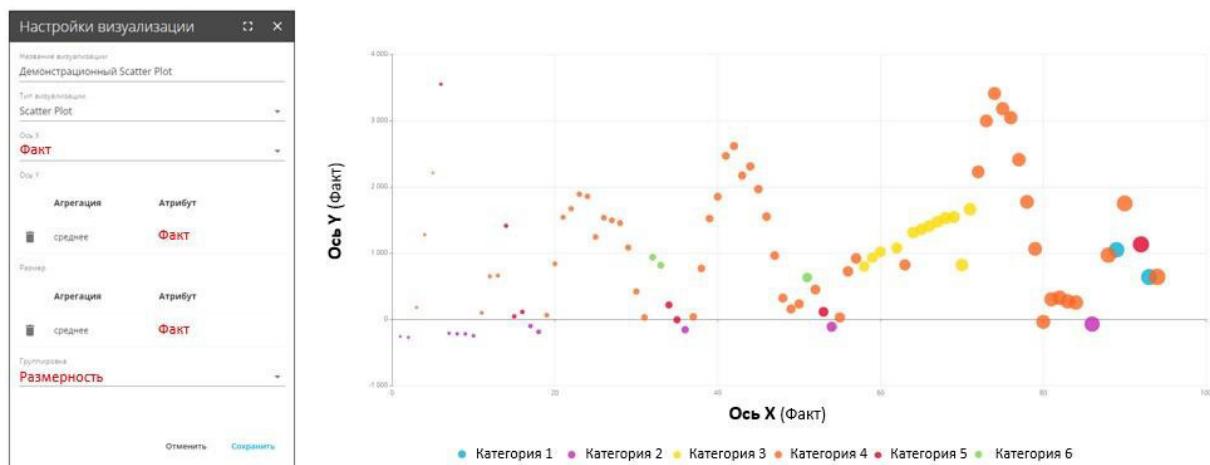
**Рисунок 20 Точечный график**

При необходимости можно масштабировать наблюдения согласно значению третьей количественной переменной (параметр **Размер**).



**Рисунок 21 Точечный график с заданным параметром Размер**

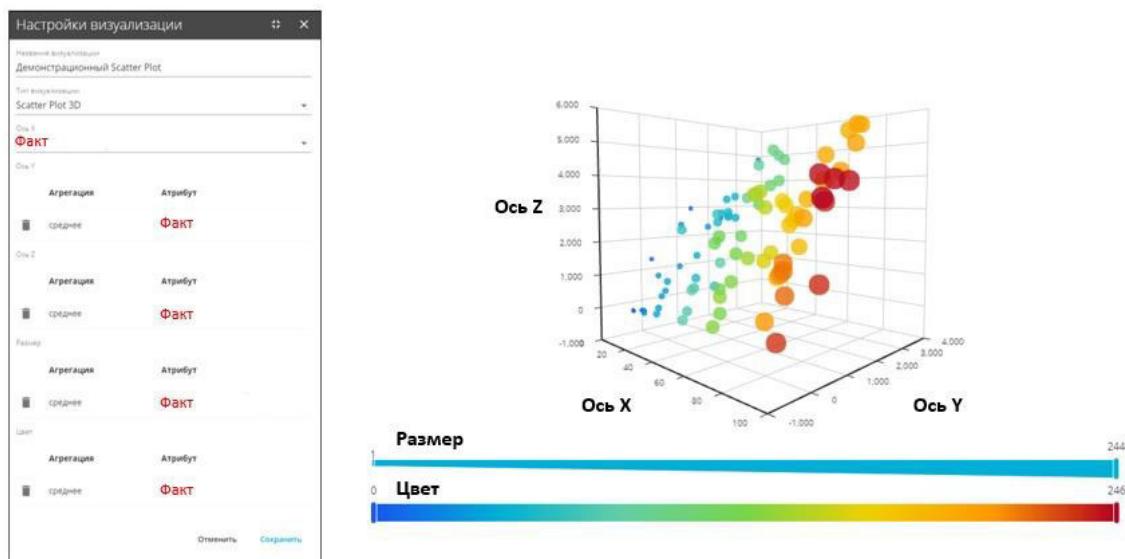
Пользователь также может присвоить цвет наблюдениям согласно значению некоторой категориальной переменной (параметр **Группировка**).



**Рисунок 22 Точечный график с заданными параметрами Размер и Группировка**

#### 2.2.4.8. Точечный 3D график (Scatter Plot 3D)

Точечный 3D график отображает отношение трех количественных переменных (параметры **Ось X**, **Ось Y** и **Ось Z**) в трехмерном пространстве, позволяет масштабировать наблюдения согласно четвертой количественной переменной (параметр **Размер**) и цветовым градиентом выделить различия по пятой количественной переменной (параметр **Цвет**).



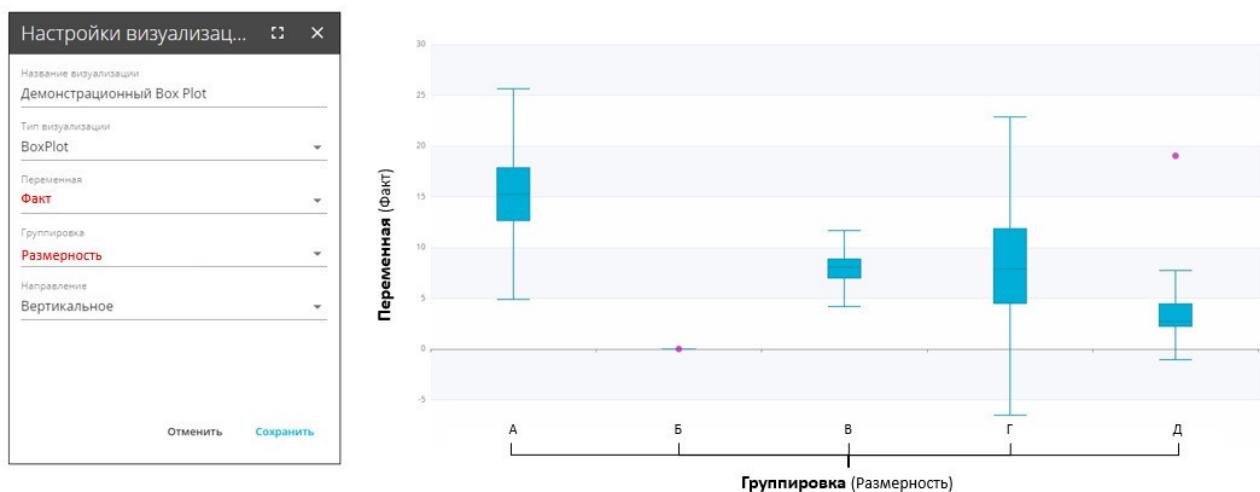
**Рисунок 23 Точечный 3D график с заданными параметрами Размер и Цвет**

#### 2.2.4.9. Диаграмма размаха (Boxplot)

Диаграмма размаха («ящик с усами») — удобный способ визуального представления непрерывных переменных (параметр **Переменная**), который позволяет:

- Отобразить значение ключевых статистик: медиана, верхний/нижний квартили, максимальное и минимальное значение выборки.
- Выявить выбросы.
- Определить степень разброса (дисперсии) и асимметрии данных.

Несколько таких диаграмм можно нарисовать бок о бок, чтобы визуально сравнивать распределение переменной в разных группах (параметр **Группировка**). Диаграммы можно располагать как горизонтально, так и вертикально (параметр **Направление**).



**Рисунок 24 Диаграмма размаха**

При наведении на ящик появляется панель со следующими статистиками:

- Min — минимальное значение;
- Q1 — нижний квартиль;
- Median — медианное значение;
- Q3 — верхний квартиль;
- Max — максимальное значение.

Помимо самих ящиков на диаграмме размаха отображаются выпадающие значения в виде фиолетовых точек (за пределами полутора межквартильных интервалов).

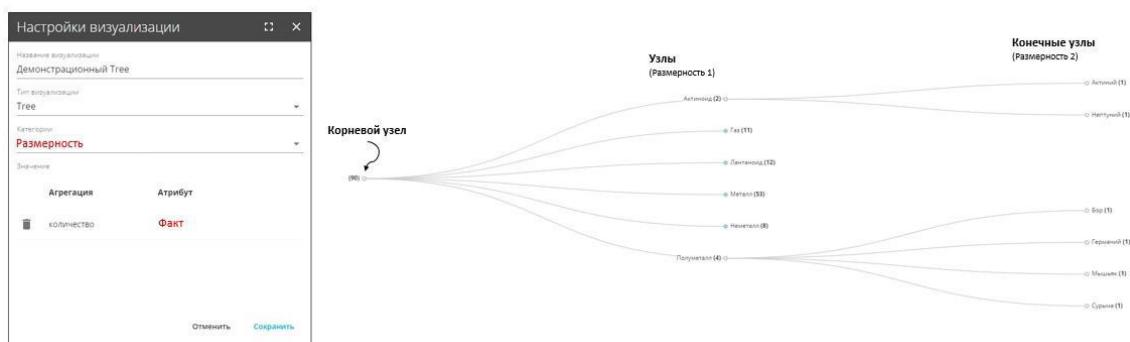
#### 2.2.4.10. Древовидная диаграмма (Tree)

Древовидная диаграмма — это метод визуального представления иерархии в древовидной структуре.

Структура древовидной диаграммы состоит из следующих элементов:

- Корневой узел — элемент, не имеющий вышестоящего элемента,
- Узлы и соединяющие их линии — ветви, которые обозначают взаимосвязи и отношения между элементами,
- Листовые узлы — элементы, у которых нет дочерних элементов.

Для создания данной диаграммы необходимо последовательно указать категориальные переменные (параметр **Категории**) и количественную переменную, которая будет отражаться в каждом узле (параметр **Значение**).

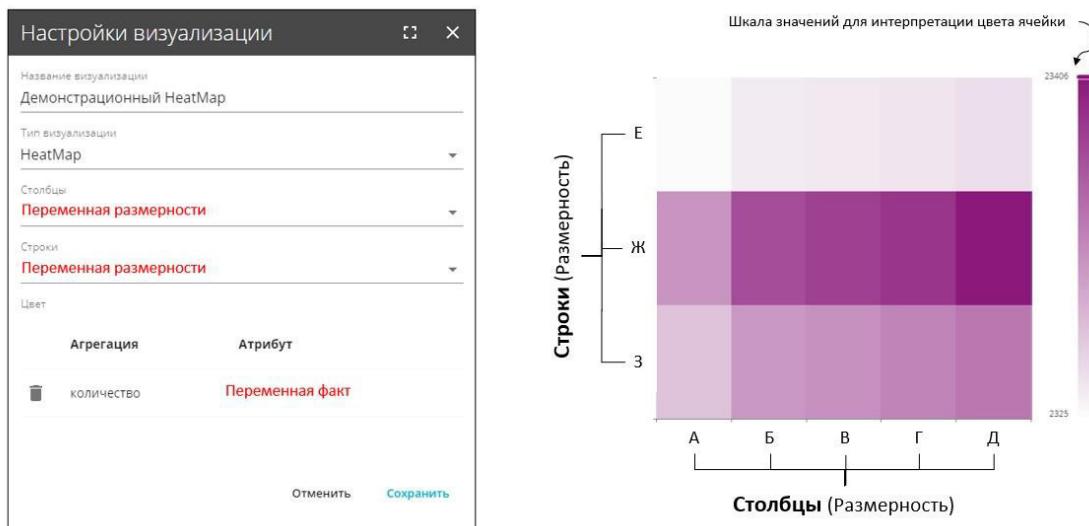


**Рисунок 25 Древовидная диаграмма**

#### 2.2.4.11. Тепловая карта (Heat Map)

Тепловые карты позволяют анализировать многомерные данные за счет распределения переменных по рядам и столбцам и закрашивания цветом ячеек таблицы.

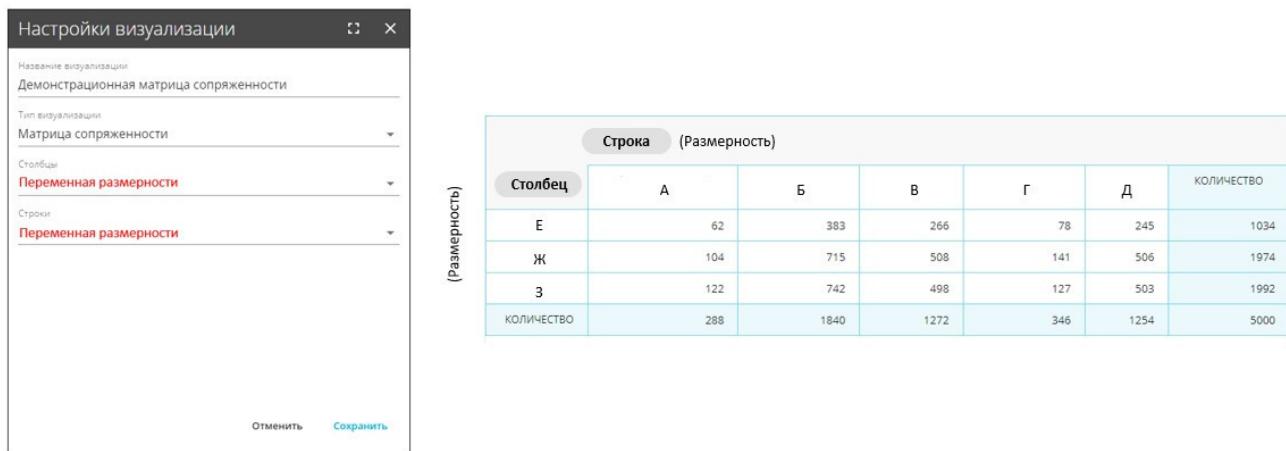
Ряды отображают одну категорию (параметр **Строки**), а все столбцы — другую (параметры **Столбцы**). Отдельные ряды и столбцы делятся на подкатегории, пересекающиеся друг с другом в рамках матрицы. В ячейках таблицы содержатся отображаемые в цвете количественные данные (параметр **Цвет**). Данные каждой ячейки основаны на взаимосвязи двух переменных: в строке и в столбце.



**Рисунок 26 Тепловая карта**

#### 2.2.4.12. Матрица сопряженности (Pivot)

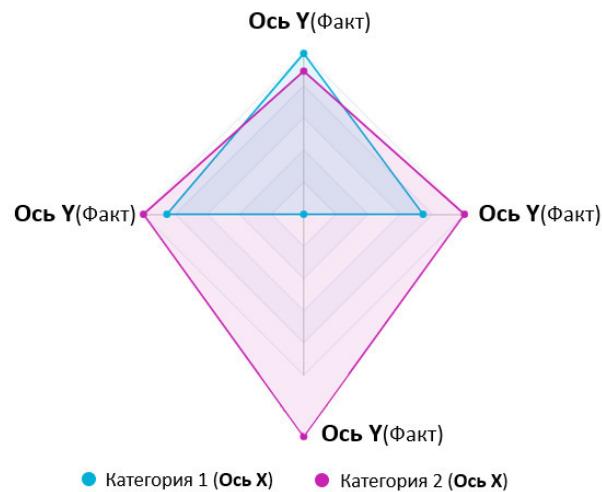
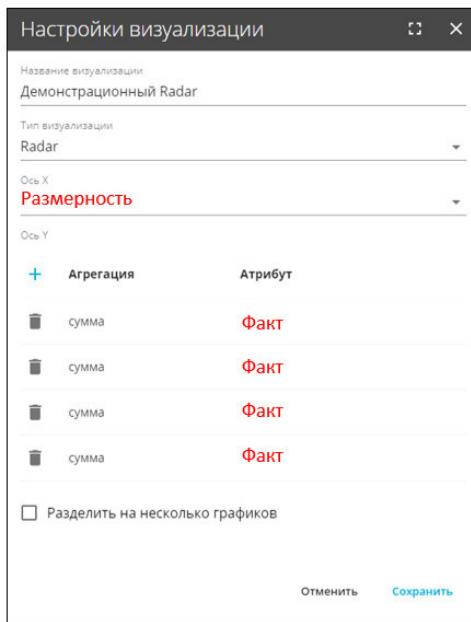
Таблица сопряженности служит для описания связи двух номинальных переменных, на пересечение столбов (параметр **Столбцы**) и строк (параметр **Строки**) которой указывается частота совместного появления соответствующих значений двух признаков.



**Рисунок 27 Матрица сопряженности**

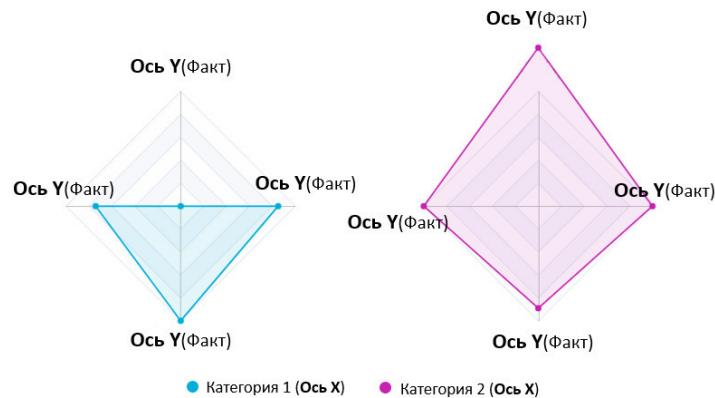
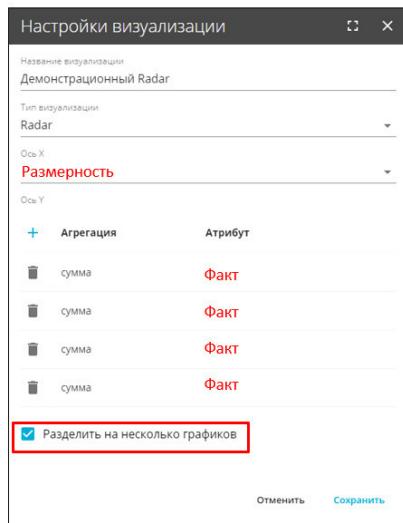
#### 2.2.4.13. Диаграмма-радар (Radar)

Диаграмма-радар представляет собой линейные графики, в которых значения **Оси X** обернуты на 360 градусов и для каждого значения x имеются значения **Оси Y**. Для корректного отображения графика необходимо задать не менее 3 агрегаций по **Оси Y**.



**Рисунок 28 Диаграмма-радар**

При выборе чекбокса **Разделить на несколько графиков** будут построены графики, количеством совпадающим с числом уникальных значений переменной **Оси X**.



**Рисунок 29 Диаграмма-радар с разделением на 2 графика**

## 2.2.5. Корреляция

Данная вкладка представляет собой список корреляционных матриц, которые будут рассчитаны в ходе Исследования.

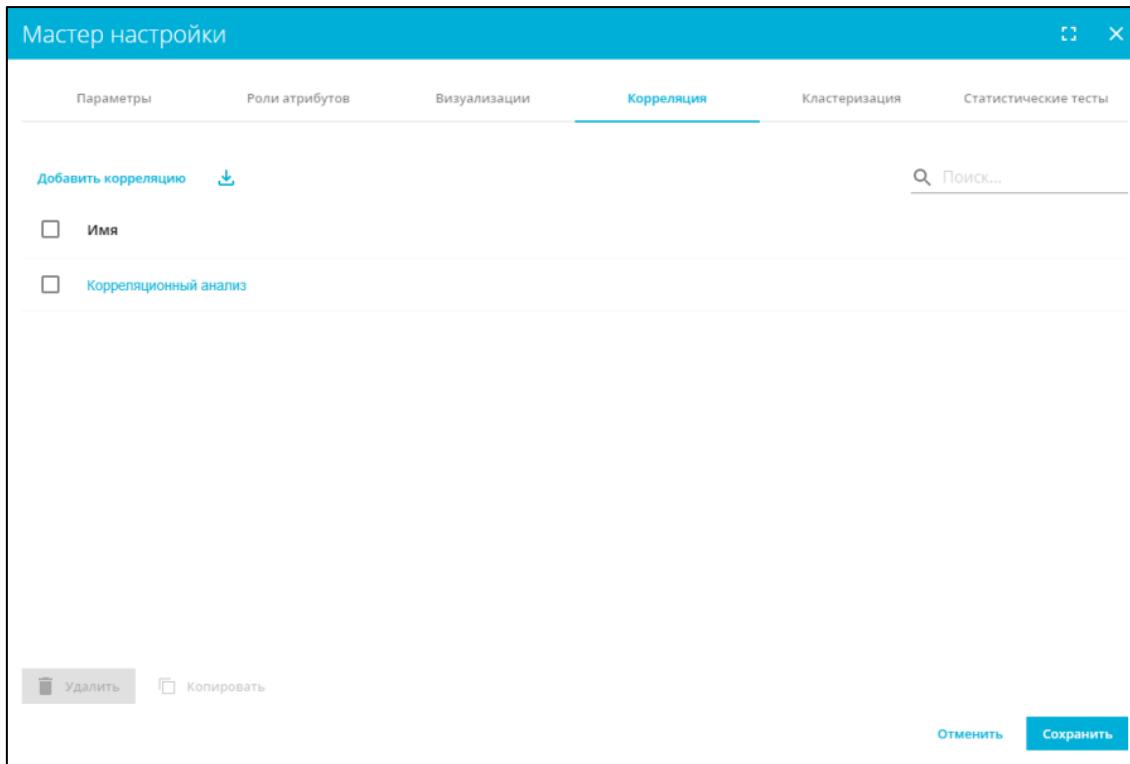
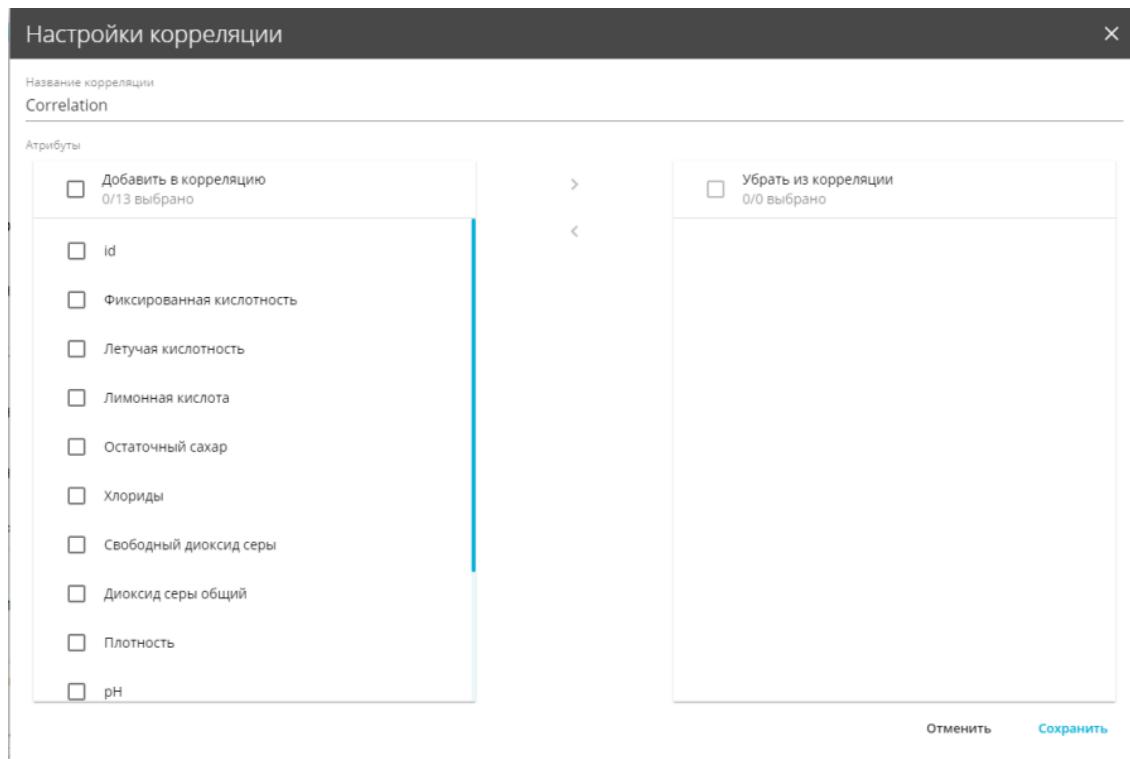


Рисунок 30 Вкладка «Корреляция» Мастера настройки

**Корреляционная матрица** — это квадратная таблица, в которой на пересечении соответствующих строк и столбца находится коэффициент корреляции для соответствующей пары признаков. **Коэффициент корреляции** — статистическая мера, которая отражает силу связи между признаками.

Для создания новой корреляционной матрицы необходимо:

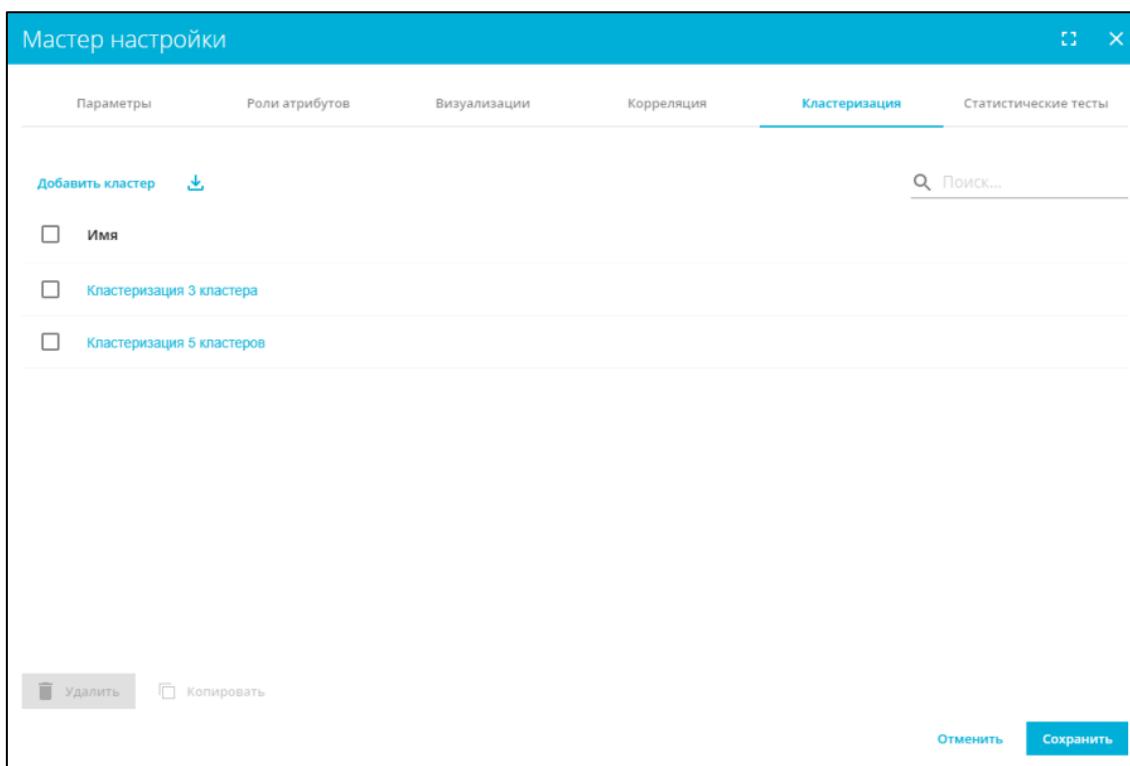
- Выбрать «**Добавить корреляцию**».
- В открывшемся окне **Настройки корреляции** задать Название корреляции.
- Добавить необходимые для анализа признаки, выбрав соответствующие чекбоксы рядом, и нажать на стрелку переноса в правое поле (для добавления всех имеющихся атрибутов выбрать чекбокс рядом с **Добавить в корреляцию** в левом поле).
- Сохранить параметры.



**Рисунок 31 Окно Настройки корреляции**

## 2.2.6. Кластеризация

Данная вкладка представляет собой список с процедурами кластерного анализа, которые будут рассчитаны в ходе Исследования.



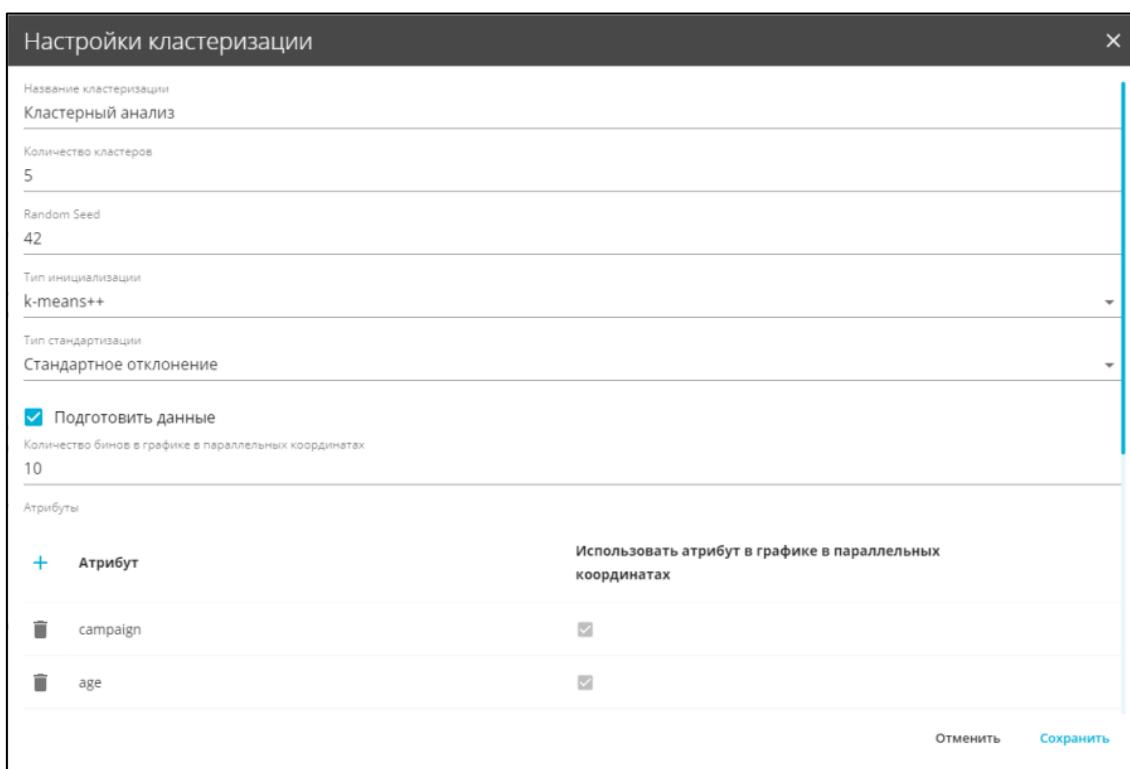
**Рисунок 32 Вкладка «Кластеризация» Мастера настройки**

**Кластеризация** — это разбиение множества объектов на подмножества (кластеры) по заданному критерию.

На этапе исследования полезно получить априорную информацию об исходных данных, чтобы на этапе построения модели k-средних указать параметр, который будет увеличивать ее предсказательную способность.

Для создания новой процедуры кластерного анализа необходимо:

- Выбрать «**Добавить кластер**».
- В открывшемся окне **Настройки кластеризации** задать параметры (подробнее в таблице ниже).
- Сохранить параметры.



**Рисунок 33 Окно Настройки кластеризации**

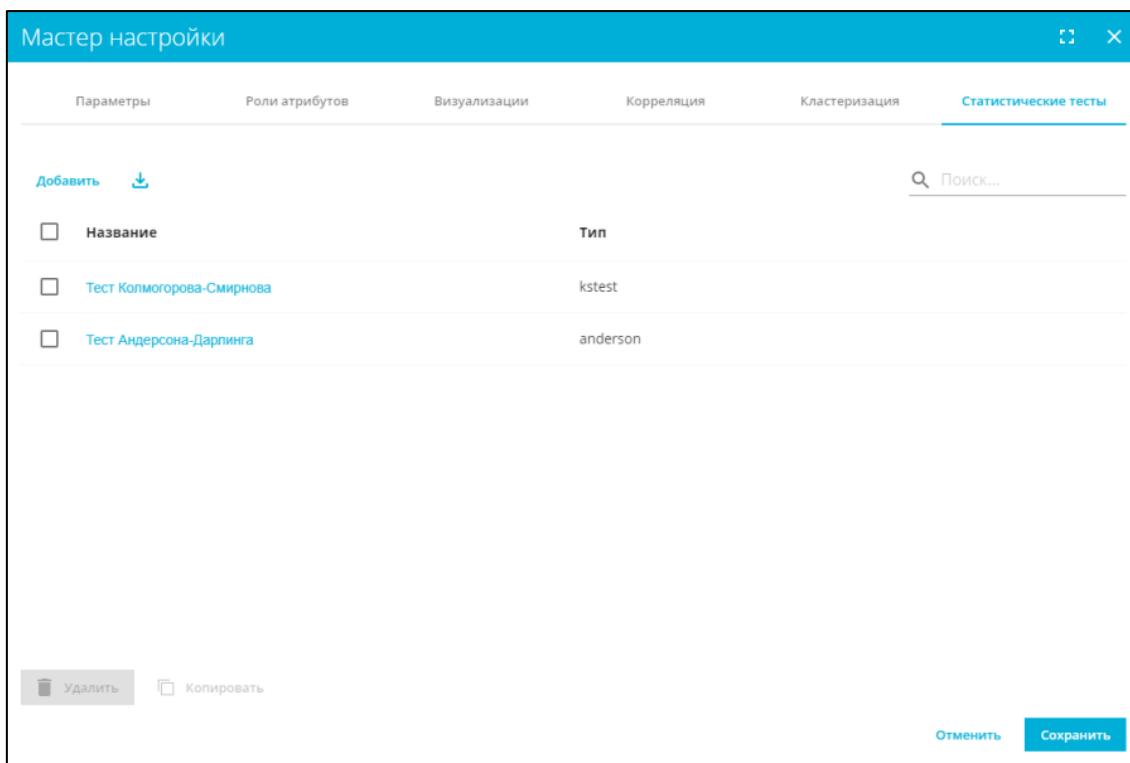
| Параметр                    | Возможные значения и ограничения                                 | Описание   |
|-----------------------------|--|--|
| <b>Название</b>             | Ручной ввод<br>Ограничений на значение нет                       | Название, которое будет отображаться в результатах исследования            |
| <b>Количество кластеров</b> | Ручной ввод целочисленного значения больше 0<br>По умолчанию — 5 | Задание числа кластеров, на которые будет делиться векторное пространство. |
| <b>Random Seed</b>          | Ручной ввод числового значения<br>По умолчанию — 42              | Начальное числовое значение для генератора случайных чисел                 |

| Параметр   | Возможные значения и ограничения  | Описание   |
|--|---|--|
| <b>Тип инициализации</b>                                     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• k-means++ (По умолчанию)</li> <li>• Forgy</li> <li>• Random</li> </ul>                         | Данный параметр отвечает за выбор метода инициализации начальных точек кластеров. Предусмотрены: <ul style="list-style-type: none"> <li>• <b>k-means++</b><br/>Идея метода k-means++ состоит в том, чтобы выбрать начальные точки, которые находятся как можно дальше друг от друга.</li> <li>• <b>Forgy</b><br/>Метод Forgy случайным образом выбирает <math>k</math> наблюдений (по числу заданных кластеров) из набора данных и использует их в качестве начальных значений.</li> <li>• <b>Random</b><br/>Метод Random сначала случайным образом назначает кластер каждому наблюдению, а затем переходит к этапу обновления, таким образом вычисляя начальное среднее значение как центроид случайно назначенных точек кластера.</li> </ul> |
| <b>Тип стандартизации</b>                                    | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Не выбрано</li> <li>• Стандартное отклонение (По умолчанию)</li> <li>• Нормализация</li> </ul> | Данный параметр задает тип стандартизации данных. Предусмотрены: <ul style="list-style-type: none"> <li>• <b>Не выбрано</b><br/>Не стандартизировать</li> <li>• <b>Стандартное отклонение</b><br/>Из каждой записи вычитается среднее значение и результат делится на стандартное отклонение</li> <li>• <b>Нормализация</b><br/>Из каждой записи вычитается минимальное значение и результат делится на разницу между максимальным и минимальным значением</li> </ul>  |
| <b>Подготовить данные</b>                                    | Чекбокс   | Выбор данного чекбокса указывает на необходимость: <ul style="list-style-type: none"> <li>• Заменить пропущенные значения количественной переменной на mean</li> <li>• Заменить пропущенные значения категориальной переменной на текстовый None</li> </ul>  |
| <b>Количество бинов в графике в параллельных координатах</b> | Ручной ввод целочисленного значения<br>Больше 0<br>По умолчанию — 10  | Данный параметр задает количество бинов, на которое делятся наблюдения для отображения на графике в параллельных координатах   |
| <b>Атрибуты</b>  | Список атрибутов, доступных в наборе данных   | Выбор атрибутов набора данных для проведения кластерного анализа   |

Таблица 2 Параметры кластеризации

## 2.2.7. Статистические тесты

Данная вкладка представляет собой список статистических тестов, которые будут рассчитаны в ходе Исследования.



**Рисунок 34 Вкладка «Статистические тесты» Мастера настройки**

Статистические тесты используются для тестирования статистических гипотез о виде распределения наблюдений:

- Нулевая гипотеза: выборка принадлежит некоторому закону распределения;
- Альтернативная гипотеза: нулевая гипотеза не верна.

В компоненте **Исследование данных** предусмотрены следующие статистические тесты.

Для непрерывных наблюдений:

- Тест Колмогорова-Смирнова.
- Тест Крамера-фон-Мизеса.
- Тест Андерсона-Дарлинга.

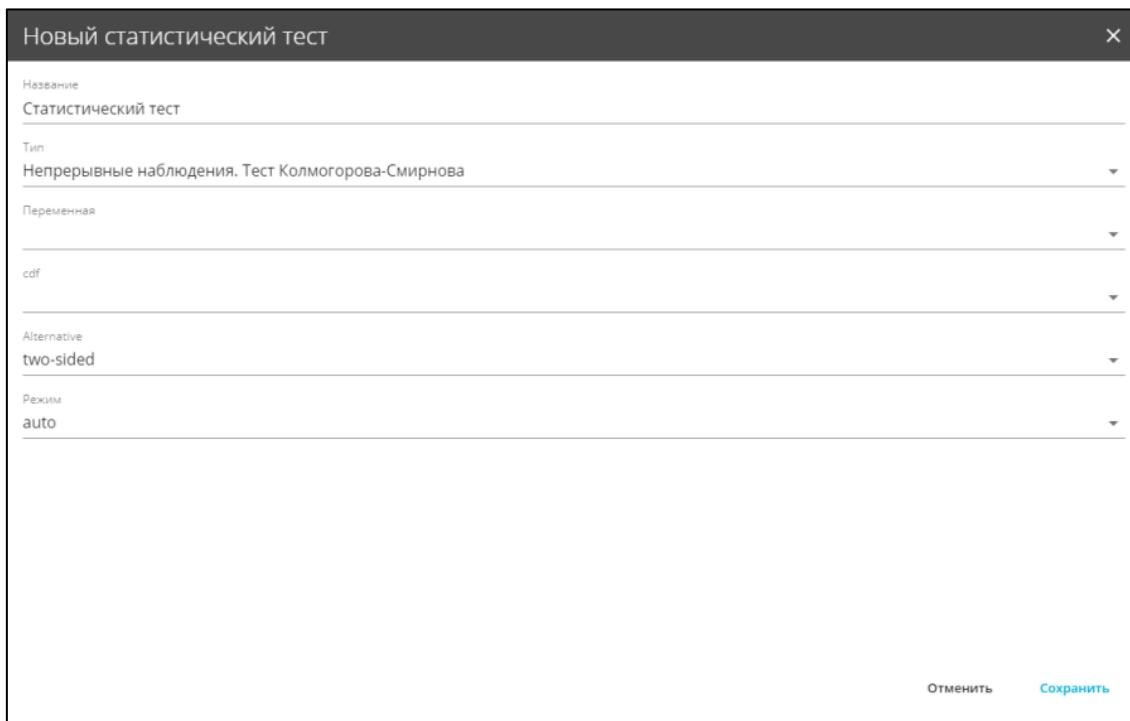
Для дискретных наблюдений:

- Критерий согласия хи-квадрат Пирсона.

Для расчета статистического теста необходимо:

- Выбрать «Добавить».
- В открывшемся окне **Новый статистический тест** задать название и тип теста (подробнее в таблице ниже), а также параметры, специфичные для каждого из статистических тестов.

- Сохранить параметры.



**Рисунок 35 Окно Новый статистический тест**

| Параметр        | Возможные значения и ограничения   | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет   | Название, которое будет отображаться в результатах исследования  |
| <b>Тип</b>      | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>Непрерывные наблюдения. Тест Колмогорова-Смирнова</li> <li>Непрерывные наблюдения. Тест Крамера-фон-Мизеса</li> <li>Непрерывные наблюдения. Тест Андерсона-Дарлинга</li> <li>Дискретные наблюдения. Критерий согласия хи-квадрат Пирсона</li> </ul> | Данный параметр задает тип рассчитываемого статистического теста |

**Таблица 3 Параметры статистического теста**

В окне **Настройки статистического теста** отображается различный набор параметров, зависящий от типа статистического теста.

Для расчета **Теста Колмогорова-Смирнова** нужно указать следующие параметры:

- Переменная, функция распределения  $F(x)$  которой будет проверяться на принадлежность некоторому закону распределения.
- $cdf$  — теоретическое распределение  $G(x)$ .
- Alternative — варианты нулевой и соответствующей альтернативной гипотезы:

- two-sided (двусторонняя) — нулевая гипотеза состоит в том, что два распределения идентичны,  $F(x) = G(x)$  для всех  $x$ ; альтернативная — они не идентичны;
- less (левосторонняя) — Нулевая гипотеза состоит в том, что  $F(x) \geq G(x)$  для всех  $x$ ; альтернативная —  $F(x) < G(x)$  хотя бы для одного  $x$ ;
- greater (правосторонняя) — Нулевая гипотеза состоит в том, что  $F(x) \leq G(x)$  для всех  $x$ ; альтернативная —  $F(x) > G(x)$  хотя бы для одного  $x$ .
- Режим — определяет распределение, используемое для расчета р-значения:
  - auto — автоматический подбор;
  - exact — использует распределение тестовой статистики;
  - approx — приближает двустороннюю вероятность с удвоенной односторонней вероятностью;
  - asympt — использует асимптотическое распределение тестовой статистики.

Для расчета **Теста Крамера-фон-Мизеса** необходимо указать следующие параметры:

- Переменная, функция распределения которой будет проверяться на принадлежность некоторому закону распределения.
- cdf — теоретическое распределение.

Для расчета **Теста Андерсона-Дарлинга** нужно указать следующие параметры:

- Переменная, функция распределения которой будет проверяться на принадлежность некоторому закону распределения.
- Dist — теоретическое распределение.

Для расчета **Критерия согласия хи-квадрат Пирсона** нужно указать следующие параметры:

- Категориальная переменная, ожидаемые частоты которой предполагаются однородными и определяются как среднее значение наблюдаемых частот.
- ddof — поправка на число степеней свободы.

## 2.3. Расписание

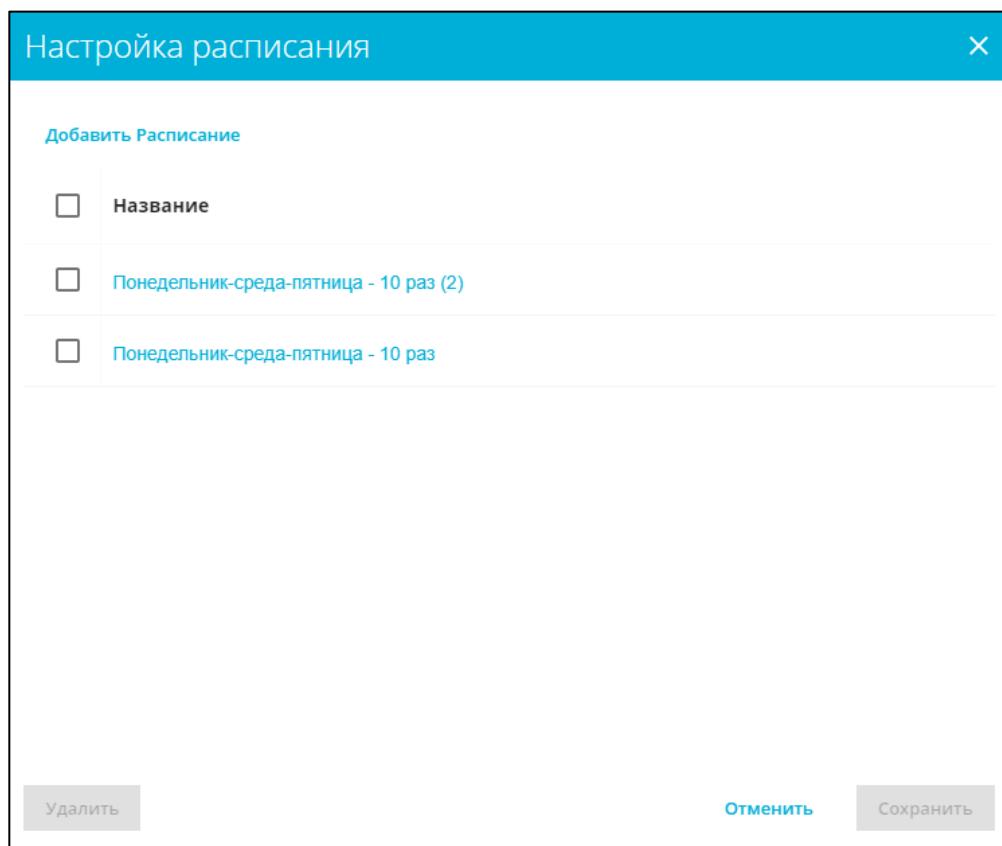
Данная функция позволяет сформировать список расписаний, согласно которым будет запускаться расчет Исследований.

| Результат                           | Описание  | Статус | Создал            | Изменил           |
|-------------------------------------|---|--------|-------------------|-------------------|
| Качество красного вина              | DEMO. Физико-химические данные о красном португальском вине региона Vinho Verde |        | 20.01.2022, 16:18 | 19.07.2022, 11:11 |
| Исследование Качество красного вина | Исследование Качество красного вина   |        | 01.11.2022, 14:47 | 01.11.2022, 14:47 |
| Результат                           | Исследование Качество красного вина   | ●      | 01.11.2022, 14:47 | 01.11.2022, 14:47 |
| Исследование Качество красного вина | Исследование Качество красного вина   |        | 01.11.2022, 14:30 | 01.11.2022, 14:30 |
| Результат                           | Исследование Качество красного вина   | ●      | 01.11.2022, 14:43 | 01.11.2022, 14:43 |
| Результат                           | Исследование Качество красного вина   | ●      | 01.11.2022, 14:30 | 01.11.2022, 14:31 |

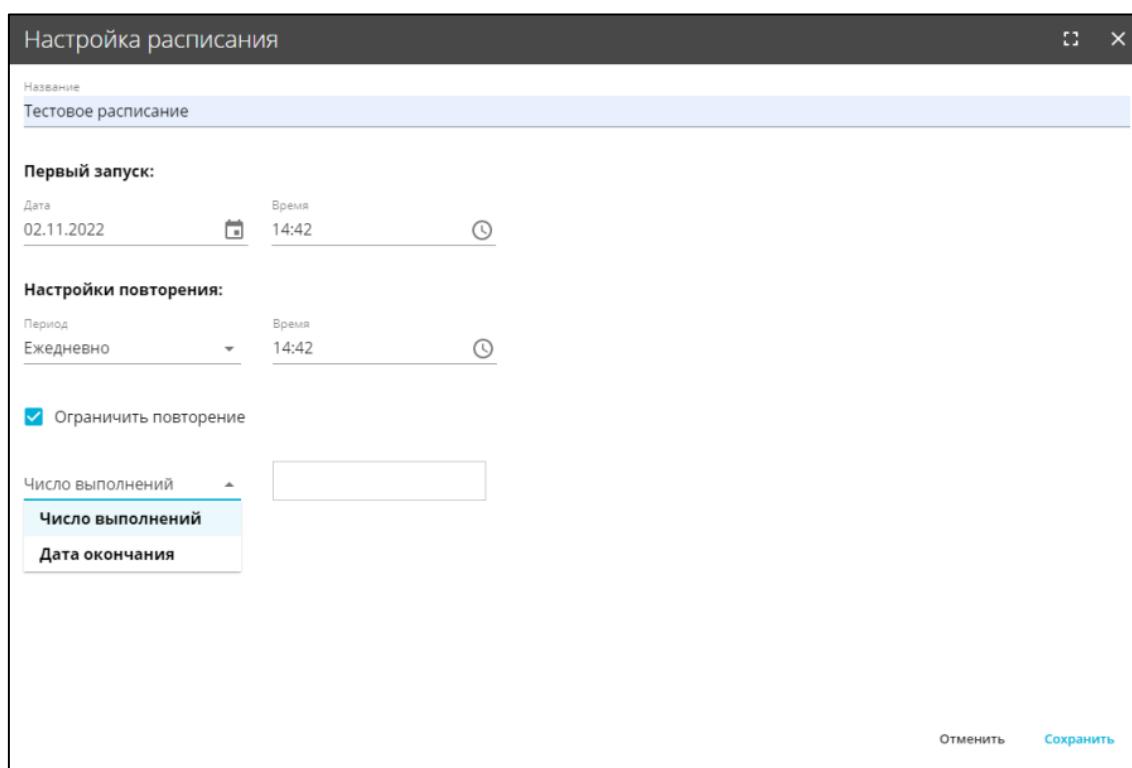
**Рисунок 36 Функция настройки расписания запуска Исследований**

Для постановки Исследования на расписание необходимо:

- Нажать кнопку «**Настройка расписания**».
- Выбрать **Добавить Правило**.
- В открывшемся окне **Настройка расписания** задать:
  - Название, которое далее будет отображаться в системе.
  - Выбрать дату и время первого запуска.
  - Настроить повторение (период — ежедневно, еженедельно, ежемесячно и время повторного запуска).
  - Окончание расписания (Всегда запускать/Число выполнений (указать число запусков)/Дата окончания (указать последнюю дату запуска)).
  - Сохранить настройки.



**Рисунок 37 Окно Настройка расписания**



**Рисунок 38 Окно Настройка расписания**

После постановки Исследования на расписание в колонке **Статус** появится статус **Запланировано** и после каждого расчета результаты будут отображаться ниже по иерархии, никак не влияя на предыдущие и последующие расчеты.

|  |   |   |                      |                      |
|--|---|---|----------------------|----------------------|
|  | Качество белого вина  | Физико-химические данные о белом португальском вине региона Vinho Verde | 15.06.2021, 18:24:52 | 15.06.2021, 18:24:52 |
|  | Исследование Качество белого вина   | Исследование Качество белого вина                                       | 16.09.2021, 17:41:11 | 16.09.2021, 17:41:11 |
|  | Результат Исследование Качество белого вина                                 | Результат Исследование Качество белого вина                             | 15.06.2021, 18:28:56 | 15.06.2021, 18:28:56 |
|  | Запуск по расписанию Запуск по расписанию Исследование Качество белого вина | Запуск по расписанию Исследование Качество белого вина                  | 09.09.2021, 14:51:03 | 09.09.2021, 14:51:03 |
|  | Запуск по расписанию Запуск по расписанию Исследование Качество белого вина | Запуск по расписанию Исследование Качество белого вина                  | 09.09.2021, 15:00:08 | 09.09.2021, 15:00:08 |

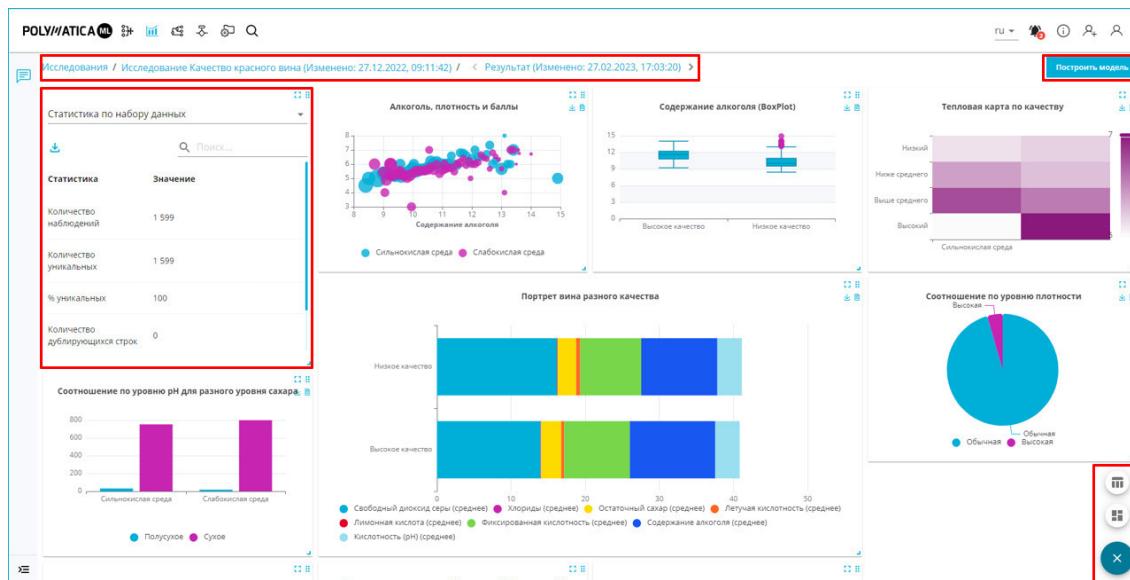
**Рисунок 39 Пример Исследования, поставленного на расписание**

## 2.4. Результаты исследования

### 2.4.1. Интерфейс экрана Результаты Исследования

На экране Результаты Исследования представлены построенные при помощи Мастера настройки визуализации, корреляционные матрицы и другие объекты.

Для перехода на экран Результаты Исследования необходимо на главном экране **компонентаИсследование данных** выбрать необходимый результат исследования.



**Рисунок 40 Пример экрана с результатами исследования и выделенными основными элементами**

На экране с результатами исследования каждый объект расположен в отдельном контейнере. Для каждого из них предусмотрены изменение размера и положения — в верхней части контейнера иконки и , соответственно. Отрегулировать размер контейнера также можно ухватив левой кнопкой мыши правый нижний угол.

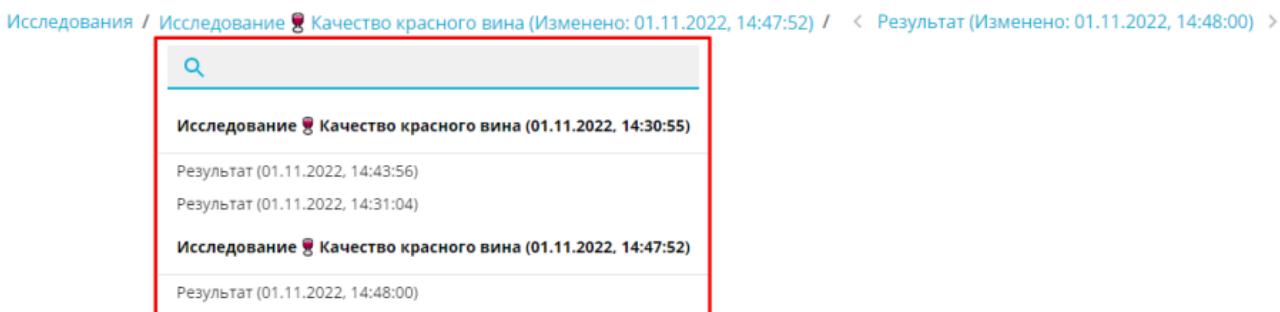
Также предусмотрена возможность сохранить визуализацию в формате .png. Для этого необходимо выбрать иконку в верхней части контейнера, либо при нажатии на контейнер правой кнопкой мыши в выпадающем меню выбрать **Сохранить картинку как...**. Копировать визуализацию в буфер обмена можно в том же выпадающем меню по пункту **Копировать картинку**.

Посмотреть на данные, которые лежат в основе визуализации, можно выбрав иконку в верхней части контейнера.

Дополнительно настроить контейнер можно в выпадающем меню в правом нижнем углу (элемент **Настройка вида** ). Там же можно ознакомиться с примером исходного набора данных (элемент **Пример данных** ) и выгрузить его локально (для этого в верхней части таблицы нужно выбрать иконку экспорта в Excel ).

Визуальные объекты результатов исследования динамически настраиваемы. Они изменяют свой размер в соответствии с размером контейнера. Для многих типов визуализаций предусмотрены скрытие категорий/элементов. При наведении на элементы графиков выводятся всплывающие подсказки с подробной информацией о содержимом.

Для перехода между результатами исследований одного набора данных предусмотрена навигационная цепочка в верхней части экрана. Первая категория возвращает Пользователя на главную страницу компонента Исследование данных, вторая – позволяет выбрать интересующее исследование в рамках одного набора данных в выпадающем меню, тут же можно выбрать и интересующий результат, третья категория позволяет выбрать результат в рамках интересующего исследования.



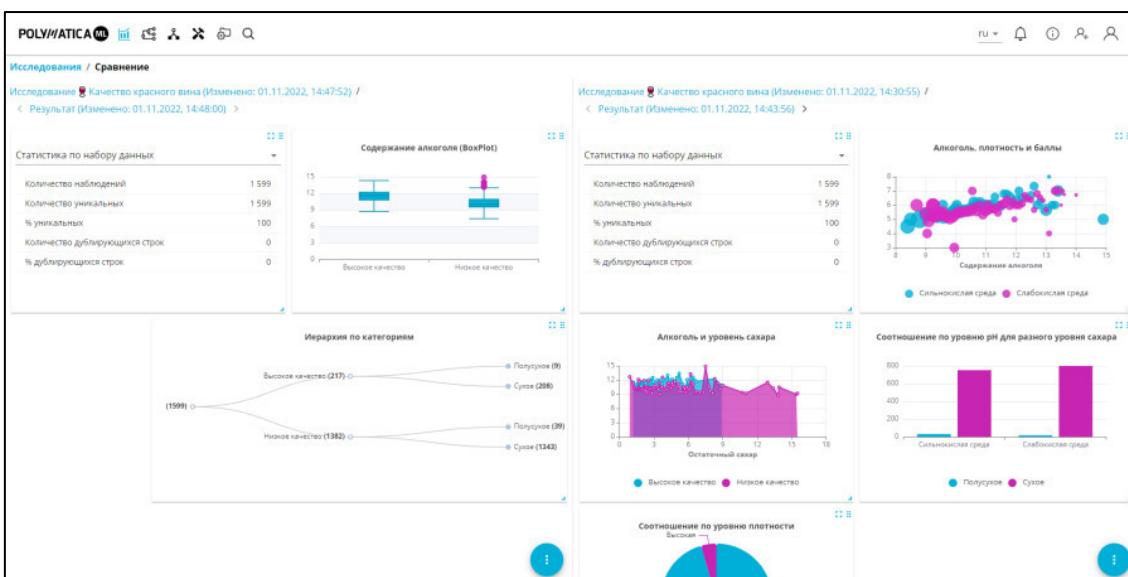
**Рисунок 41 Экран с результатами нескольких исследований**

На экране с результатами Исследования можно отобразить несколько результатов.

Для этого на Главном экране компонента необходимо выбрать чекбоксы  рядом с двумя интересующими расчетами и выбрать «**Сравнить**» в нижней части таблицы (Рисунок 39). Таким образом, откроется экран, содержащий два результата исследования.

The screenshot shows a table titled 'Результаты (3)' (Results) with three items listed. The first item is 'Качество красного вина' (Red Wine Quality) with a status of 'Изменено' (Changed). The second item is 'Исследование Качество красного вина' (Research Red Wine Quality) with a status of 'Создан' (Created). The third item is another 'Исследование Качество красного вина' (Research Red Wine Quality) with a status of 'Создан'. Below the table are two buttons: 'Удалить' (Delete) and 'Сравнить' (Compare), with the 'Сравнить' button highlighted by a red box.

**Рисунок 42 Пример выбора двух результатов для сравнения**



**Рисунок 43 Экран с результатами нескольких исследований**

## 2.4.2. Профилирование

На экране с результатами исследования также представлен посчитанный в автоматическом режиме **Профиль данных**.

| Статистика по набору данных    |       |
|--------------------------------|-------|
| Количество наблюдений          | 1 599 |
| Количество уникальных          | 1 599 |
| % уникальных                   | 100   |
| Количество дублирующихся строк | 0     |
| % дублирующихся строк          | 0     |

**Рисунок 44 Пример Профиля данных**

**Профилирование** представляет собой анализ данных с целью выяснения статистических характеристик переменных, таких как характер распределения, наличие выбросов, параметры выборки, пропущенные значения и другие.

Статистики считаются по всему набору данных и по каждой из переменных и зависят от ее типа.

### 2.4.2.1. Набор данных

Результат профилирования для всего набора данных зависит от размера анализируемого файла.

Так для набора данных размером менее 1 Гб считается следующий набор статистик:

- Количество наблюдений.
- Количество уникальных наблюдений.
- Процент уникальных наблюдений (в процентах).
- Количество дублирующих строк.
- Процент дублирующих строк (в процентах).
- Количество количественных переменных.
- Количество категориальных переменных.
- Количество переменных дат.

Для набора данных размером более 1 Гб:

- Количество наблюдений.
- Количество количественных переменных.
- Количество категориальных переменных.
- Количество переменных дат.

#### 2.4.2.2. Категориальные переменные (Тип данных String)

Для категориальных переменных считается набор статистик в соответствии с таблицей ниже.

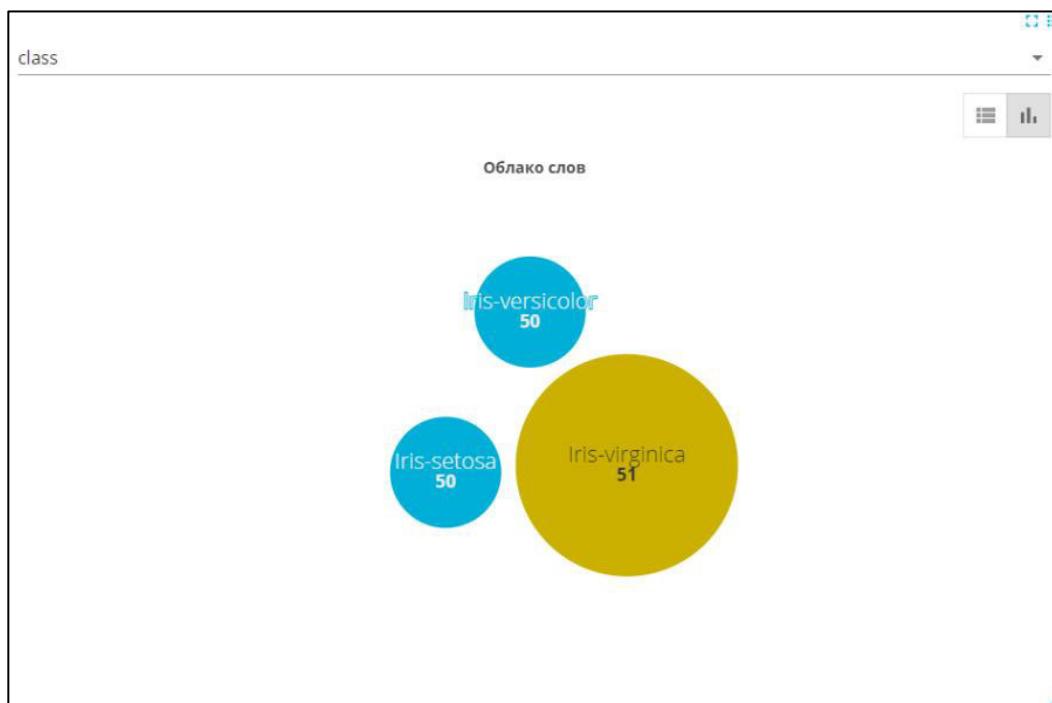
| Статистика                             | Описание   |
|--|--|
| <b>Количество уникальных значений</b>  | Количество неповторяющихся значений                        |
| <b>Процент уникальных значений</b>     | Процент уникальных значений от общего количества значений  |
| <b>Количество пропущенных значений</b> | Количество пропущенных значений                            |
| <b>Процент пропущенных значений</b>    | Процент пропущенных значений от общего количества значений |
| <b>Максимальная длина</b>              | Количество знаков максимального по длине значения          |
| <b>Минимальная длина</b>               | Количество знаков минимального по длине значения           |
| <b>Top</b>                             | 3 значения с максимальной частотой                         |
| <b>Bottom</b>                          | 3 значения с минимальной частотой                          |

**Таблица 4 Набор статистик, рассчитываемых для категориальных переменных**

Для удобства интерпретации и анализа также предусмотрено построение **Облака слов**.

**Облако слов** (или облако тегов) демонстрирует частотность появления значения переменной. Размер облака отражает частоту появления значения. Цветовая гамма не несет в себе дополнительного смысла и выполняет исключительно эстетическую функцию.

Посмотреть **Облако слов** можно в том же контейнере с профилированием, выбрав в правом верхнем углу иконку  .



**Рисунок 45 Пример Облака слов**

#### 2.4.2.3. Качественные переменные (Тип данных Numeric)

Для количественных переменных считается набор статистик в соответствии с таблицей ниже.

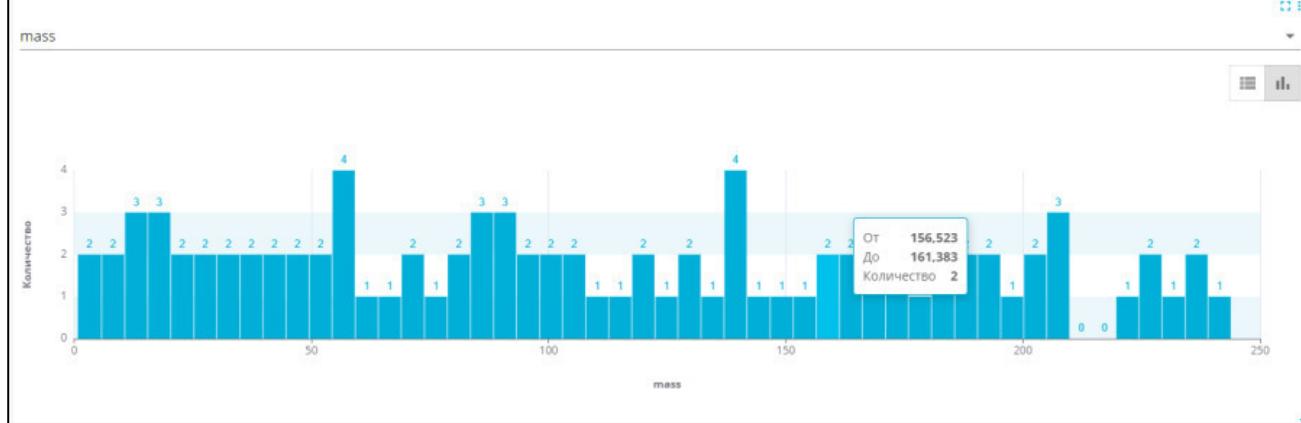
| Статистика                             | Описание   |
|--|--|
| <b>Количество уникальных значений</b>  | количество неповторяющихся значений переменной   |
| <b>Процент уникальных значений</b>     | процент неповторяющихся значений переменной  |
| <b>Количество пропущенных значений</b> | количество пропущенных значений переменной   |
| <b>Процент пропущенных значений</b>    | процент пропущенных значений   |
| <b>Минимальное значение</b>            | наименьшее значение переменной   |
| <b>Максимальное значение</b>           | наибольшее значение переменной   |
| <b>Среднее значение</b>                | сумма всех значений переменной, разделенная на число этих значений   |
| <b>5-я персентиль</b>                  | это некоторое значение Х из данного ряда, которое делит все имеющиеся в нем значения на две группы: 5% значений, которые меньше Х, и оставшиеся значения (то есть 95%), которые превышают Х.   |
| <b>95-я персентиль</b>                 | это некоторое значение Х из данного ряда, которое делит все имеющиеся в нем значения на две группы: 95% значений, которые меньше Х, и оставшиеся значения (то есть 5%), которые превышают Х.   |
| <b>1-я квартиль</b>                    | 25-я персентиль  |
| <b>3-я квартиль</b>                    | 75-я персентиль  |
| <b>Медиана</b>                         | значение, которое делит распределение пополам (его площадь в т.ч.): половина значений больше медианы, половина — не больше.  |
| <b>Межквартильный размах</b>           | разница между 3-м и 1-м квартилями   |
| <b>Количество выбросов</b>             | Выбросами считаются наблюдения, которые отклоняются от своего математического ожидания более чем на три среднеквадратических отклонения (правило трех сигм).   |
| <b>Коэффициент вариации</b>            | величина, равная отношению стандартного (среднеквадратичного) отклонения случайной величины к ее математическому ожиданию. Он применяется для сравнения вариативности одного и того же признака в нескольких совокупностях с различным средним арифметическим. Если значение коэффициента вариации не превышает 33%, то совокупность считается однородной, а если больше 33%, то — неоднородной. |
| <b>Коэффициент эксцесса</b>            | характеризует меру высоты графика. Если коэффициент больше нуля, то распределение является более высоким («островершинным») относительно «эталонного» нормального распределения. Если коэффициент ниже нуля, то более низким и пологим.  |

| Статистика                             | Описание   |
|--|--|
| <b>Медианное абсолютное отклонение</b> | вычисляется как медиана абсолютного значения для каждого значения минус медианное значение группы.<br>Является статистикой, более устойчивой к выбросам в наборе данных, чем стандартное отклонение  |
| <b>Асимметрия</b>                      | характеризует меру скошенности графика влево/вправо. Если коэффициент асимметрии отрицателен, то скос левосторонний. Если коэффициент положителен, то скос правосторонний. И чем коэффициент больше по модулю, тем сильнее скос распределения.   |
| <b>Стандартное отклонение</b>          | статистическая характеристика распределения случайной величины, показывающая среднюю степень разброса значений величины относительно математического ожидания. Большее значение среднеквадратического отклонения показывает больший разброс наблюдаемых значений признака относительно среднего; меньшее значение, соответственно, показывает, что величины в множестве сгруппированы вокруг среднего. |
| <b>Дисперсия</b>                       | величина, которая характеризует меру разброса значений случайной величины относительно ее математического ожидания.  |
| <b>T-статистика</b>                    | T-статистика — это разница между средним по выборке и гипотетическим средним (предполагаемым равным нулю), деленная на расчетную стандартную ошибку среднего.  |
| <b>Пи-значение</b>                     | Уровень значимости — вероятность получить Т-значение, равное или превышающее то значение, которое мы в действительности рассчитали по имеющимся выборочным данным (при условии, что нулевая гипотеза верна)  |

**Таблица 5 Набор статистик, рассчитываемых для количественных переменных**

Для удобства интерпретации и анализа предусмотрено построение **Гистограммы**.

**Гистограмма** визуализирует распределение данных в рамках непрерывного интервала. Каждая полоса представляет в табличной форме частотность за определенный бин. Количество бинов не изменяется и по умолчанию равно 50. Посмотреть гистограмму можно в том же контейнере, выбрав в правом верхнем углу иконку .



**Рисунок 46 Пример Гистограммы**

### 2.4.3. Результат расчета корреляции

Результатом расчета является корреляционная матрица, представляющая собой квадратную тепловую карту, в которой на пересечении выбранных признаков (одинаковый набор для столбцов и строк) находится значение коэффициента корреляции, закодированного при помощи градиентной шкалы. Синий цвет отражает максимальное значение коэффициента корреляции (равный 1), красный — минимальное (равный -1).

Матрица является симметричной, с единичной диагональю. При наведении указателя мыши на ячейку можно получить детализированную информацию:

- Коэффициент корреляции.
- Коэффициент ковариации.
- p-value — уровень значимости.

**Коэффициент ковариации** — это мера линейной зависимости двух случайных величин. Она является ненормированной версией коэффициента корреляции.

**Коэффициент корреляции** — статистическая мера, которая отражает силу связи между двумя порядковыми признаками. Коэффициент корреляции может принимать значения от -1 до +1. Если значение по модулю находится ближе к 1, то это означает наличие сильной связи, а если ближе к 0 — связь слабая или вообще отсутствует.

Вычисление коэффициента корреляции производится тремя методами:

- Корреляция Пирсона.
- Корреляция Спирмена.
- Корреляция Кендалла.

Для каждого метода предусмотрено построение своей корреляционной матрицы. Выбрать требуемую можно в выпадающем меню в верхней части панели. Ползунки боковой градиентной шкалы позволяют фильтровать ячейки корреляционной матрицы.

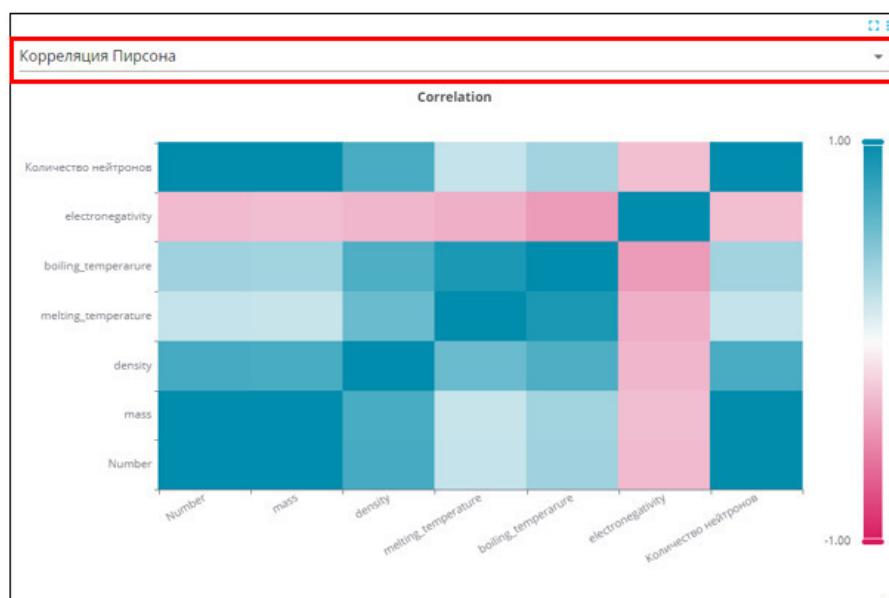


Рисунок 47 Результат расчета корреляционной матрицы

## 2.4.4. Результат расчета кластеризации

Результатом расчета является разбиение совокупности наблюдений на однородные группы, или кластеры.

Для удобства интерпретации результаты кластерного анализа представлены в отдельном контейнере со следующими объектами:

- **Круговая диаграмма** с количеством наблюдений в каждом кластере.

При наведении курсора мыши на сектор кластера можно узнать количество наблюдений в нем.

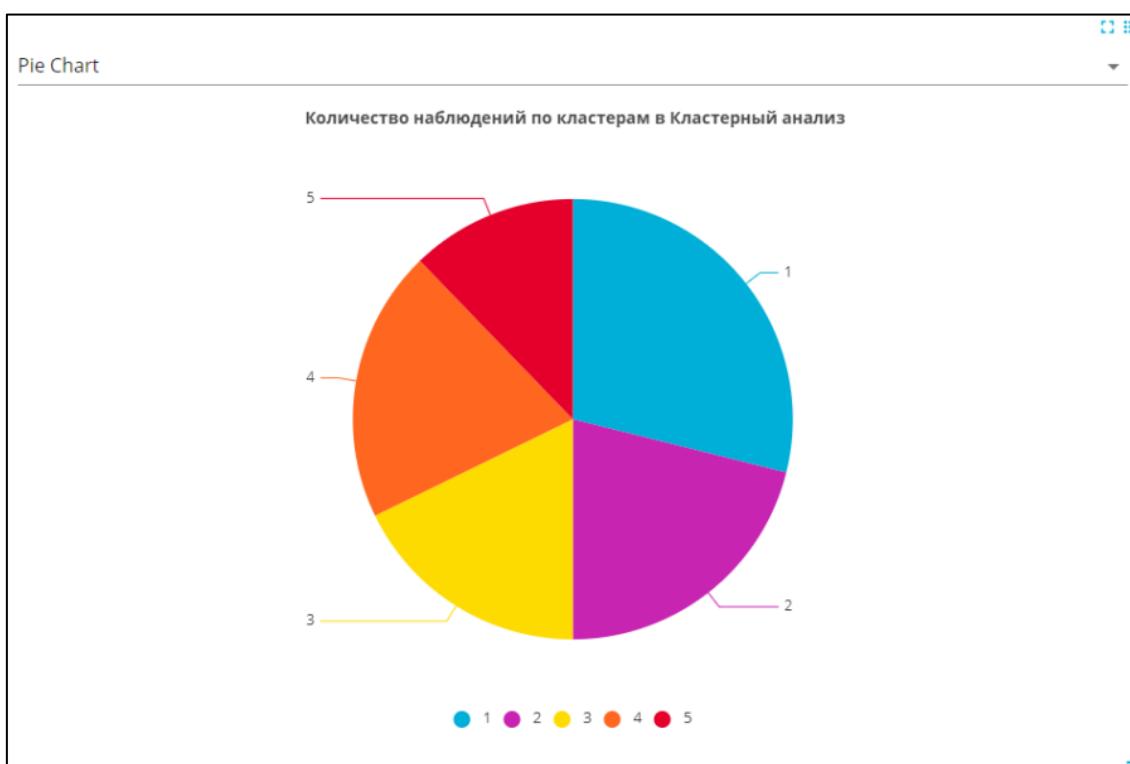
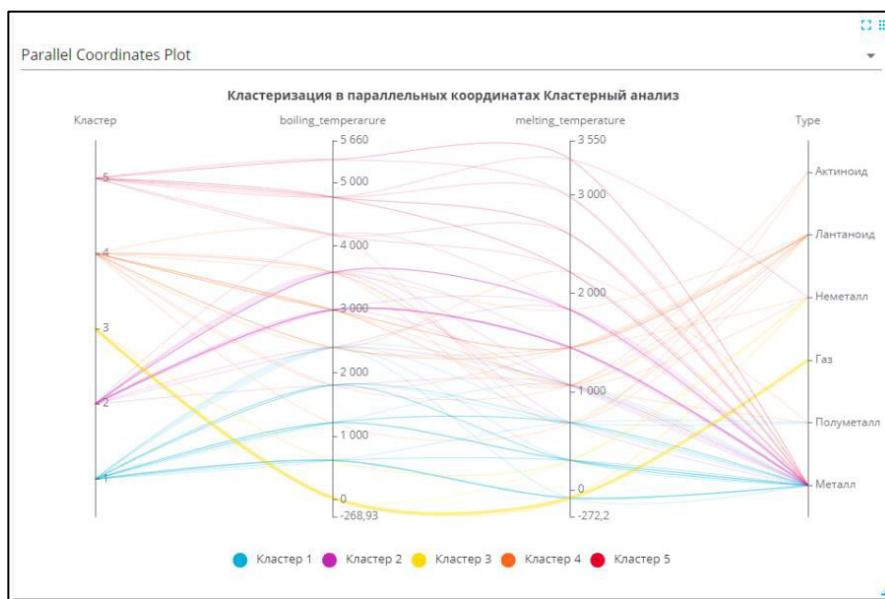


Рисунок 48 Пример Круговой диаграммы с результатами кластеризации

- **График в параллельных координатах** — Parallel Coordinates Plot.

Диаграмма с параллельными координатами позволяет интерпретировать построенные кластеры.

На диаграмме с параллельными координатами каждой переменной присваивается собственная ось. Оси располагаются параллельно друг другу, и каждая имеет свою собственную шкалу. Начальная ось отражает кластер, к которому модель отнесла наблюдение. Каждое наблюдение наносится на график в виде линии, пересекающейся с каждой из осей. Таким образом, пользователь может выявить паттерны и корреляции между разными переменными.



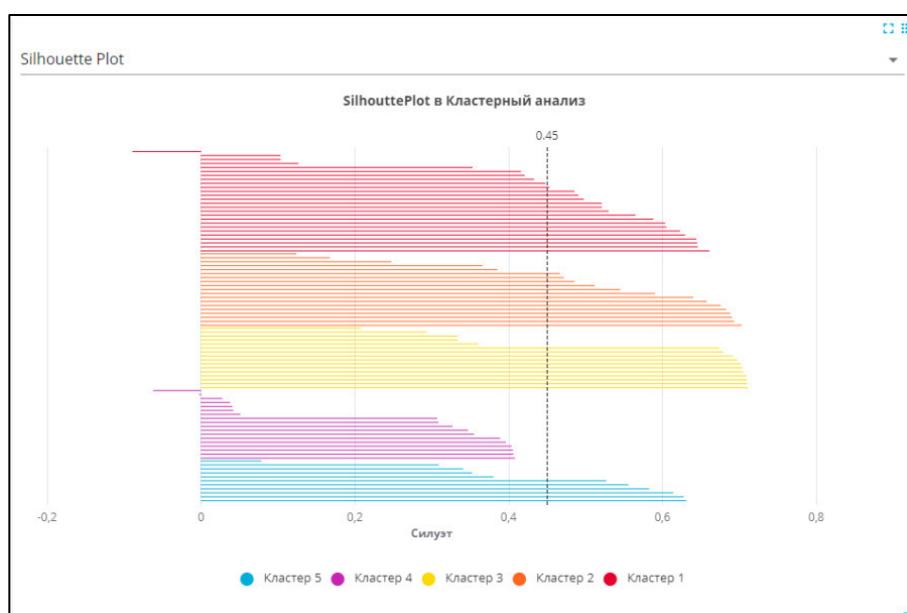
**Рисунок 49 Пример Графика в параллельных координатах**

- **Силуэт — Silhouette Plot.**

**Значение Silhouette** для каждого наблюдения является мерой того, насколько это наблюдение похоже на наблюдения в собственном кластере по сравнению с наблюдениями в других кластерах.

**Значение Silhouette** находится в диапазоне от -1 до 1. Высокое значение указывает на то, что наблюдение хорошо соответствует собственному кластеру и плохо соответствует другим кластерам.

Если большинство наблюдений имеют низкое или отрицательное значение Silhouette, тогда пользователь должен перестроить кластеризацию с большим или меньшим количеством кластеров.



**Рисунок 50 Пример Silhouette Plot**

- **Таблица с координатами центроидов**, где в качестве строк выступают номера кластеров и значения переменных, в которых находятся центроиды этих кластеров.

| Координаты центроидов |                     |                     |
|-----------------------|---------------------|---------------------|
| Кластер               | boiling_temperature | melting_temperature |
| 1                     | 1 356,275           | 422,277             |
| 2                     | 3 149,105           | 1 478,281           |
| 3                     | -13,97              | -98,588             |
| 4                     | 2 860,389           | 1 239,689           |
| 5                     | 4 857,182           | 2 762,273           |

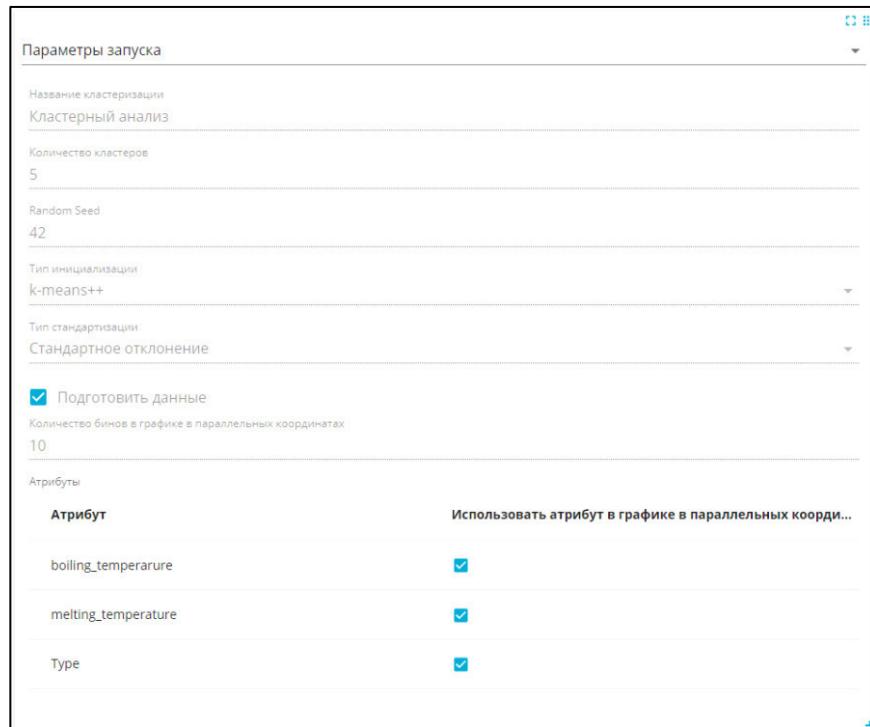
**Рисунок 51 Пример таблицы с координатами центроидов**

- **Таблица со статистиками по кластерам**. Содержит следующие статистики:
  - Номер кластера.
  - Количество наблюдений.
  - Среднеквадратичное расстояние между наблюдениями внутри кластера.
  - Сумма расстояний между наблюдениями внутри кластера.
  - Расстояние между центроидом и ближайшим наблюдением.
  - Расстояние между центроидом и наиболее удаленным наблюдением.
  - Расстояние между центроидом и вторым по удаленности наблюдением.
  - Расстояние между центроидом и третьим по удаленности наблюдением.
  - Ближайший кластер.
  - Расстояние до ближайшего центроида.
  - Среднее расстояние между центроидом и наблюдениями в кластере.
  - Сумма расстояний между наблюдениями и центроидом.

| Статистика по кластерам |                       |  |   |   |  |   |  |                   |                                    |  |  |
|-------------------------|-----------------------|--|---|---|--|---|--|-------------------|------------------------------------|--|--|
| Номер кластера          | Количество наблюдений | Среднеквадратичное расстояние между наблюдениями внутри кластера | Сумма расстояний между наблюдениями внутри кластера | Расстояние между центроидом и ближайшим наблюдением | Расстояние между центроидом и наиболее удаленным наблюдением | Расстояние между центроидом и вторым по удаленности наблюдением | Расстояние между центроидом и третьим по удаленности наблюдением | Ближайший кластер | Расстояние до ближайшего центроида | Среднее расстояние между центроидом и наблюдением в кластере | Сумма расстояний между наблюдениями и центроидом |
| 1                       | 26                    | 0,362  | 6,558   | 0,143   | 0,802  | 0,792   | 0,783  | 3                 | 1,029                              | 0,464  | 12,056   |
| 2                       | 19                    | 0,37   | 4,931   | 0,08  | 0,909  | 0,872   | 0,789  | 4                 | 0,315                              | 0,447  | 8,49   |
| 3                       | 16                    | 0,17   | 0,665   | 0,013   | 0,556  | 0,368   | 0,26   | 1                 | 1,029                              | 0,194  | 3,098  |
| 4                       | 18                    | 0,429  | 6,255   | 0,207   | 1,144  | 1,057   | 0,921  | 2                 | 0,315                              | 0,519  | 9,343  |
| 5                       | 11                    | 0,451  | 4,062   | 0,22  | 0,861  | 0,848   | 0,775  | 2                 | 1,751                              | 0,571  | 6,277  |

**Рисунок 52 Пример таблицы со статистиками кластеров**

- **Информация по параметрам запуска.** Соответствует интерфейсу окна **Параметры кластеризации** в **Мастере настройки**.



**Рисунок 53 Пример вкладки с параметрами запуска кластерного анализа**

- **Таблица со статистиками по переменным кластера.** По каждому кластеру отражены среднее и стандартное отклонение для каждой переменной.

| Имя параметра       | Стандартное отклонение | Среднее |
|---------------------|------------------------|---------|
| ▼ Номер кластера: 1 | boiling_temperature    | 596,7   |
|                     | melting_temperature    | 322,285 |
| ▼ Номер кластера: 2 | boiling_temperature    | 557,935 |
|                     | melting_temperature    | 360,444 |
| ▼ Номер кластера: 3 | boiling_temperature    | 277,814 |
|                     | melting_temperature    | 152,825 |
| ➤ Номер кластера: 4 |                        |         |
| ➤ Номер кластера: 5 |                        |         |

**Рисунок 54 Пример таблицы со статистиками переменных по кластерам**

## 2.4.5. Результат расчета статистических тестов

Результаты всех рассчитанных в исследовании статистических тестов отображаются в едином контейнере.

Каждый из статистических тестов имеет свои результаты. Выбрать результаты интересующего можно в выпадающем меню в верхней части контейнера.

| Тест Колмогорова-Смирнова |   |
|---------------------------|---|
| Название                  | Тест Колмогорова-Смирнова                         |
| Тип                       | Непрерывные наблюдения. Тест Колмогорова-Смирнова |
| Переменная                | sepal_length_in_cm                                |
| cdf                       | norm  |
| Alternative               | two-sided   |
| Режим                     | auto  |
| statistic                 | 1   |
| pvalue                    | 0   |

**Рисунок 55 Контейнер с результатами статистических тестов**

Помимо расчетных параметров, заданных при создании статистического теста, в данном контейнере отражаются следующие результаты:

- **Тест Колмогорова-Смирнова.**
  - statistic — статистика Колмогорова-Смирнова.
  - p-value — уровень значимости.
- **Тест Крамера-фон-Мизеса.**
  - statistic — статистика теста Крамера-фон-Мизеса.
  - p-value — уровень значимости.
- **Тест Андерсона-Дарлинга.**
  - statistic — статистика теста Андерсона-Дарлинга.
  - critical\_values — критические значения для этого распределения.
  - significance\_level — уровни значимости для соответствующих критических значений в процентах. Функция возвращает критические значения для различного набора уровней значимости в зависимости от тестируемого распределения. (Если возвращаемая статистика превышает эти критические значения, то для соответствующего уровня

значимости нулевая гипотеза о том, что данные принадлежат выбранному распределению, может быть отклонена.).

- normal/exponential — 15%, 10%, 5%, 2.5%, 1%.
- logistic — 25%, 10%, 5%, 2.5%, 1%, 0.5%.
- Gumbel — 25%, 10%, 5%, 2.5%, 1%.

- **Критерий согласия хи-квадрат Пирсона.**

- statistic — статистика критерий согласия хи-квадрат Пирсона.
- p-value — уровень значимости.

## 2.4.6. Автоматическое построение модели в MD

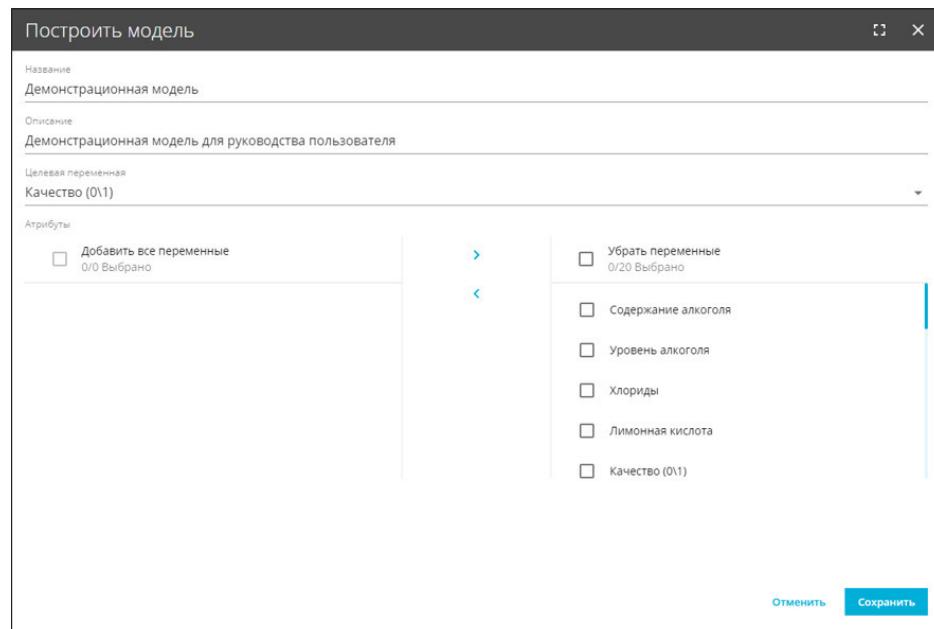
В рамках модуля DD существует возможность быстрого построения пайплайна моделирования в MD.

Для этого на главном экране необходимо выбрать результат интересующего набора данных и на открывшемся экране выбрать кнопку **Построить модель**.



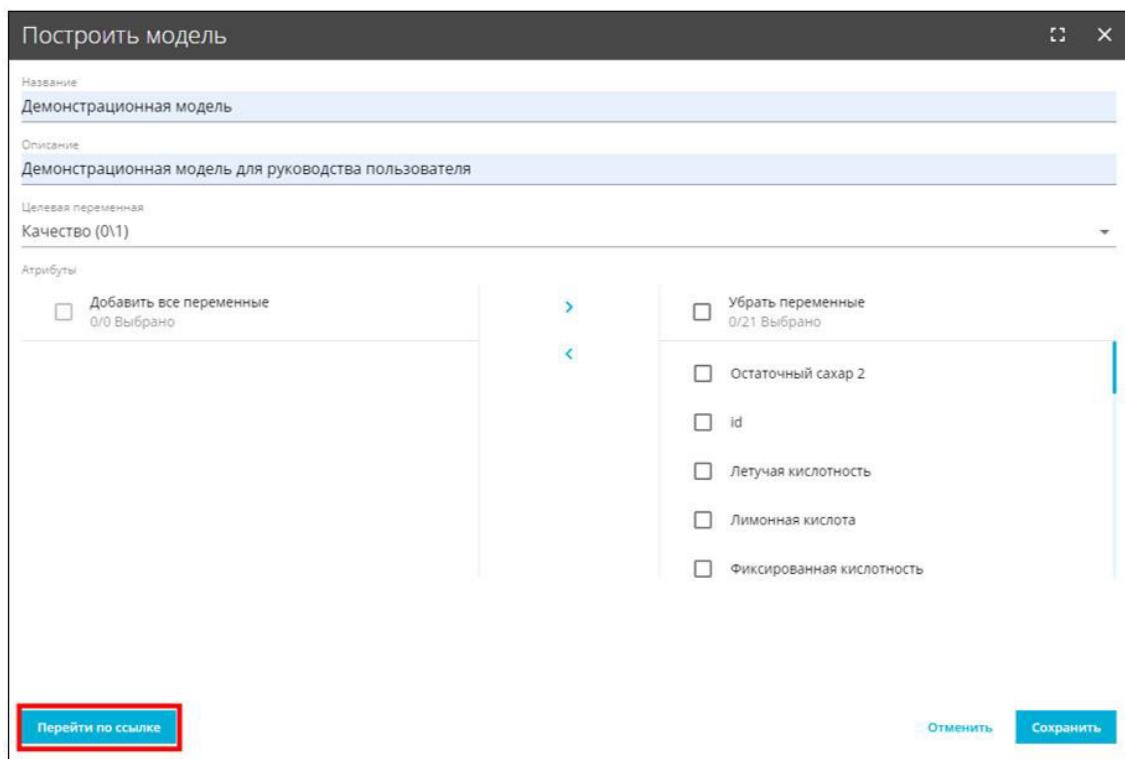
**Рисунок 56 Расположение кнопки Построить модель**

В окне **Построить модель** необходимо задать Название, Описание будущего проекта MD, а также выбрать Целевую переменную и отобрать необходимые для моделирования атрибуты из исходного набора данных. Нажать кнопку **Сохранить**.



**Рисунок 57 Окно Построить модель**

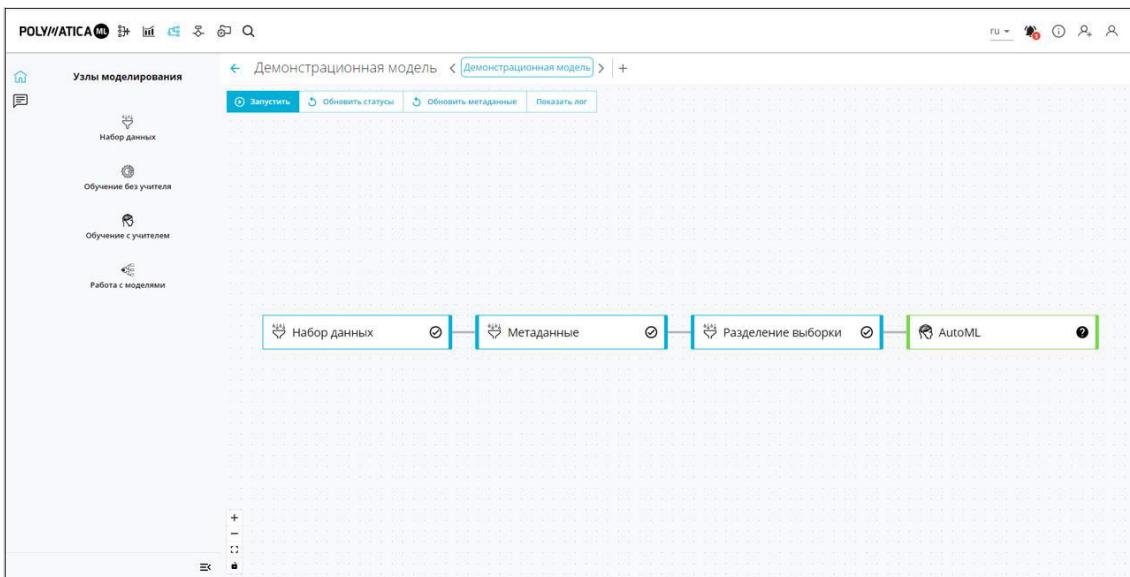
После выполнения необходимых расчетов будет доступна кнопка **Перейти по ссылке**, которая переведет Пользователя к созданному проекту MD.



**Рисунок 58 Кнопка Перейти по ссылке**

В автоматически созданном проекте предусмотрены следующие узлы моделирования:

- Узел "Набор данных" (подробнее [Узел «Набор данных»](#)).
- Узел "Метаданные" (подробнее [Узел «Метаданные»](#)).
- Узел "Разделение выборки" (подробнее [Узел «Разделение выборки»](#)).
- Узел "AutoML" (подробнее [Узел «AutoML»](#)).



**Рисунок 59 Пример созданного проекта МД**

### 3. Компонент Построение моделей (Model Designer, MD)

Этап построения моделей включает в себя обширный пласт работ. Первоначально необходимо нормализовать и доработать данные согласно полученной на этапе их исследования информации и решаемой задачи. Сам процесс обучения происходит итерационно — пробуются разные модели, перебираются гиперпараметры, сравниваются значения выбранной метрики качества и выбирается лучшая комбинация. Также перед переходом к внедрению нужно убедиться, что результат моделирования понятен и логичен.

С учетом этих шагов в компоненте Построение моделей предусмотрены инструменты:

- Подготовки данных (обработка пропусков, расчет новых показателей, фильтрация и т.д.).
- Алгоритмы машинного обучения (включающие в себя широкий список методов обучения с учителем и обучения без учителя).
- Разделение набора данных на выборки для обучения, валидации и тестирования.
- Кросс-валидация.
- Автоподбор параметров.
- Оценки результатов (метрики и диаграммы оценки с учетом типа задачи).
- Сравнения моделей.
- Интерпретации моделей.
- Регистрация моделей в Репозиторий.

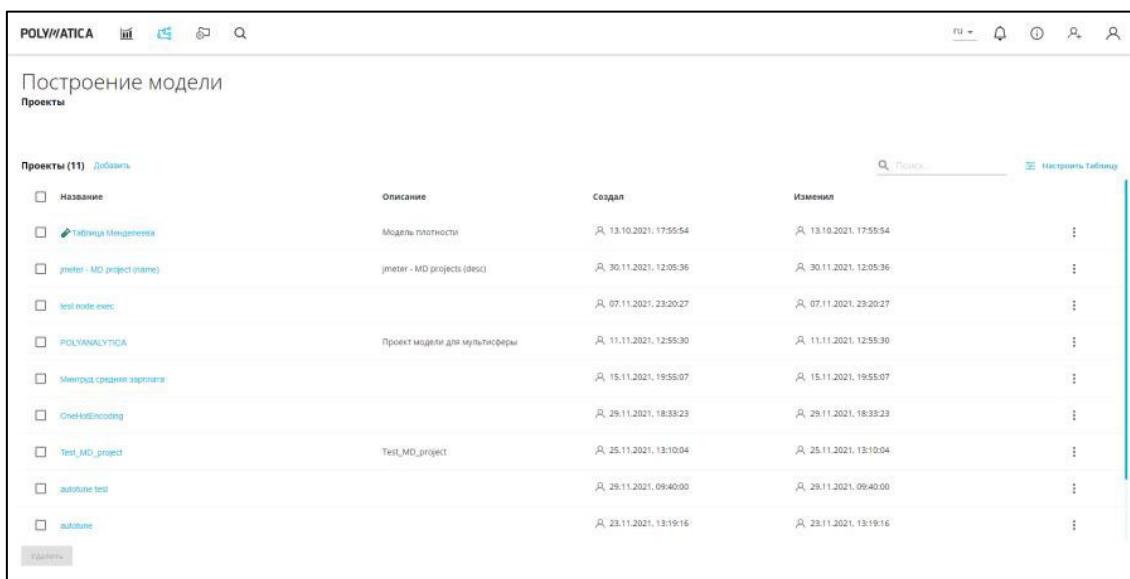
Компонент включает в себя:

- Главный экран со списком доступных проектов моделирования (проектов MD).
- Конструктор сценариев.
- Вспомогательные окна настройки вида, создания проекта и т.д.

#### 3.1. Главный экран MD

##### 3.1.1. Интерфейс главного экрана MD

Главный экран компонента **Построение моделей** открывается при выборе иконки  в левом верхнем меню и представляет собой список доступных пользователю **Проектов** в табличном виде.



**Рисунок 60 Главный экран компонента Построение моделей (Model Designer)**

**Проект MD** — это сценарий моделирования, настраиваемый пользователем для решения конкретной задачи.

Таблица с доступными проектами имеет гибкие настройки отображения. Так, пользователь может:

- Изменить ширину любого столбца (для этого необходимо перетащить границу его заголовка до нужной ширины).
- Сортировать таблицу (для этого необходимо выбрать иконку ↑ рядом с заголовком сортируемого столбца).
- Скрывать/отображать столбцы и изменять их порядок в окне **Вид таблицы** (для открытия окна необходимо выбрать в правом верхнем углу таблицы; при выборе иконки столбец скрывается, при наведении на иконку активируется возможность перемещения столбца).
- Сбросить внесенные изменения также в окне **Вид таблицы** (для этого выбрать кнопку «**Сбросить**»).

Для быстрого поиска объекта в таблице предусмотрено поле в правой верхней части таблицы.

Объекты таблицы можно выгрузить в формате Excel. Для этого нужно выбрать иконку **Экспорта в excel** .

### 3.1.2. Создание проекта MD

Для создания нового проекта необходимо выполнить следующие шаги:

- Выбрать кнопку «**Добавить**» в верхней части таблицы с Проектами.
- В открывшемся окне **Создание проекта** задать название и описание.
- Сохранить изменения.

### 3.1.3. Удаление проекта MD

Для удаления проекта необходимо выполнить следующие шаги:

- Выбрать чекбокс рядом с удаляемым проектом.
- Выбрать кнопку «**Удалить**».

### 3.1.4. Копирование проекта MD

Для копирования проекта необходимо выполнить следующие шаги:

- С правой стороны от проекта выбрать иконку в виде трех вертикальных точек и в выпадающем меню выбрать **Копировать**.
- В открывшемся окне **Копирование проекта** задать название и описание нового проекта.
- Сохранить изменения.

### 3.1.5. Редактирование проекта MD

Для переименования и изменения описания проекта необходимо выполнить следующие шаги:

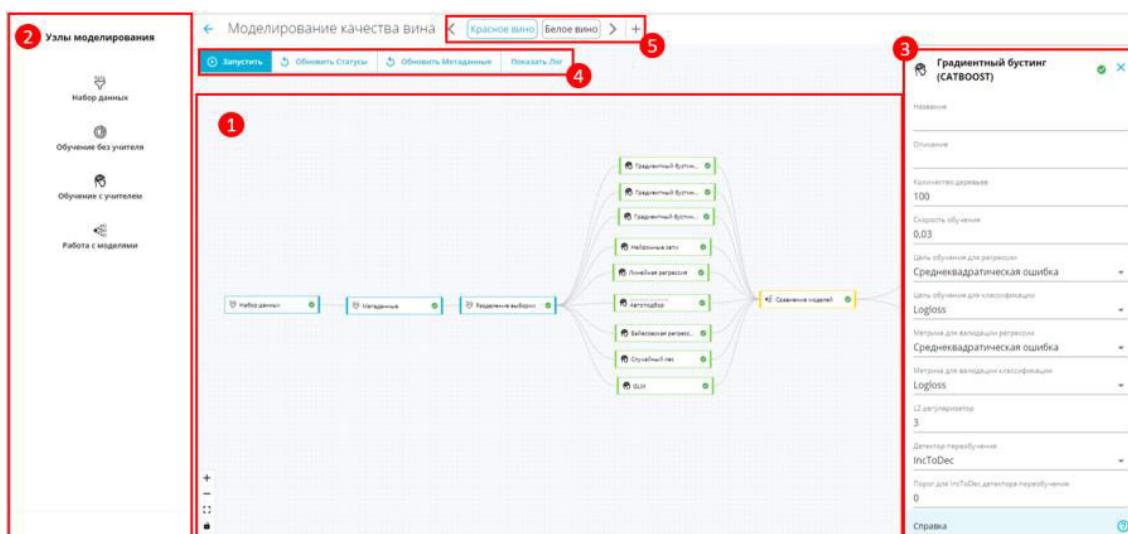
- С правой стороны от проекта выбрать иконку в виде трех вертикальных точек и в выпадающем меню выбрать **Редактировать**.
- В открывшемся окне **Редактирование проекта** внести необходимые изменения в названии и описании проекта.
- Сохранить изменения.

### 3.1.6. Работа с проектом MD

Для начала работы с проектом необходимо выбрать его из списка доступных.

## 3.2. Конструктор сценариев

При выборе проекта на главной странице компонента открывается **Конструктор сценариев**.



**Рисунок 61 Элементы Конструктора сценариев**

**Сценарий** — это комбинация узлов моделирования, настраиваемая пользователем в зависимости от решаемой задачи.

### 3.2.1. Интерфейс конструктора сценариев

Интерфейс **Конструктора сценариев** состоит из следующих элементов:

1. Рабочей области, в которую помещаются узлы моделирования.
2. Левой боковой панели с узлами моделирования.
3. Правой боковой панели с настройками узла (открывается при выборе узла, размещенного в рабочем поле).
4. Верхней панели с кнопками запуска узлов, обновления статусов и метаданных и раскрытием журнала логирования.
5. Панель создания нескольких сценариев в рамках одного проекта

Для удобства пользования рабочей областью в левом нижнем углу предусмотрены кнопки масштабирования (кнопки  и  для приближения и удаления, соответственно, кнопка  для масштабирования на объектах, расположенных в рабочем поле) и кнопка, блокирующая перемещение объектов в рабочем пространстве ().

Левую боковую панель можно скрыть, выбрав кнопку  в нижней части панели. Раскрыть скрытую панель можно выбрав  также в нижней части панели. Правая боковая панель открывается при выборе узла из рабочего поля. Скрыть эту панель можно выбрав иконку  в правом верхнем углу.

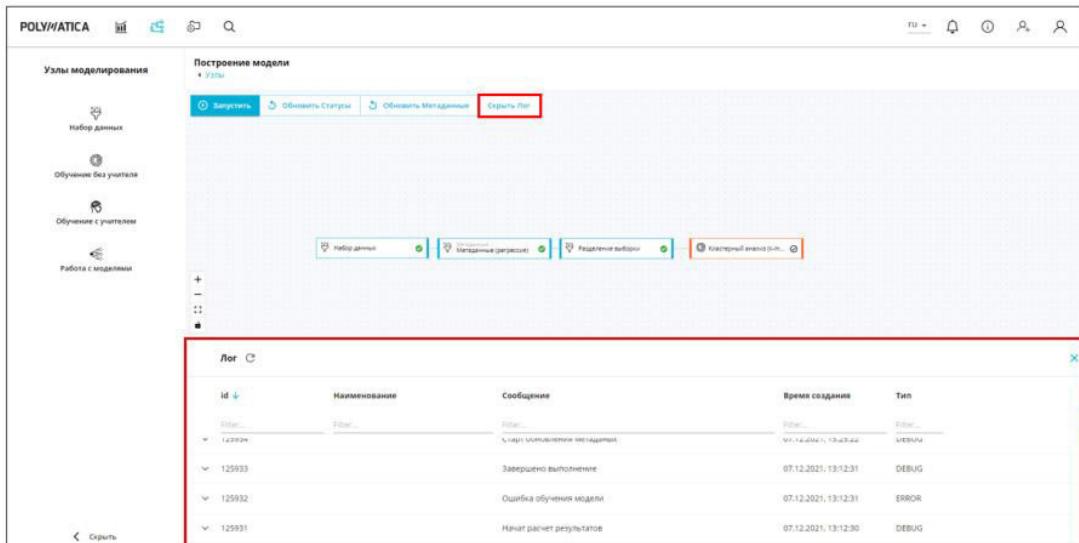
Панель с кнопками включает в себя:

- Кнопку **«Запустить»**, которая запускает процесс расчета сценария.
- Кнопку **«Обновить Статусы»**, которая обновляет статусы узлов.
- Кнопку **«Обновить Метаданные»**, которая обновляет метаданные узлов.
- Кнопку **«Показать Лог»**, которая раскрывает панель с логами сценария.

### 3.2.2. Логирование

Для удобства работы в компоненте предусмотрен просмотр лога событий, связанных с расчетом сценария моделирования.

При выборе кнопки «**Показать Лог**» в нижней части страницы откроется панель.



**Рисунок 62 Журнал логирования**

В этой панели отображена информация о событиях сценария, а именно:

- id события.
- Сообщение — описание события.
- Наименование — наименование узла MD.
- Дата и время события.
- Тип события. В лог записываются события типов INFO, DEBUG и ERROR.

В случае событий типа ERROR в описание также включен traceback. Для отображения необходимо выбрать иконку рядом с событием.

| id     | Наименование | Сообщение                | Время создания       | Тип   |
|--------|--------------|--------------------------|----------------------|-------|
| 125932 | Flow...      | Ошибка обучения модели   | 07.12.2021, 13:12:31 | ERROR |
| 125931 | Flow...      | Начал расчет результатов | 07.12.2021, 13:12:30 | DEBUG |

**Рисунок 63 Пример отображения traceback**

Для удобства поиска записей в логе можно использовать фильтры. Фильтры предусмотрены для каждой из колонок таблицы.

### 3.2.3. Создание сценария

Для создания сценария необходимо:

- Поместить в рабочую область **узел Набор данных**. Он является первоначальным узлом любого сценария, т.к. подгружает набор данных, на основе которого будет строиться модель.
- Выбрать необходимый набор данных, задать остальные параметры узла и запустить расчет, выбрав кнопку «**Запустить**».
- Поместить следующий необходимый узел, создать связь с предыдущим, задать параметры и запустить.
- Последовательно проделать аналогичные шаги со всеми необходимыми узлами (о последовательности узлов подробнее в разделе Пример базового сценария).

### 3.2.4. Создание нескольких сценариев

В рамках моделирования может возникнуть необходимость построить несколько сценариев с разными настройками. Для создания нового сценария в рамках одного проекта необходимо:

- В верхней части конструктора рядом с названием проекта выбрать иконку 
- В открывшемся окне **Создание сценария** задать название и описание нового сценария и сохранить изменения.
- В панели рядом с названием проекта появится название созданного сценария.

Переключение между сценариями производится в панели рядом с названием проекта моделирования.

### 3.2.5. Узлы

Каждый **узел** отвечает за определенный этап машинного обучения.

Для работы с узлом необходимо выбрать его левой кнопкой мыши в панели с узлами и поместить в рабочее поле.

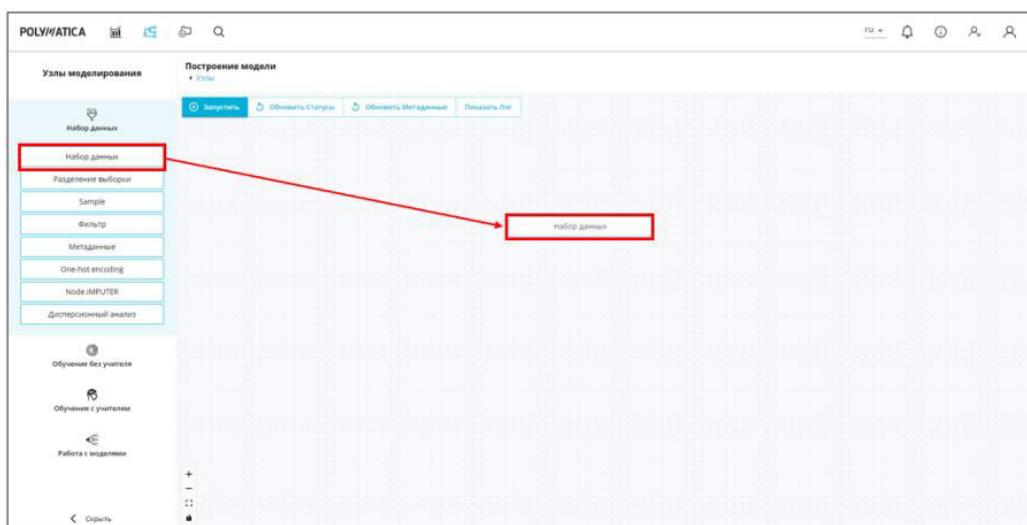
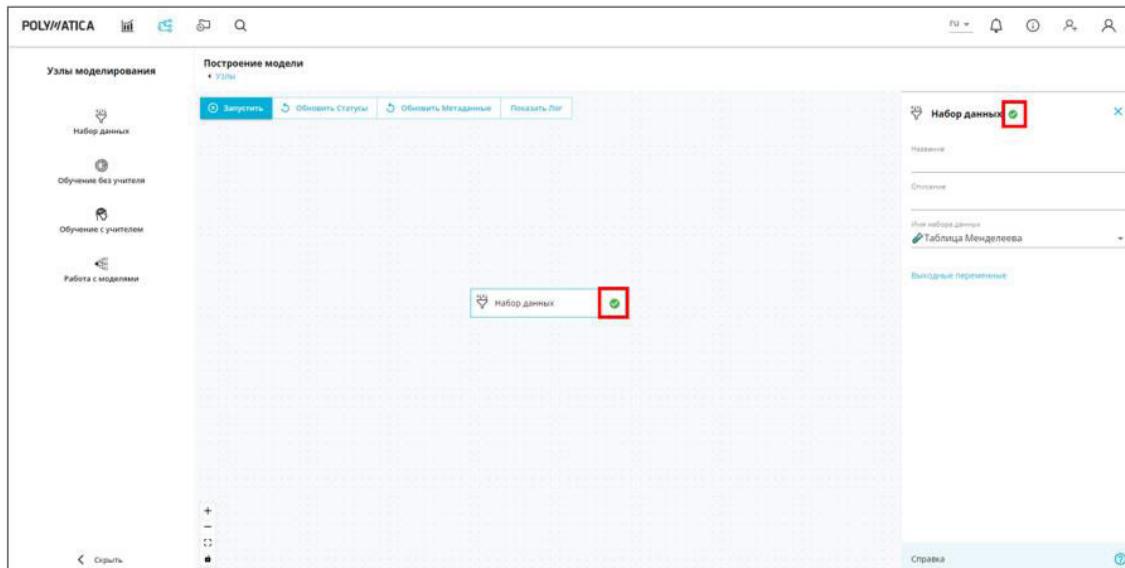


Рисунок 64 Размещение узла в рабочее поле

### 3.2.5.1. Статус Узла

Каждый узел имеет статус. Отображается он с правой стороны узла, а также в верхней части боковой панели с настройками узла.



**Рисунок 65 Пример статуса узла Набор данных**

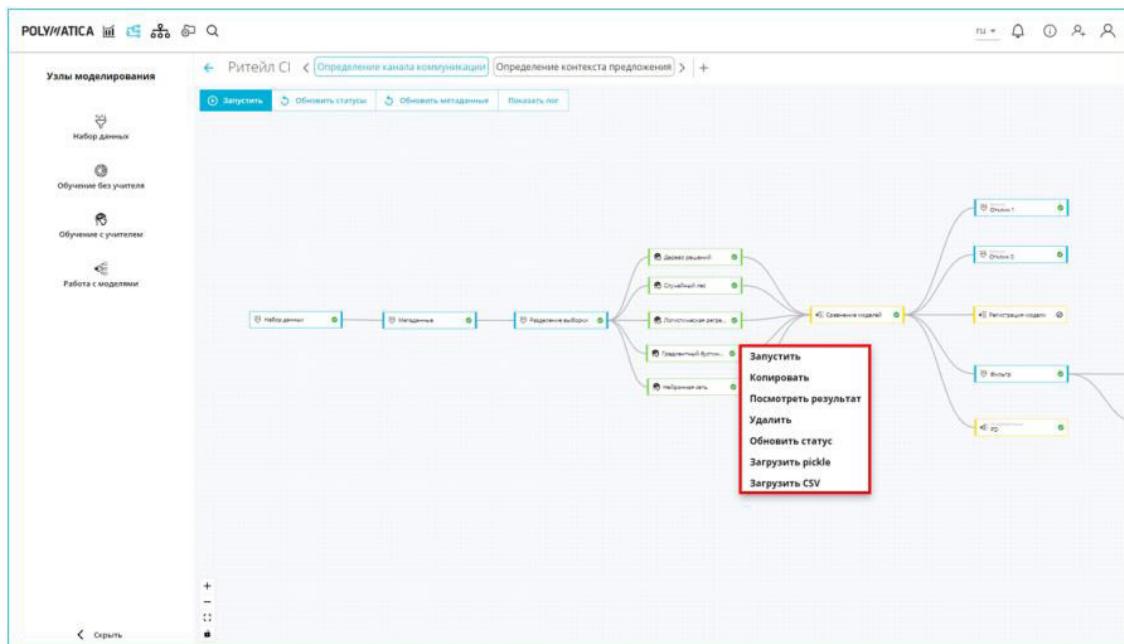
Узлы могут иметь следующие статусы:

- ⓘ — необходимо обратить внимание на узел (при наведении на иконку появится всплывающая подсказка).
- ⚡ — Узел готов к запуску.
- ✖ — Ошибка выполнения узла (понять причину ошибки можно в журнале логирования).
- ✓ — Узел выполнен.
- ⏪ — Узел в процессе выполнения.

### 3.2.5.2. Действия над узлом

После размещения узла в рабочем поле над ним можно совершить действия. Для этого нужно нажать правой кнопкой мыши на узел и в выпадающем меню выбрать необходимое. Узел можно:

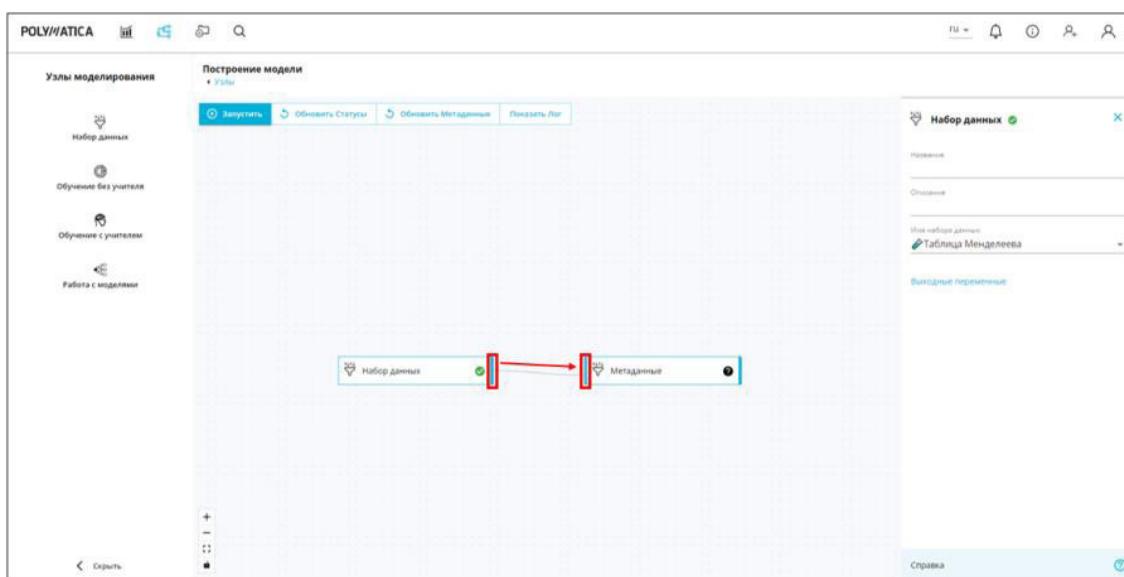
- Запустить.
- Копировать.
- Посмотреть результат выполнения (В результате выбора откроется окно с результатами расчета. Для каждого узла предусмотрены свои результаты).
- Удалить.
- Обновить статус.
- Загрузить pickle (Локальная выгрузка pickle файла).
- Загрузить CSV (Локальная выгрузка результатов узла в виде csv-файла).



**Рисунок 66 Пример выпадающего меню**

### 3.2.5.3. Связь между узлами

Для создания связи между узлами необходимо правой кнопкой мыши выбрать область с правой стороны узла и протянуть появившуюся линию связи до левой стороны следующего узла.



**Рисунок 67 Создание связи между узлами**

Для удаления связи необходимо левой кнопкой мыши выбрать линию связи и нажать на клавиатуре клавишу **Delete**.

В таблице ниже представлен список узлов, между которыми можно создавать связь.

| <b>Исходный узел</b>      | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>  | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|---------------------------|---|--|
| <b>Набор данных</b>       | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов  | –  |
| <b>Разделение выборки</b> | Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM | Набор данных<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off) |
| <b>Sample</b>             | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия  | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия  |

| <b>Исходный узел</b> | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|----------------------|--|--|
|                      | Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM  | Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Интерпретация<br>Подбор отсечки (Cut off)   |
| <b>Фильтр</b>        | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Интерпретация | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off) |
| <b>Метаданные</b>    | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов   | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Кластерный анализ (k-means)  |

| <b>Исходный узел</b>        | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>  | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|-----------------------------|---|--|
|                             | Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM   | Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off)  |
| <b>One-hot encoding</b>     | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Интерпретация |
| <b>Заполнение пропусков</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)   | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация   |

| <b>Исходный узел</b>    | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|-------------------------|--|--|
|                         | Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей   | Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Интерпретация<br>Регистрация модели                           |
| <b>Трансформация</b>    | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Биннинг/энкодинг<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Сравнение моделей | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Биннинг/энкодинг<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Сравнение моделей<br>Подбор отсечки (Cut off) |
| <b>Биннинг/энкодинг</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия   | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Дисперсионный анализ<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия  |

| <b>Исходный узел</b>        | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>  | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>  |
|-----------------------------|---|---|
|                             | Линейные модели<br>Нейронная сеть<br>LDA<br>Сравнение моделей   | Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Сравнение моделей<br>Подбор отсечки (Cut off)  |
| <b>Дисперсионный анализ</b> | Бининг/энкодинг<br>Фильтр   | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Сравнение моделей  |
| <b>Стандартизация</b>       | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off) | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off) |
| <b>Веса классов</b>         | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Кластерный анализ (k-means)<br>Иерархическая  | Набор данных<br>Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Кластерный анализ (k-means)   |

| <b>Исходный узел</b>               | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|------------------------------------|--|--|
|                                    | кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off)   | Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off)   |
| <b>Кластерный анализ (k-means)</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Регистрация модели | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Подбор отсечки (Cut off) |
| <b>Иерархическая кластеризация</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Дисперсионный анализ<br>Стандартизация  | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Стандартизация<br>Веса классов  |

| <b>Исходный узел</b>  | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>  | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|-----------------------|---|--|
|                       | Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM   | Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Подбор отсечки (Cut off)  |
| <b>Дерево решений</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off) | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Биннинг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Подбор отсечки (Cut off) |
| <b>Случайный лес</b>  | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные  | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные   |

| <b>Исходный узел</b>         | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>  |
|------------------------------|--|---|
|                              | One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off) | One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Подбор отсечки (Cut off) |
| <b>Байесовская регрессия</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)   | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)              |

| <b>Исходный узел</b>           | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|--------------------------------|--|--|
|                                | Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off)   | GLM<br>Сравнение моделей<br>Интерпретация<br>Подбор отсечки (Cut off)  |
| <b>Линейная регрессия</b>      | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off) |
| <b>Логистическая регрессия</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес  | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия   |

| <b>Исходный узел</b>   | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|------------------------|--|--|
|                        | Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off)  | Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off)   |
| <b>Линейные модели</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off) | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Трансформация<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off) |
| <b>Нейронная сеть</b>  | Sample<br>Фильтр<br>Метаданные<br>Трансформация<br>Бининг/энкодинг   | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>Трансформация  |

| <b>Исходный узел</b>                 | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>  |
|--------------------------------------|--|---|
|                                      | Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off)  | Бининг/энкодинг<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>Сравнение моделей  |
| <b>LDA</b>                           | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>Трансформация<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off)                   | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>Трансформация<br>Бининг/энкодинг<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>LDA<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация  |
| <b>Градиентный бустинг (XGBoost)</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей |

| <b>Исходный узел</b>                      | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>  | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>  |
|---|---|---|
|   | Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off)  | Интерпретация<br>Подбор отсечки (Cut off)   |
| <b>Градиентный бустинг<br/>(LightGBM)</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off) | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off) |
| <b>Градиентный бустинг<br/>(CatBoost)</b> | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>GLM<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off) | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off)   |

| <b>Исходный узел</b>      | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|---------------------------|--|--|
| <b>GLM</b>                | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>Сравнение моделей<br>Регистрация моделей<br>Интерпретация<br>Подбор отсечки (Cut off) | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Заполнение пропусков<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>Сравнение моделей<br>Подбор отсечки (Cut off) |
| <b>Сравнение моделей</b>  | Разделение выборки<br>Sample<br>Фильтр<br>Метаданные<br>One-hot encoding<br>Трансформация<br>Заполнение пропусков<br>Бининг/энкодинг<br>Дисперсионный анализ<br>Стандартизация<br>Веса классов<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Регистрация моделей<br>Интерпретация   | One-hot encoding<br>Трансформация<br>Заполнение пропусков<br>Бининг/энкодинг<br>Стандартизация<br>Веса классов<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Подбор отсечки (Cut off)  |
| <b>Регистрация модели</b> | -  | Кластерный анализ (k-means)<br>Дерево решений  |

| <b>Исходный узел</b>            | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b>   | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b>   |
|---------------------------------|--|--|
|                                 |  | Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей                                    |
| <b>Интерпретация</b>            | Sample<br>Фильтр<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей<br>Подбор отсечки (Cut off)   | Фильтр<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>LDA<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM<br>Сравнение моделей |
| <b>Подбор отсечки (Cut off)</b> | Разделение выборки<br>Sample<br>Метаданные<br>Фильтр<br>Трансформация<br>Бининг/Энкодинг<br>Стандартизация<br>Веса классов<br>Кластерный анализ (k-means)<br>Иерархическая кластеризация<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Линейная регрессия<br>Логистическая регрессия<br>Линейные модели<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM | Стандартизация<br>Веса классов<br>Дерево решений<br>Случайный лес<br>Байесовская регрессия<br>Логистическая регрессия<br>Линейные модели<br>Нейронная сеть<br>Градиентный бустинг (XGBoost)<br>Градиентный бустинг (LightGBM)<br>Градиентный бустинг (CatBoost)<br>GLM                           |

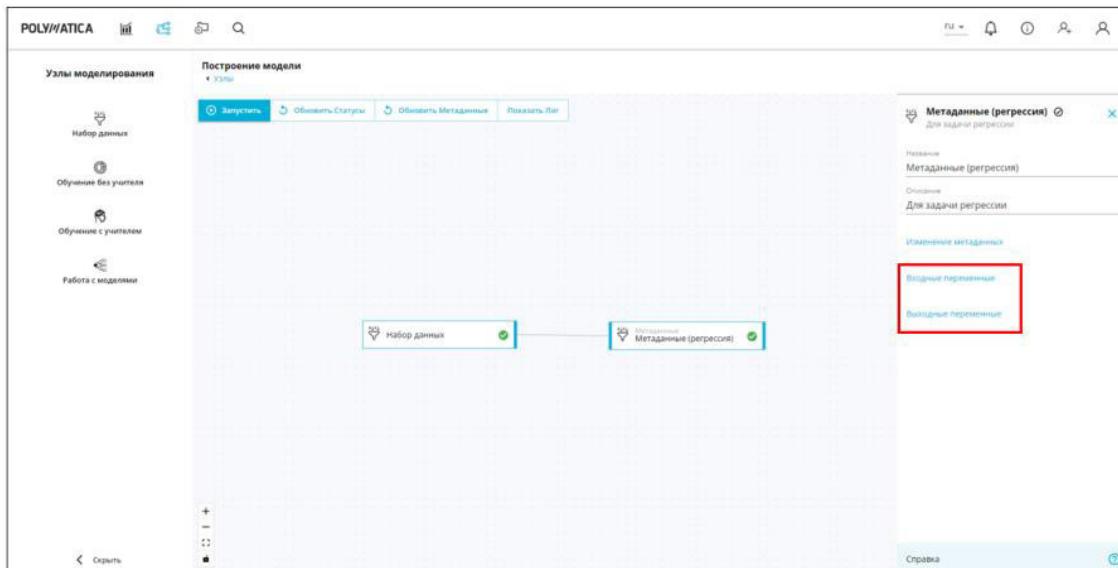
| <b>Исходный узел</b> | <b>Исходящий узел<br/>(в который входит связь из исходного узла)</b> | <b>Входящий узел<br/>(из которого выходит связь в исходный узел)</b> |
|----------------------|--|--|
|                      | (CatBoost)<br>GLM<br>Сравнение моделей                               |  |

**Таблица 6 Список возможной связи узлов**

### 3.2.5.4. Входные и выходные данные узла

Каждый узел получает данные на вход, преобразует их и подает на выход. Входом текущего узла является выход предыдущего. Ознакомиться с входными и выходными параметрами можно в правой панели с настройками узла, щелкнув по ссылкам «**Входные переменные**» и «**Выходные переменные**».

Для узла «**Набор данных**» не предусмотрены **ссылки** (и соответственно окна) «**Входные переменные**» и «**Выходные переменные**». Ознакомиться с входными переменными **узла «Набор данных»** можно в **окне «Конфигурация переменных»** (подробнее в разделе [Узел «Набор данных»](#))



**Рисунок 68 Расположение ссылок для просмотра Входных и Выходных переменных узла**

При выборе ссылки «**Входные переменные**» откроется одноименное окно, в котором будут отражены подающиеся на вход узла переменные и их метаданные. При выборе же ссылки «**Выходные переменные**» откроется аналогичное окно, но с измененными и созданными переменными на выходе узла.

The screenshot shows two panels of a node configuration interface:

- Входные переменные (Input Variables):** A table with columns: Пользовательское и., Роль (Role), Тип (Type), and Целевой класс (Target Class). The 'density' row is highlighted with a red border.
- Выходные переменные (Output Variables):** A table with columns: Пользовательское и., Роль (Role), Тип (Type), and Целевой класс (Target Class). The 'density' row is highlighted with a red border.

**Рисунок 69 Пример Входных и выходных переменных узла «Метаданные»**

К метаданным переменной относятся **Роль**, **Тип** и **Целевой класс**.

В Модуле предусмотрены следующие **Роли** переменной:

- **Target — Целевая переменная** — зависимая переменная, используемая в качестве цели моделирования.
- **Predictor — Предиктор** — независимая переменная, используемая в качестве прогнозирующей целевую переменную.
- **Excluded — Исключен** — данная роль позволяет исключить переменную из процесса моделирования.
- **Weight — Вес** — данная роль задается переменной, которая указывает модели на важность наблюдений. Пример: для прогноза могут быть важны тенденции за последний период (например, сезон).
- **Cluster\_ID — Сегмент** — переменная с данной ролью рассчитывается при решении задачи кластеризации.
- **Prediction — Прогноз** — переменная с данной ролью рассчитывается при решении задачи регрессии.
- **Classification — Прогноз класс** — переменная с данной ролью рассчитывается при решении задачи классификации.
- **Partition — Разделение** — переменная с данной ролью задает разделение набора данных на обучающую, валидационную и тестовую выборки (рассчитывается в узле «Разделение выборки»).
- **ClassWeight — Вес класса** — переменная с данной ролью корректирует несбалансированность классов (задается в узле «Веса классов»).

Первоначально при выборе набора данных в узле «Набор данных» все атрибуты имеют роль Предиктор.

Предусмотрены следующие **Типы** переменной:

- **Номинальный (Nominal)** — дискретные значения без числовой связи между категориями.

- **Порядковый (Ordinal)** — дискретные значения, которые можно ранжировать или сортировать (Пример — хорошо, великолепно, превосходно).
- **Двоичный (Binary)** — это дискретные данные, которые могут относиться только к одной из двух категорий: 0 и 1.
- **Интервальный (Interval)** — числовые значения.

**Целевой класс** используется в процессе моделирования и может быть указан только для Целевой переменной.

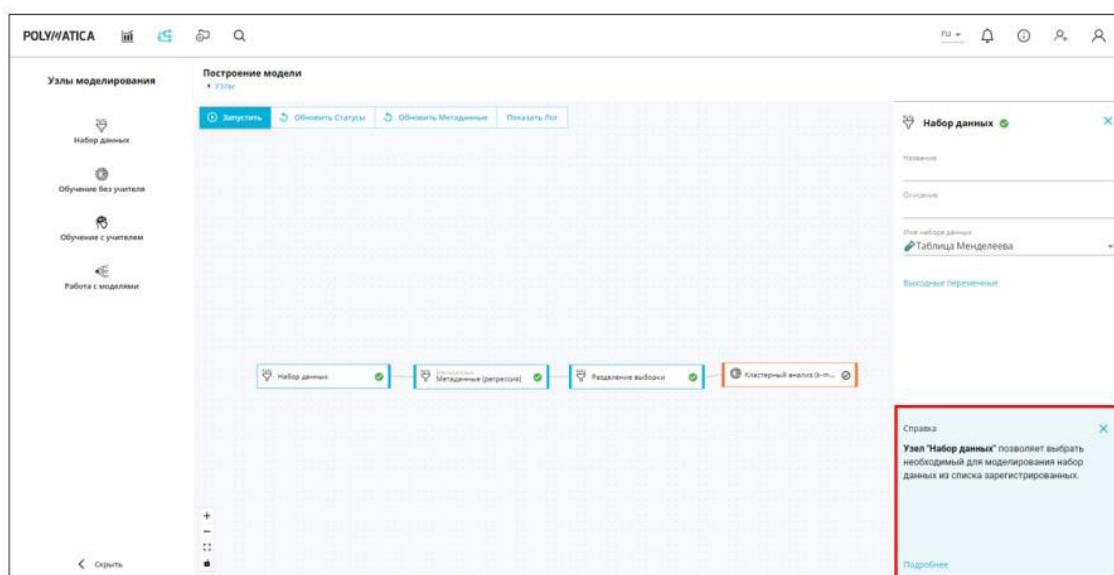
Вручную изменить метаданные атрибутов можно при помощи узла «Метаданные».

### 3.2.5.5. Описание узлов

Каждый **Узел** выполняет отдельную операцию. Для удобства все узлы разделены на группы в зависимости от выполняемых функций:

- Группа «**Набор данных**» включает в себя узлы для подготовки и преобразования данных перед построением моделей.
- Группы «**Обучение с учителем**» и «**Обучение без учителя**» представлены узлами-алгоритмами машинного обучения.
- Узлы группы «**Работа с моделями**» используются на завершающих этапах моделирования для интерпретации, отсечки, сравнения моделей и регистрации в Репозиторий.

Сведения о конкретном узле можно получить в его **Справке**, которая находится в нижней части боковой панели с настройками узла. Для получения подробной информации об узле необходимо щелкнуть ссылку «**Подробнее**». В результате откроется новая вкладка в браузере с описанием узла из настоящего руководства пользователя.



**Рисунок 70 Справка узла «Набор данных»**

В текущей версии Модуля предусмотрены 37 узлов. Ознакомиться с кратким описанием каждого узла можно в таблице ниже.

### Краткое описание узлов

| Группа узлов                | Название узла                             | Краткое описание   |
|-----------------------------|---|--|
| <b>Набор данных</b>         | <b>Узел «Набор данных»</b>                | Данный узел позволяет выбрать необходимый для моделирования набор данных из списка зарегистрированных  |
|                             | <b>Узел «Разделение выборки»</b>          | Данный узел разбивает набор данных на обучающую, валидационную и тестовую выборки  |
|                             | <b>Узел «Sample»</b>                      | Данный узел корректирует неравномерное распределение классов в исходном наборе данных  |
|                             | <b>Узел «Фильтр»</b>                      | Данный узел позволяет по заданным условиям удалить наблюдения из процесса моделирования  |
|                             | <b>Узел «Метаданные»</b>                  | Данный узел позволяет изменить метаданные переменных   |
|                             | <b>Узел «One-hot encoding»</b>            | Данный узел преобразует категориальные переменные в числовые данные  |
|                             | <b>Узел «Заполнение пропусков»</b>        | Данный узел обрабатывает пропущенные значения  |
|                             | <b>Узел «Трансформация»</b>               | Данный узел позволяет рассчитать новые переменные  |
|                             | <b>Узел «Биннинг/энкодинг»</b>            | Данный узел включает в себя методы бинаризации интервальных переменных и кодирования категориальных переменных.  |
|                             | <b>Узел «Дисперсионный анализ»</b>        | Данный узел позволяет исследовать значимость различия между средними значениями зависимой количественной переменной по группам фактора (независимой переменной).                                 |
|                             | <b>Узел «Стандартизация»</b>              | Данный узел приводит признаки в разных единицах измерения и диапазонах значений к единому виду, который позволит сравнивать их между собой или использовать для расчета схожести объектов.       |
|                             | <b>Узел «Веса классов»</b>                | Данный узел корректирует несбалансированность классов (в обучающей выборке доли объектов разных классов существенно различаются)   |
|                             | <b>Узел «Автоэнкодер (PyTorch)»</b>       |  |
|                             | <b>Узел «PCA»</b>                         |  |
|                             | <b>Узел «Профилирование»</b>              |  |
| <b>Обучение без учителя</b> | <b>Узел «Кластерный анализ (k-means)»</b> | Данный узел группирует наблюдения в подмножества (кластеры) таким образом, чтобы наблюдения внутри одного кластера были похожи друг на друга, но различались с наблюдениями из других кластеров. |
|                             | <b>Узел «Иерархическая кластеризация»</b> | Данный узел создает иерархии вложенных подмножеств (кластеров).  |

| <b>Группа узлов</b>        | <b>Название узла</b>                                | <b>Краткое описание</b>  |
|----------------------------|---|--|
| <b>Обучение с учителем</b> | <b>Узел «Детекция аномалий (PyTorch)»</b>           |  |
|                            | <b>Узел «Ассоциативные правила»</b>                 |  |
|                            | <b>Узел «Дерево решений»</b>                        | Данный узел обобщает наблюдения правилами вида «Если..., то...» в иерархическую, последовательную структуру в виде дерева.<br>Используется для решения задач классификации и регрессии   |
|                            | <b>Узел «Случайный лес»</b>                         | В основе данного узла лежит алгоритм, который представляет собой ансамбль деревьев решений.<br>Используется для решения задач классификации и регрессии  |
|                            | <b>Узел «Байесовская регрессия»</b>                 | Данный узел представляет собой линейную регрессию с применением распределения вероятностей параметров, а не точечных оценок<br>Используется для решения задач регрессии  |
|                            | <b>Узел «Линейная регрессия»</b>                    | В результате данного узла строится модель зависимости между входными и выходными переменными с линейной функцией связи<br>Используется для решения задач регрессии   |
|                            | <b>Узел «Логистическая регрессия»</b>               | В основе данного узла лежит метод построения линейного классификатора, позволяющий оценивать апостериорные вероятности принадлежности объектов классам.<br>Используется для решения задач классификации  |
|                            | <b>Узел «Линейные модели»</b>                       | Данный узел объединяет в себе линейные классификаторы и регрессоры   |
|                            | <b>Узел «Нейронная сеть»</b>                        | В основе данного узла лежит упрощенная модель биологической нейронной сети.<br>Используется для решения задач классификации и регрессии  |
|                            | <b>Узел «LDA» (Линейный дискриминантный анализ)</b> | Данный узел применяется для нахождения линейных комбинаций признаков, наилучшим образом разделяющих два или более класса объектов или событий.   |
|                            | <b>Узел «Градиентный бустинг (XGBOOST)»</b>         | В основе данного узла лежит алгоритм градиентного бустинга на деревьях поиска решений.<br>Используется для решения задач классификации и регрессии.  |
|                            | <b>Узел «Градиентный бустинг (XGBOOST)»</b>         | В основе узла лежит реализация алгоритма градиентного бустинга на деревьях поиска решений, который включает в себя две ключевые идеи: Градиентная односторонняя выборка (GOSS) и Объединение взаимоисключающих признаков (EFB).<br>Используется для решения задач классификации и регрессии. |

| Группа узлов      | Название узла                               | Краткое описание   |
|-------------------|---|--|
| Работа с моделями | <b>Узел «Градиентный бустинг (XGBOOST)»</b> | В основе узла лежит реализация алгоритма градиентного бустинга, которая оптимизирована под работу с категориальными признаками и хорошо работает с параметрами по умолчанию. Используется для решения задач классификации и регрессии. |
|                   | <b>Узел «GLM»</b>                           | Данный узел обобщает линейную регрессию и допускает наличие у зависимой переменной распределения, отличающегося от нормального. GLM связывает зависимую переменную с факторами посредством задаваемой функции связи.                   |
|                   | <b>Узел «Нейронная сеть (PyTorch)»</b>      |  |
|                   | <b>Узел «AutoML»</b>                        |  |
| Работа с моделями | <b>Узел «Сравнение моделей»</b>             | Данный узел оценивает построенные модели и выбирает лучшую.  |
|                   | <b>Узел «Регистрация модели»</b>            | Данный узел сохраняет построенную модель в выбранном проекте репозитория Model Manager.  |
|                   | <b>Узел «Интерпретация»</b>                 | Данный узел, включает в себя методы, которые позволяют объяснить принципы и закономерности, которые использует модель в ходе прогнозирования.  |
|                   | <b>Узел «Подбор отсечки (Cut off)»</b>      | Данный узел позволяет определить оптимальный порог отсечения для высокого соотношения истинных и ложных срабатываний модели  |

### 3.2.5.6. Группа узлов «Набор данных»

#### 3.2.5.6.1. Узел «Набор данных»

**Узел «Набор данных»** является начальным узлом любого сценария моделирования, позволяет выбрать необходимый для моделирования набор данных из списка зарегистрированных и задать конфигурацию переменных.

При конфигурировании переменных Пользователь может указать необходима ли переменная далее в процессе моделирования, а также указать необходимо ли регистрировать переменную при регистрации модели.

Помимо этого, Пользователь может отфильтровать, а также случайным образом отобрать наблюдения, которые далее будут использоваться в сценарии моделирования.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения           | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе |

| Параметр   | Возможные значения и ограничения                             | Описание   |
|--|--|--|
| <b>Описание</b>                                  | Ручной ввод<br>Ограничений на значение нет                   | Описание узла  |
| <b>Имя набора данных</b>                         | Раскрывающийся список с названиями наборов данных            | Список доступных пользователю наборов данных   |
| <b>Конфигурация переменных</b>                   | Кнопка   | При выборе кнопки откроется окно <b>Конфигурация переменных</b> (подробное описание ниже).                     |
| <b>Фильтр</b>                                    | Поле ввода условия фильтрации                                | При выборе поля ввода откроется окно <b>Формула</b> (подробное описание ниже).                                 |
| <b>Доля наблюдений для выборки</b>               | Ручной ввод<br>Значение больше 0 и меньше или равно 1        | Данный параметр задает долю наблюдений исходного набора данных, которые будут включены в итоговый набор данных |
| <b>Максимальное количество загружаемых строк</b> | Ручной ввод целочисленного значения<br>По умолчанию – 100000 | Данный параметр ограничивает количество строк, которые будут использоваться в ходе моделирования               |

Таблица 7 Параметры узла "Набор данных"

### Окно Конфигурация переменных

В окне **Конфигурация переменных** Пользователь может ознакомиться с входными переменными узла и указать на необходимость регистрировать переменную при сохранении модели в репозиторий (столбец **Регистрация**) и необходимость данной переменной при построении дальнейшего сценария (столбец **Дроп**).

При выборе **Регистрация** (значение Да) данная переменная будет зарегистрирована при сохранении модели в репозиторий ММ.

При выборе **Дроп** (значение Да) далее в сценарии данная переменная использоваться не будет и вернуть ее в процесс будет невозможно. Для временного исключения атрибута необходимо назначить ему роль Исключен (подробнее узел «**Метаданные**»).

**ВАЖНО:** Целевой переменной на данном шаге необходимо задать значение Нет в столбце **Регистрация**. Это необходимо, чтобы далее данная переменная не фигурировала при регистрации и при публикации в ММ.

**Роль** и **Тип переменной** можно изменить в узле "**Метаданные**". **Тип источника** задается автоматически: если переменная из исходного набора данных, то указано значение **Base**, если рассчитана в системе - значение **Computed**.

| Конфигурация переменных       |           |          |               |             |      |
|-------------------------------|-----------|----------|---------------|-------------|------|
| Пользовательское имя атрибута | Роли      | Тип      | Тип источника | Регистрация | Дроп |
| id                            | Предиктор | Interval | Base          | да          | нет  |
| filial                        | Предиктор | Nominal  | Base          | да          | нет  |
| округ                         | Предиктор | Nominal  | Base          | да          | нет  |
| YEAR                          | Предиктор | Interval | Base          | да          | нет  |
| model_id                      | Предиктор | Interval | Base          | да          | нет  |
| TYPE                          | Предиктор | Nominal  | Base          | да          | нет  |
| vendor                        | Предиктор | Nominal  | Base          | да          | нет  |
| model                         | Предиктор | Nominal  | Base          | да          | нет  |
| quarter                       | Предиктор | Interval | Base          | да          | нет  |
| purchase_type                 | Предиктор | Nominal  | Base          | да          | нет  |
| supplier_id                   | Предиктор | Interval | Base          | да          | нет  |
| supplier_name                 | Предиктор | Nominal  | Base          | да          | нет  |
| purchaser_id                  | Предиктор | Interval | Base          | да          | нет  |

**Рисунок 71 Окно Конфигурация переменных**

### Окно Формула

В окне **Формула** задается условие отбора наблюдений из исходного набора данных.

Основными элементами окна Формула являются:

Формула

1 Σ ('Год'>"1996") and true  
2        
3 Переменные >  
 Операторы >  
 Функции >

#Год > 1996 AND TRUE

Отменить      Сохранить

**Рисунок 72 Окно Формула**

1. Стока ввода текстового представления условия.
2. Рабочее поле, в котором строится графическое представление условия.
3. Вкладка, которая открывается при щелчке правой кнопкой мыши в строке ввода, и включает в себя:
  - Переменные из набора данных.
  - Операторы и функции, представленные в таблице ниже.

Позволяет просматривать переменные набора данных и строить конструкции условий.

Для **задания условия** необходимо:

- В строку текстового представления ввести условие, используя операторы, функции и переменные набора данных.
  - Название переменной должно быть указано в обратных одинарных кавычках (` `). Пример: `Год` > 1996.
  - В качестве десятичного разделителя используется точка.
  - Стока должна указываться в двойных кавычках. Пример: `Пол` = "мужской".
- Посмотреть список операторов, функций и переменных, а также добавить их в условие можно и в панели, которая открывается при щелчке правой кнопкой мыши по строке ввода.

| Оператор/Функция | Описание                  | Текстовое представление (формула)  |
|------------------|---------------------------|--|
| <b>And</b>       | Логическая операция «И»   | true <b>and</b> true<br>Вместо true нужно вставить условие<br><b>Пример:</b> `Год` > 1996 <b>and</b> `Пол` = "Женский"                       |
| <b>Or</b>        | Логическая операция «ИЛИ» | true <b>or</b> true<br>Вместо true нужно вставить условие<br><b>Пример:</b> `Год` > 1996 <b>or</b> `Пол` = "Женский"                         |
| <b>+</b>         | Операция сложения         | 0 <b>+</b> 0<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` <b>+</b> `Цена продукта M` |
| <b>-</b>         | Операция вычитания        | 0 <b>-</b> 0<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` <b>-</b> 1000              |
| <b>*</b>         | Операция умножения        | 0 <b>*</b> 0<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` <b>*</b> `Скидка`          |
| <b>/</b>         | Операция деления          | 0 <b>/</b> 0<br>Вместо 0 нужно вставить числовое значение/числовую переменную  |

| <b>Оператор/Функция</b> | <b>Описание</b>   | <b>Текстовое представление (формула)</b>   |
|-------------------------|---|--|
| <                       | Меньше  | "" < ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Выручка за первый квартал` / `Выручка за год` |
| ≤                       | Меньше или равно  | "" <= ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Возраст` <= 32                               |
| >                       | Больше  | "" > ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Возраст` > "32"                               |
| ≥                       | Больше или равно  | "" >= ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Температура` >= 120                          |
| =                       | Равно   | "" = ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Температура` = 120                            |
| ≠                       | Не равно  | "" != ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Температура` != 120                          |
| <b>like</b>             | Проверяет, удовлетворяет ли символьная строка заданному образцу, который может содержать поисковые символы. Учитывает регистр | "" like ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` like "Москва"                       |
| <b>ilike</b>            | Проверяет, удовлетворяет ли символьная строка заданному образцу, который может содержать поисковые символы.                   | "" ilike ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` ilike "Москва"                     |
| <b>startswith</b>       | Проверяет, есть ли в начале одной текстовой строки другая текстовая строка  | "" startswith ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` startswith "Сан"              |
| <b>endswith</b>         | Проверяет, есть ли в конце одной текстовой строки другая текстовая строка   | "" endswith ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` endswith "бург"                 |
| <b>contains</b>         | Проверяет, встречается ли указанная строка внутри другой строки   | "" contains ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Название` contains "consulting"        |

| Оператор/Функция         | Описание   | Текстовое представление (формула)   |
|--------------------------|--|---|
| <b>between</b>           | Проверяет, входит ли значение в заданный диапазон  | 0 <b>between</b> (0,0)<br>Вместо 0 нужно вставить числовую переменную/значение<br><b>Пример:</b> `ID` between (1, 100500)                                   |
| <b>in</b>                | Проверяет наличие элемента в последовательности  | "" <b>in</b> ("")<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Регион` in ("Москва", "Московская область")                            |
| <b>not in</b>            | Проверяет отсутствие элемента в последовательности   | "" <b>not_in</b> ("")<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Регион` not_in ("Москва", "Московская область")                    |
| <b>coalesce</b>          | Возвращает данные из первого столбца, содержащего значение, отличное от NULL                         | <b>coalesce</b> ("")<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> coalesce(`ProductNumber`, `ProductName`)                             |
| <b>not</b>               | Задает противоположное условие   | <b>not</b> (true)<br>Вместо true нужно вставить условие<br><b>Пример:</b> not(`Год` > 1996)   |
| <b>create_date</b>       | Создает переменную даты из последовательно введенных года, месяца и дня                              | <b>create_date</b> (0,0,0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> create_date(2016,10,25)                       |
| <b>current_date</b>      | Возвращает текущую дату  | <b>current_date</b> ()<br>Вводится без дополнительных параметров  |
| <b>current_timestamp</b> | Возвращает текущие дату и время  | <b>current_timestamp</b> ()<br>Вводится без дополнительных параметров   |
| <b>now</b>               | Возвращает текущие дату и время  | <b>now</b> ()<br>Вводится без дополнительных параметров   |
| <b>create_datetime</b>   | Создает переменную даты времени из последовательно введенных года, месяца, дня, часа, минут и секунд | <b>create_datetime</b> (0,0,0,0,0,0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> create_datetime(2016,10,25,12,18,0) |
| <b>char_length</b>       | Возвращает длину строки  | <b>char_length</b> ("")<br>Вместо "" нужно вставить строку/строковую переменную<br><b>Пример:</b> char_length(`Код_продукта`)                               |
| <b>random</b>            | Возвращает случайное число   | <b>random</b> ()<br>Вводится без дополнительных параметров  |
| <b>In</b>                | Натуральный логарифм   | <b>In</b> (0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> In(`Числовая переменная`)                                  |
| <b>exp</b>               | Экспонента   | <b>exp</b> (0)  |

| Оператор/Функция             | Описание   | Текстовое представление (формула)  |
|------------------------------|--|--|
|                              |  | Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> exp(`Числовая переменная`)   |
| <b>power</b>                 | Возведение в степень   | <code>power(0,0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> power(`Числовая переменная`,2)  |
| <b>sqrt</b>                  | Квадратный корень  | <code>sqrt(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> sqrt(`Числовая переменная`)  |
| <b>abs</b>                   | Абсолютное значение  | <code>abs(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> abs(`Числовая переменная`)  |
| <b>ceil</b>                  | Возвращает наименьшее целое число, которое больше или равно текущему значению                | <code>ceil(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> ceil(25.1)<br>Вернет значение 26   |
| <b>floor</b>                 | Возвращает наибольшее целое число, которое меньше или равно текущему значению                | <code>floor(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> floor(25.1)<br>Вернет значение 25   |
| <b>extract_from_datetime</b> | Получает указанную часть (день, месяц, год, час, минута, секунда) из значения даты и времени | <code>extract_from_datetime("", "")</code><br>Вместо первых кавычек нужно указать необходимую часть для извлечения:<br><ul style="list-style-type: none"> <li>· SECOND</li> <li>· MINUTE</li> <li>· HOUR</li> <li>· DAY</li> <li>· MONTH</li> <li>· YEAR</li> </ul><br>Вместо вторых кавычек нужно указать переменную типа datetime<br><b>Пример:</b><br><code>extract_from_datetime("DAY", `datetime`)</code> |
| <b>extract_from_date</b>     | Получает указанную часть (день, месяц, год) из значения даты                                 | <code>extract_from_date("", "")</code><br>В первые кавычки нужно вписать необходимую часть для извлечения:<br><ul style="list-style-type: none"> <li>· DAY</li> <li>· MONTH</li> <li>· YEAR</li> </ul><br>Вместо вторых кавычек нужно указать переменную типа date<br><b>Пример:</b><br><code>extract_from_date("DAY", `date`)</code>  |

| Оператор/Функция | Описание                                      | Текстовое представление (формула)   |
|------------------|---|---|
| <b>concat</b>    | Объединяет в единую строку указанные значения | <code>concat("")</code><br>Вместо "" нужно вставить строки/строковые переменные<br><b>Пример:</b><br><code>concat(`Фамилия` , " ", `Имя` )</code> |

**Таблица 8 Операторы и функции**

**Результаты выполнения узла:**

- **Таблица с примером данных.** Отображаются первые 100 наблюдений.

| id   | Филиал    | Округ       | Номер отчётного периода | model_id | Тип устр |
|------|-----------|-------------|-------------------------|----------|----------|
| 246  | Филиал 3  | Южный       | 2020                    | 2        | МФУ      |
| 3194 | Филиал 5  | Южный       | 2020                    | 2        | МФУ      |
| 11   | Филиал 6  | Приволжский | 2020                    | 1        | Струйный |
| 12   | Филиал 12 | Сибирский   | 2020                    | 2        | МФУ      |
| 13   | Филиал 6  | Приволжский | 2020                    | 5        | МФУ      |
| 14   | Филиал 4  | Южный       | 2020                    | 3        | МФУ      |
| 15   | Филиал 5  | Южный       | 2020                    | 3        | МФУ      |
| 16   | Филиал 12 | Сибирский   | 2020                    | 3        | МФУ      |
| 17   | Филиал 8  | Приволжский | 2018                    | 4        | Лазерный |

**Рисунок 73 Таблица с примером данных**

- **Таблица с результатами сэмплирования.**

| Данные   | Количество наблюдений | % относительно входных данных |
|----------|-----------------------|-------------------------------|
| Исходные | 11162                 | 100                           |
| Выборка  | 11162                 | 100                           |

**Рисунок 74 Таблица с результатами сэмплирования**

### 3.2.5.6.2. Узел «Разделение выборки»

**Узел «Разделение выборки»** разбивает набор данных на части: **обучающую** (используемую в процессе обучения модели), **валидационную** (используемую для подбора оптимального набора гиперпараметров модели) и **тестовую** согласно заданным Пользователем пропорциям.



**Рисунок 75 Принцип работы узла «Разделение выборки»**

Разбиение можно произвести двумя способами:

- **Простая случайная выборка** — все наблюдения имеют одинаковую вероятность быть отобранными.
- **Выборка со стратификацией** — случайный отбор наблюдений производится в пределах каждого класса (т.е. с помощью стратификации можно избежать "непредставительной" выборки (когда в выборку попадают наблюдения только одной страты/класса) что не гарантируется простой случайной выборкой)

**Список параметров узла** представлен в таблице ниже.

| Параметр                           | Возможные значения и ограничения  | Описание  |
|------------------------------------|---|---|
| <b>Название</b>                    | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>                    | Ручной ввод<br>Ограничений на значение нет  | Описание узла   |
| <b>Разделение выборки на части</b> | Ручной ввод доли (в %) для каждой части<br>Сумма долей должна быть равна 100%   | Доли обучающей, валидационной и тестовой выборок в исходном наборе данных   |
| <b>Метод разбиения</b>             | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"><li>• Simple Random (по умолчанию)</li><li>• Stratified</li></ul> | Метод разбиения исходного набора данных.<br>Предусмотрены: <ul style="list-style-type: none"><li>• <b>Simple Random</b> — все наблюдения имеют одинаковый шанс быть отобранными</li><li>• <b>Stratified</b> — случайный отбор наблюдений выполняется в пределах каждого класса (при выборе данного метода появится поле «<b>Список входных переменных</b>» для указания переменной, по которой будет проводиться стратификация)</li></ul> |

| Параметр    | Возможные значения и ограничения                       | Описание   |
|-------------|--|--|
| <b>Seed</b> | Ручной ввод числового значения<br>По умолчанию — 12345 | Начальное числовое значение для генератора случайных чисел<br>Используется для воспроизведения результатов при повторном запуске |

**Таблица 9 Параметры узла «Разделение выборки»**

#### Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

The screenshot shows a data preview interface with the following columns: Поставщик (Supplier), purchaser\_id, Закупщик (Buyer), Цена за единицу (Unit Price), Количество (Quantity), Итоговая цена (Total Price), Год (Year), and \_partid\_0. The data includes records from PAO "ГОСПРИНТ" and AO "РОСОБЛПЕЧАТЬ".

| Поставщик                 | purchaser_id | Закупщик      | Цена за единицу | Количество | Итоговая цена | Год  | _partid_0 |
|---------------------------|--------------|---------------|-----------------|------------|---------------|------|-----------|
| PAO "ГОСПРИНТ"            | 11           | Иванков Л.Д.  | 11420           | 9          | 102780        | 2020 | 2         |
| PAO "ГОСПРИНТ"            | 15           | Рыбаков А.А.  | 11360           | 8          | 90880         | 2020 | 1         |
| АО "РОСОБЛПЕЧАТЬ"         | 17           | Воронов Р.О.  | 2960            | 8          | 23680         | 2020 | 1         |
| АО "ЦВЕТНАЯ ПЕЧАТЬ ГРУПП" | 28           | Селезнев Л.А. | 20810           | 8          | 166480        | 2020 | 1         |

**Рисунок 76 Таблица с примером данных**

В результате выполнения узла будет рассчитана новая переменная, по которой далее будет производиться разделение набора данных на выборки (переменная **\_partid\_0**).

- Таблица с указанием долей и количества наблюдений, попавших в соответствующую выборку.

The screenshot shows a sampling table with three rows: Валидационная (Validation), Тестовая (Test), and Обучающая (Training). Each row contains the number of observations (27, 30, 36) and the percentage (30%, 40%) of the total dataset.

| Разделение выборки    |    |               |
|-----------------------|----|---------------|
| Количество наблюдений | %  | Выборка       |
| 27                    | 30 | Валидационная |
| 27                    | 30 | Тестовая      |
| 36                    | 40 | Обучающая     |

**Рисунок 77 Пример таблицы с указанием выборок и количества наблюдений**

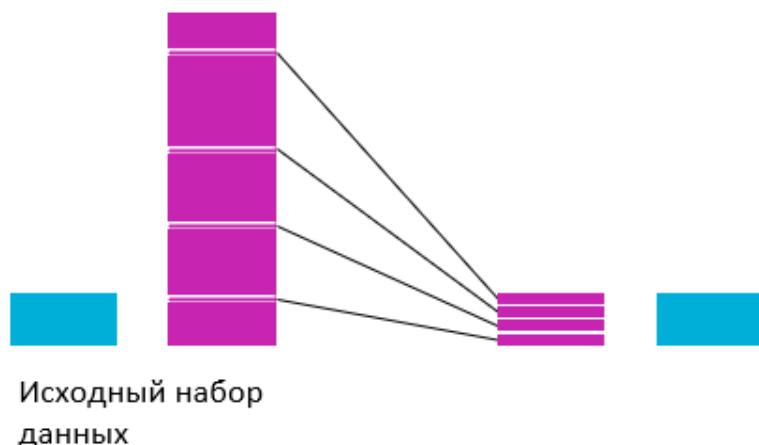
### 3.2.5.6.3. Узел «Sample»

**Узел «Sample»** позволяет сформировать репрезентативную выборку из исходного набора данных, а также скорректировать неравномерное распределение классов. Предусмотрены два метода построения выборки:

- **Простая случайная выборка** — все наблюдения имеют одинаковую вероятность быть отобранными.
- **Выборка со стратификацией** — случайный отбор наблюдений производится в пределах каждого класса (т.е. с помощью стратификации можно избежать "непредставительной" выборки (когда в выборку попадают наблюдения только одной страты/класса) что не гарантируется простой случайной выборкой).

В **задаче классификации** данные называются **несбалансированными**, когда в обучающей выборке доли объектов разных классов существенно различаются.

Для корректировки неравномерного распределения классов предусмотрена **субдискретизация (undersampling)**, которая заменяет больший класс подвыборкой по мощности равной малому классу.



**Рисунок 78** Принцип работы метода Undersampling

При выборе чекбокса **Использовать undersampling** параметр **Доля наблюдений в выборке** игнорируется.

Если выбраны **Выборка со стратификацией** и чекбокс **Использовать undersampling**, то стратификация делается отдельно для выборки с целевым классом и отдельно для выборки с прочими классами. При этом может получиться в целом нестратифицированная выборка, т.к. распределения переменных в выборке по целевому классу может отличаться от распределений во всей выборке.

Целевой класс задается **в узле «Метаданные»**.

Если выбрать чекбокс **Использовать Undersampling** и **Долю целевого класса в выборке** задать меньше доли целевого класса во входных данных, то возникнет ошибка. Это связано с тем, что возможна ситуация, когда в выборке должно быть больше наблюдений с прочими классами, чем их есть во всем исходном наборе.

**Список параметров узла** представлен в таблице ниже.

| Параметр  | Возможные значения и ограничения   | Описание  |
|---|--|---|
| <b>Название</b>   | Ручной ввод<br>Ограничений на значения нет   | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>   | Ручной ввод<br>Ограничений на значения нет   | Описание узла   |
| <b>Метод построения выборки</b>                         | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Простая случайная выборка</li> <li>• Выборка со стратификацией</li> </ul> | Метод разбиения исходного набора данных. Предусмотрены: <ul style="list-style-type: none"> <li>• <b>Простая случайная выборка</b> — все наблюдения имеют одинаковую вероятность быть отобранными</li> <li>• <b>Выборка со стратификацией</b> — случайный отбор наблюдений производится в пределах каждого класса</li> </ul> |
| <b>Доля наблюдений в выборке</b>                        | Ручной ввод числового значения<br>Значение не должно быть больше 1 и меньше или равно 0<br>По умолчанию — 0,5  | Данный параметр задает долю наблюдений, которые попадут в выборку из входных данных   |
| <b>Seed</b>   | Ручной ввод числового значения<br>По умолчанию — 42  | Начальное числовое значение для генератора случайных чисел.<br>Используется для воспроизведения результатов при повторном запуске   |
| <b>Переменные для стратификации</b>                     | Поле выбора со списком переменных  | В случае выбора Выборки со стратификацией нужно указать переменную, по которой будет проводиться стратификация  |
| <b>Использовать Undersampling</b>                       | Чекбокс  | Выбор данного чекбокса указывает на стратегию сэмплирования — <b>субдискретизация</b> , которая удаляет наблюдения из большего класса.  |
| <b>Доля целевого класса относительно входных данных</b> | Ручной ввод числового значения<br>Значение не должно быть больше 1 и меньше или равно 0<br>По умолчанию — 1  | Данный параметр задает долю целевого класса относительно входных данных   |
| <b>Доля целевого класса в выборке</b>                   | Ручной ввод числового значения<br>Значение не должно быть больше 1 и меньше или равно 0<br>По умолчанию — 0,2  | Данный параметр задает долю целевого класса в выборке   |

**Таблица 10 Параметры узла «Sample»**

## Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

Пример данных [Скачать](#)

| Поставщик                 | purchaser_id | Закупщик      | Цена за единицу | Количество | Итоговая цена | Год  | _partid_0 |
|---------------------------|--------------|---------------|-----------------|------------|---------------|------|-----------|
| ПАО "ГОСПРИНТ"            | 11           | Иванков Л.Д.  | 11420           | 9          | 102780        | 2020 | 2         |
| ПАО "ГОСПРИНТ"            | 15           | Рыбаков А.А.  | 11360           | 8          | 90880         | 2020 | 1         |
| АО "РОСОБЛПЕЧАТЬ"         | 17           | Воронов Р.О.  | 2960            | 8          | 23680         | 2020 | 1         |
| АО "ЦВЕТНАЯ ПЕЧАТЬ ГРУПП" | 28           | Селезнев Л.А. | 20610           | 8          | 16480         | 2020 | 1         |

- Таблица с результатами сэмплирования.

Выборка

| Данные ↑  | Количество ↑ | Доля ↑    |
|-----------|--------------|-----------|
| Filter... | Filter...    | Filter... |
| Входные   | 90           | 100       |
| Выборка   | 45           | 50        |

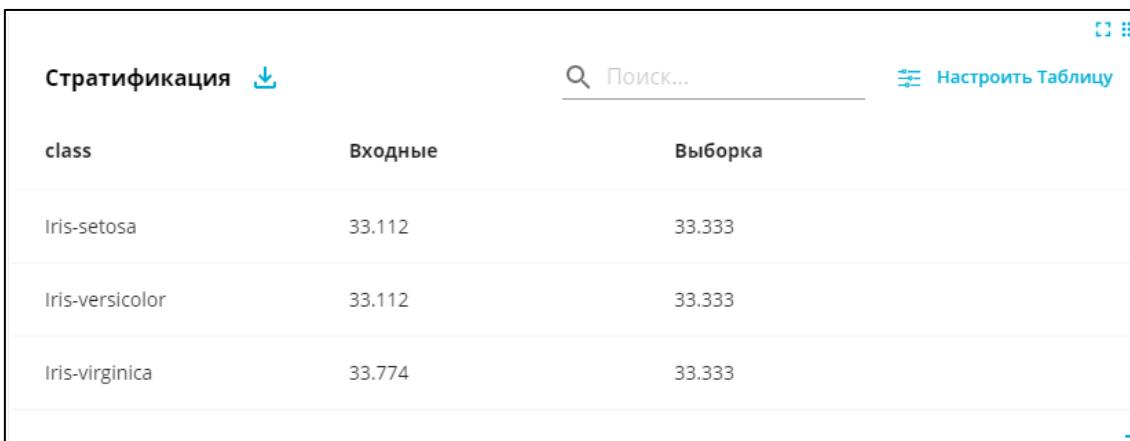
Рисунок 79 Пример таблицы с результатами сэмплирования

Выборка [Скачать](#)

| Классы                      | Наблюдений (входные) | Доля (входные) | Наблюдений (выборка) | Доля (выборка) | % от входных |
|-----------------------------|----------------------|----------------|----------------------|----------------|--------------|
| Целевой класс = Iris-setosa | 50                   | 33.112         | 30                   | 71.428         | 60           |
| Прочие классы               | 101                  | 66.887         | 12                   | 28.571         | 11.881       |
| Все классы                  | 151                  | 100            | 42                   | 100            | 27.814       |

Рисунок 80 Пример таблицы с результатами сэмплирования (с undersampling)

- Таблица с результатами стратификации (отображается при выборе метода построения выборки **Выборка со стратификацией**).



Стратификация

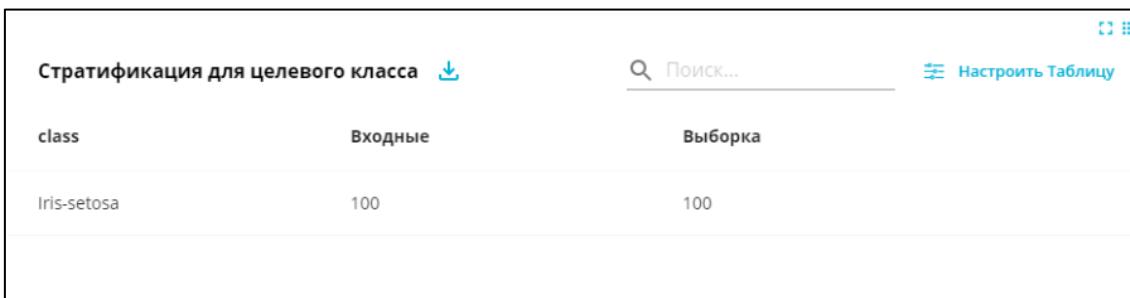
Поиск...

Настройте Таблицу

| class           | Входные | Выборка |
|-----------------|---------|---------|
| Iris-setosa     | 33.112  | 33.333  |
| Iris-versicolor | 33.112  | 33.333  |
| Iris-virginica  | 33.774  | 33.333  |

**Рисунок 81 Пример таблицы с результатами стратификации**

- Таблица с результатами стратификации для целевого класса (отображается при выборе чекбокса **Использовать Undersampling и Выборке со стратификацией**).



Стратификация для целевого класса

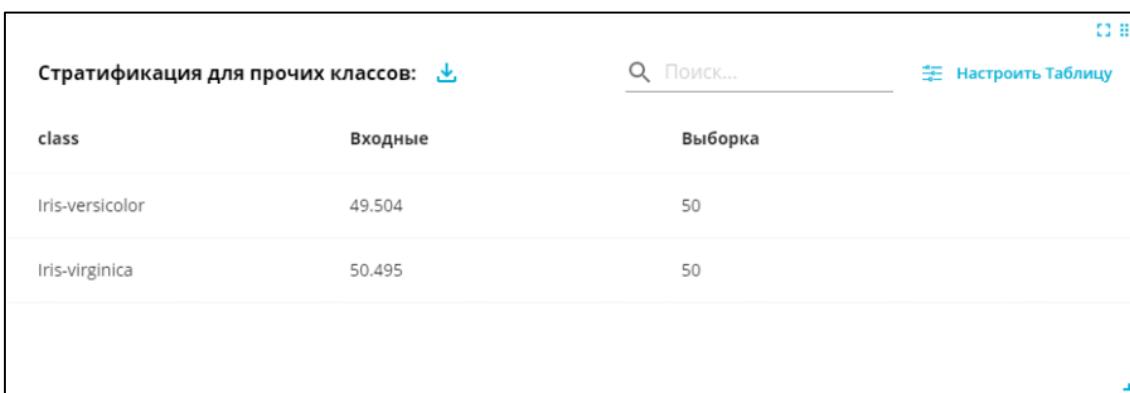
Поиск...

Настройте Таблицу

| class       | Входные | Выборка |
|-------------|---------|---------|
| Iris-setosa | 100     | 100     |

**Рисунок 82 Пример таблицы с результатами стратификации для целевого класса**

- Таблица с результатами стратификации для прочих классов (отображается при выборе чекбокса **Использовать Undersampling и Выборке со стратификацией**).



Стратификация для прочих классов:

Поиск...

Настройте Таблицу

| class           | Входные | Выборка |
|-----------------|---------|---------|
| Iris-versicolor | 49.504  | 50      |
| Iris-virginica  | 50.495  | 50      |

**Рисунок 83 Пример таблицы с результатами стратификации для прочих классов**

### 3.2.5.6.4. Узел «Фильтр»

**Узел «Фильтр»** позволяет по заданным условиям удалить часть наблюдений из процесса моделирования. Отфильтрованные (не соответствующие критериям фильтрации) данные не попадут на вход последующих узлов.

**Список параметров узла** представлен в таблице ниже.

| Параметр                      | Возможные значения и ограничения           | Описание   |
|-------------------------------|--|--|
| <b>Название</b>               | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе                         |
| <b>Описание</b>               | Ручной ввод<br>Ограничений на значение нет | Описание узла  |
| <b>Поле для ввода формулы</b> | Поле ввода условия фильтрации              | При выборе поля ввода откроется окно <b>Формула</b> (подробное описание ниже). |

Таблица 11 Параметры узла «Фильтр»

#### Окно Формула

В окне **Формула** задается условие отбора наблюдений из исходного набора данных.

Основными элементами окна Формула являются:

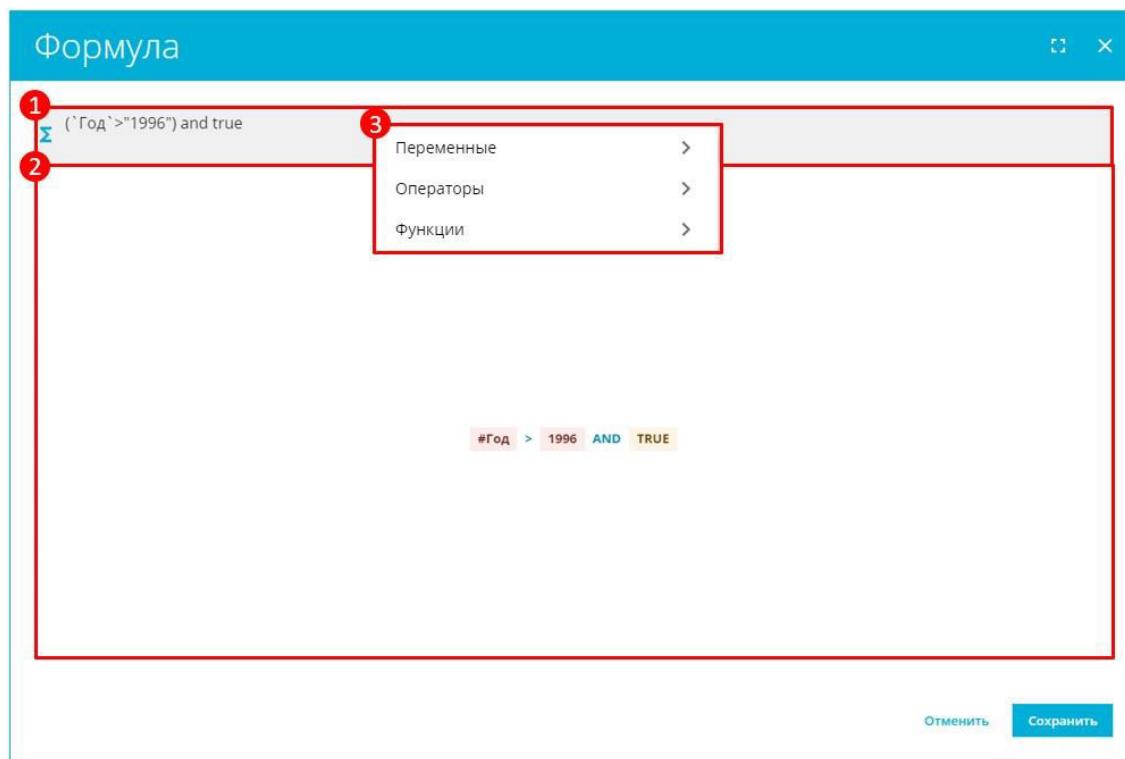


Рисунок 84 Окно Формула

1. Стока ввода текстового представления условия.
2. Рабочее поле, в котором строится графическое представление условия.
3. Вкладка, которая открывается при щелчке правой кнопкой мыши в строке ввода, и включает в себя:
  - Переменные из набора данных.
  - Операторы и функции, представленные в таблице ниже.

Позволяет просматривать переменные набора данных и строить конструкции условий.

**Для задания условия** необходимо:

- В строку текстового представления ввести условие, используя операторы, функции (Таблица 8) и переменные набора данных.
  - Название переменной должно быть указано в обратных одинарных кавычках (` `). Пример: `Год` > 1996.
  - В качестве десятичного разделителя используется точка.
  - Стока должна указываться в двойных кавычках. Пример: `Пол` = "мужской".
- Посмотреть список операторов, функций и переменных, а также добавить их в условие можно и в панели, которая открывается при щелчке правой кнопкой мыши по строке ввода.

### Список функций и операторов

| Оператор/Функция | Описание                  | Текстовое представление (формула)  |
|------------------|---------------------------|--|
| <b>And</b>       | Логическая операция «И»   | true <b>and</b> true<br>Вместо true нужно вставить условие<br><b>Пример:</b> `Год` > 1996 and `Пол` = "Женский"                |
| <b>Or</b>        | Логическая операция «ИЛИ» | true <b>or</b> true<br>Вместо true нужно вставить условие<br><b>Пример:</b> `Год` > 1996 or `Пол` = "Женский"                  |
| <b>+</b>         | Операция сложения         | 0 + 0<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` + `Цена продукта M` |
| <b>-</b>         | Операция вычитания        | 0 - 0<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` - 1000              |
| <b>*</b>         | Операция умножения        | 0 * 0<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` * `Скидка`          |
| <b>/</b>         | Операция деления          | 0 / 0  |

| Оператор/Функция  | Описание  | Текстовое представление (формула)  |
|-------------------|---|--|
|                   |   | Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Выручка за первый квартал` / `Выручка за год` |
| <                 | Меньше  | "" < ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Возраст` < 32                                       |
| ≤                 | Меньше или равно  | "" <= ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Температура` <= 120                                |
| >                 | Больше  | "" > ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Возраст` > "32"                                     |
| ≥                 | Больше или равно  | "" >= ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Температура` >= 120                                |
| =                 | Равно   | "" = ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Температура` = 120                                  |
| ≠                 | Не равно  | "" != ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Температура` != 120                                |
| <b>like</b>       | Проверяет, удовлетворяет ли символьная строка заданному образцу, который может содержать поисковые символы. Учитывает регистр | "" like ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` like "Москва"                             |
| <b>ilike</b>      | Проверяет, удовлетворяет ли символьная строка заданному образцу, который может содержать поисковые символы.                   | "" ilike ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` ilike "Москва"                           |
| <b>startswith</b> | Проверяет, есть ли в начале одной текстовой строки другая текстовая строка  | "" startswith ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` startswith "Сан"                    |
| <b>endswith</b>   | Проверяет, есть ли в конце одной текстовой строки другая текстовая строка   | "" endswith ""<br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Город` endswith "бург"                       |
| <b>contains</b>   | Проверяет, встречается ли указанная строка внутри другой строки   | "" contains ""<br>Вместо "" нужно вставить переменную/значение   |

| Оператор/Функция         | Описание   | Текстовое представление (формула)  |
|--------------------------|--|--|
|                          |  | <b>Пример:</b> `Название` contains "consulting"  |
| <b>between</b>           | Проверяет, входит ли значение в заданный диапазон  | <b>0 between(0,0)</b><br>Вместо 0 нужно вставить числовую переменную/значение<br><b>Пример:</b> `ID` between (1, 100500)                                   |
| <b>in</b>                | Проверяет наличие элемента в последовательности  | <b>"" in("")</b><br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Регион` in ("Москва", "Московская область")                            |
| <b>not in</b>            | Проверяет отсутствие элемента в последовательности   | <b>"" not_in("")</b><br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> `Регион` not_in ("Москва", "Московская область")                    |
| <b>coalesce</b>          | Возвращает данные из первого столбца, содержащего значение, отличное от NULL                     | <b>coalesce("")</b><br>Вместо "" нужно вставить переменную/значение<br><b>Пример:</b> coalesce(`ProductNumber`, `ProductName`)                             |
| <b>not</b>               | Задает противоположное условие   | <b>not(true)</b><br>Вместо true нужно вставить условие<br><b>Пример:</b> not(`Год` > 1996)   |
| <b>create_date</b>       | Создает переменную даты из последовательно введенных года, месяца и дня                          | <b>create_date(0,0,0)</b><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> create_date(2016,10,25)                       |
| <b>current_date</b>      | Возвращает текущую дату  | <b>current_date()</b><br>Вводится без дополнительных параметров  |
| <b>current_timestamp</b> | Возвращает текущие дату и время  | <b>current_timestamp()</b><br>Вводится без дополнительных параметров   |
| <b>now</b>               | Возвращает текущие дату и время  | <b>now()</b><br>Вводится без дополнительных параметров   |
| <b>create_datetime</b>   | Создает переменную даты времени из последовательно введенных года месяца дня часа минут и секунд | <b>create_datetime(0,0,0,0,0,0)</b><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> create_datetime(2016,10,25,12,18,0) |
| <b>char_length</b>       | Возвращает длину строки  | <b>char_length("")</b><br>Вместо "" нужно вставить строку/строковую переменную<br><b>Пример:</b> char_length(`Код_продукта`)                               |
| <b>random</b>            | Возвращает случайное число   | <b>random()</b><br>Вводится без дополнительных параметров  |
| <b>In</b>                | Натуральный логарифм   | <b>In(0)</b><br>Вместо 0 нужно вставить числовое значение/числовую переменную  |

| Оператор/Функция             | Описание   | Текстовое представление<br>(формула)  |
|------------------------------|--|---|
| <b>exp</b>                   | Экспонента   | <b>Пример:</b> ln(`Числовая переменная`)<br>exp(0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> exp(`Числовая переменная`)  |
| <b>power</b>                 | Возведение в степень   | power(0,0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> power(`Числовая переменная`,2)  |
| <b>sqrt</b>                  | Квадратный корень  | sqrt(0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> sqrt(`Числовая переменная`)  |
| <b>abs</b>                   | Абсолютное значение  | abs(0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> abs(`Числовая переменная`)  |
| <b>ceil</b>                  | Возвращает наименьшее целое число, которое больше или равно текущему значению                | ceil(0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> ceil(25.1)<br>Вернет значение 26   |
| <b>floor</b>                 | Возвращает наибольшее целое число, которое меньше или равно текущему значению                | floor(0)<br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> floor(25.1)<br>Вернет значение 25   |
| <b>extract_from_datetime</b> | Получает указанную часть (день, месяц, год, час, минута, секунда) из значения даты и времени | extract_from_datetime("","", "")<br>Вместо первых кавычек нужно указать необходимую часть для извлечения:<br>· SECOND<br>· MINUTE<br>· HOUR<br>· DAY<br>· MONTH<br>· YEAR<br>Вместо вторых кавычек нужно указать переменную типа datetime<br><b>Пример:</b><br>extract_from_datetime("DAY", `datetime`) |
| <b>extract_from_date</b>     | Получает указанную часть (день, месяц, год) из значения даты                                 | extract_from_date("","", "")<br>В первые кавычки нужно вписать необходимую часть для извлечения:<br>· DAY<br>· MONTH<br>· YEAR<br>Вместо вторых кавычек нужно указать переменную типа date<br><b>Пример:</b><br>extract_from_date("DAY", `date`)  |

| Оператор/Функция | Описание                                      | Текстовое представление (формула)   |
|------------------|---|---|
| <b>concat</b>    | Объединяет в единую строку указанные значения | concat("")<br>Вместо "" нужно вставить строки/строковые переменные<br><b>Пример:</b><br>concat(`Фамилия` , " ", `Имя` ) |

Таблица 12 Операторы и функции

### Результаты выполнения узла:

- Таблица с примером отфильтрованного набора данных. Отображаются первые 100 наблюдений.

| Результаты выполнения узла |        |        |          |        |         |              |              |             |         |         |         |         |         |         |    |           |           |           |
|----------------------------|--------|--------|----------|--------|---------|--------------|--------------|-------------|---------|---------|---------|---------|---------|---------|----|-----------|-----------|-----------|
| Табличные результаты       |        |        |          |        |         |              |              |             |         |         |         |         |         |         |    |           |           |           |
| Пример данных              |        |        |          |        |         |              |              |             |         |         |         |         |         |         |    |           |           |           |
| id                         | Number | symbol | name     | mass   | density | melting_t... | boiling_t... | electron... | Level 1 | Level 2 | Level 3 | Level 4 | Level 5 | Level 6 | Le | Filter... | Filter... | Filter... |
| 1                          | 1      | Н      | Водород  | 1.007  | 0       | -259.14      | -252.87      | 2.02        | 1       | 0       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 2                          | 2      | Не     | Гелий    | 4.002  | 0       | -272.2       | -268.93      | 12.3        | 2       | 0       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 3                          | 3      | Li     | Литий    | 6.941  | 0.534   | 180.54       | 1347         | 0.98        | 2       | 1       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 4                          | 4      | Be     | Бериллий | 9.012  | 1.85    | 1278         | 2970         | 1.57        | 2       | 2       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 5                          | 5      | B      | Бор      | 10.81  | 2.34    | 2210         | 2600         | 2.04        | 2       | 3       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 6                          | 6      | C      | Углерод  | 12.011 | 2.265   | 3550         | 4827         | 2.55        | 2       | 4       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 7                          | 7      | N      | Азот     | 14.006 | 0.001   | -209.86      | -195.8       | 3.04        | 2       | 5       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 8                          | 8      | O      | Кислород | 15.999 | 0.001   | -218.4       | -182.96      | 3.44        | 2       | 6       | 0       | 0       | 0       | 0       | 0  |           |           |           |
| 9                          | 9      | F      | Фтор     | 18.998 | 0.001   | -219.62      | -188.11      | 3.98        | 2       | 7       | 0       | 0       | 0       | 0       | 0  |           |           |           |

Рисунок 85 Пример таблицы с отфильтрованным набором данных (условие id <10)

### 3.2.5.6.5. Узел «Метаданные»

Узел «Метаданные» позволяет изменить метаданные переменных.

К метаданным переменной относятся **Роль**, **Тип** и **Целевой класс**.

В Модуле предусмотрены следующие **Роли** переменной:

- Target — Целевая переменная** — зависимая переменная, используемая в качестве цели моделирования.
- Predictor — Предиктор** — независимая переменная, используемая в качестве прогнозирующей целевую переменную.
- Excluded — Исключен** — данная роль позволяет исключить переменную из процесса моделирования. Вернуть ее обратно в процесс можно также при помощи узла «Метаданные», установленного далее в сценарии.
- Weight — Вес** — данная роль задается переменной, которая указывает модели на важность наблюдений. Пример: для прогноза могут быть важны тенденции за последний период (например, сезон).
- Cluster\_ID — Сегмент** — переменная с данной ролью рассчитывается при решении задачи кластеризации.

- **Prediction — Прогноз** — переменная с данной ролью рассчитывается при решении задачи регрессии.
- **Classification — Прогноз класс** — переменная с данной ролью рассчитывается при решении задачи классификации.
- **Partition — Разделение** — переменная с данной ролью задает разделение набора данных на обучающую, валидационную и тестовую выборки (рассчитывается в узле «Разделение выборки»).
- **ClassWeight — Вес класса** — переменная с данной ролью корректирует несбалансированность классов (задается в узле «Веса классов»).

Первоначально при выборе набора данных **в узле «Набор данных»** все атрибуты имеют роль **Предиктор**.

Предусмотрены следующие **Типы** переменной:

- **Номинальный (Nominal)** — дискретные значения без числовой связи между категориями.
- **Порядковый (Ordinal)** — дискретные значения, которые можно ранжировать или сортировать (Пример — хорошо, великолепно, превосходно).
- **Двоичный (Binary)** — это дискретные данные, которые могут относиться только к одной из двух категорий: 0 и 1.
- **Интервальный (Interval)** — числовые значения.

**Целевой класс** используется в процессе моделирования (например, при стратификации) и может быть указан только для **Целевой переменной**.

**Список параметров узла** представлен в таблице ниже.

| Параметр             | Возможные значения и ограничения           | Описание   |
|----------------------|--|--|
| Название             | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе         |
| Описание             | Ручной ввод<br>Ограничений на значение нет | Описание узла  |
| Изменение метаданных | Кнопка                                     | При выборе кнопки откроется окно <b>Изменение метаданных</b> . |

**Таблица 13 Параметры узла «Метаданные»**

#### Окно Изменение метаданных

В окне **Изменение метаданных** Пользователь имеет возможность изменить **Роль**, **Тип** переменной и указать **Целевой класс**. Для этого необходимо:

- Рядом с интересующей переменной выбрать иконку
- В выпадающих меню выбрать **Роль** и/или **Тип**, указать **Целевой класс** (доступно только для переменной типа **Целевая переменная -Target**).
- Сохранить изменения, выбрав иконку

**Изменение метаданных**

| Пользовательское имя атрибута | Роли               | Тип      | Целевой класс |
|-------------------------------|--------------------|----------|---------------|
| sepal_length_in_cm            | Предиктор          | Interval |               |
| sepal_width_in_cm             | Предиктор          | Interval |               |
| petal_length_in_cm            | Предиктор          | Interval |               |
| petal_width_in_cm             | Предиктор          | Interval |               |
| class                         | Целевая переменная | Nominal  | Iris-setosa   |

Поиск... Настройка Таблицы

Отмена Сохранить

**Рисунок 86 Окно Изменение метаданных**

### Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

**Пример данных**

| sepal_length... | sepal_width_i... | petal_length_in... | petal_width_in... | class       |
|-----------------|------------------|--------------------|-------------------|-------------|
| 5.1             | 3.5              | 1.4                | 0.2               | Iris-setosa |
| 4.9             | 3                | 1.4                | 0.2               | Iris-setosa |
| 4.7             | 3.2              | 1.3                | 0.2               | Iris-setosa |
| 4.6             | 3.1              | 1.5                | 0.2               | Iris-setosa |
| 5               | 3.6              | 1.4                | 0.2               | Iris-setosa |
| 5.4             | 3.9              | 1.7                | 0.4               | Iris-setosa |

**Рисунок 87 Таблица с примером данных**

### 3.2.5.6.6. Узел «One-hot encoding»

**Узел «One-hot encoding»** преобразует категориальные данные в числовую форму.

Многие алгоритмы не могут напрямую работать с категориальными переменными. Для этого предусмотрен метод **One-hot encoding**, который преобразует категориальные данные в числовую форму. Для этого создаются дополнительные столбцы-индикаторы наличия/отсутствия категории с помощью значений 1 или 0 соответственно. Таким образом, если категориальная переменная имеет k возможных значений, то на выходе получится k столбцов для ее представления. Алгоритмы машинного обучения могут принимать эти столбцы в качестве входных данных.

| № | Категория |                  | № | Категория_A | Категория_B |
|---|-----------|------------------|---|-------------|-------------|
| 1 | А         |                  | 1 | 1           | 0           |
| 2 | Б         | One-hot encoding | 2 | 0           | 1           |
| 3 | Б         |                  | 3 | 0           | 1           |
| 4 | А         |                  | 4 | 1           | 0           |

**Рисунок 88 Принцип работы узла «One-hot encoding»**

В кодировании One-Hot есть некоторая избыточность. Например, переменная «Пол» может принимать два значения - мужчина или женщина. При кодировании достаточно использовать в качестве предиктора лишь одну из этих двух фиктивных переменных. Для этого

Пользователю нужно выбрать чекбокс **Исключить первую категорию** или вручную указать исключаемую переменную в окне **Выбор категории**.

**Список параметров узла** представлен в таблице ниже.

| Параметр                                 | Возможные значения и ограничения           | Описание   |
|--|--|--|
| <b>Название</b>                          | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                          | Ручной ввод<br>Ограничений на значение нет | Описание узла  |
| <b>Исключить первую категорию</b>        | Чекбокс                                    | Выбор данного чекбокса указывает методу на необходимость исключить первую категорию  |
| <b>Исключить оригинальную переменную</b> | Чекбокс                                    | Выбор данного чекбокса указывает методу на необходимость исключить оригинальную переменную, на основе которой были вычислены фиктивные, из дальнейшего процесса моделирования. В окне <b>Выходные переменные</b> Роль данного атрибута изменится на Исключен (Excluded) (про Роли переменной подробнее Узел «Метаданные»). |
| <b>Выбор категории</b>                   | Кнопка                                     | При выборе кнопки откроется окно <b>Выбор категории</b> .  |

**Таблица 14 Параметры узла «One-hot encoding»**

## Окно Выбор категории

В окне **Выбор категории** Пользователь имеет возможность выбрать категориальную переменную для кодирования (Рисунок 83). Для этого необходимо:

- Рядом с интересующей переменной выбрать иконку .
- В столбце **Выбрать** выбрать чекбокс и при необходимости указать исключаемые из кодирования значения (перечислить через запятую).
- Сохранить изменения, выбрав иконку .

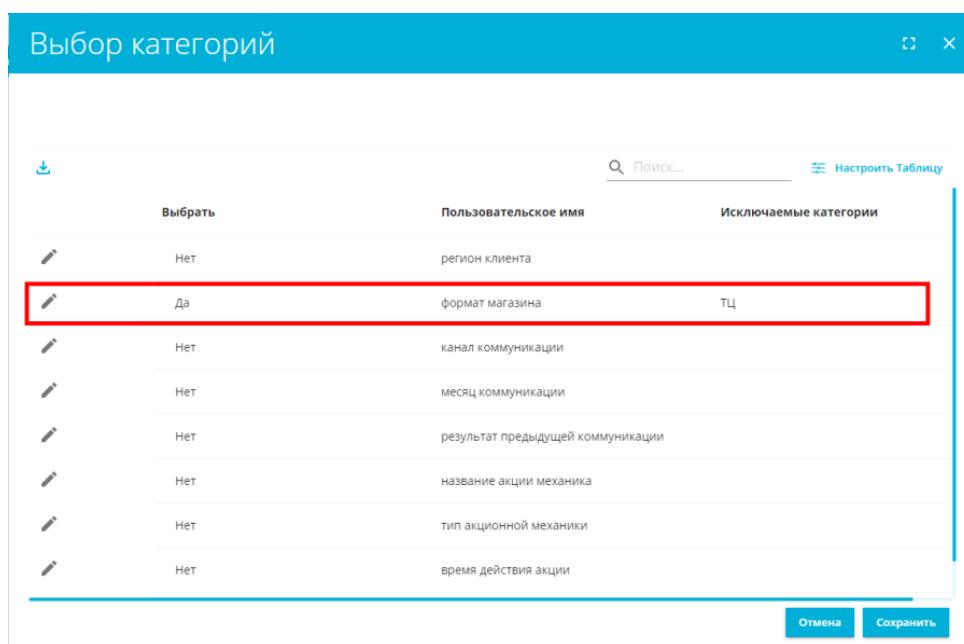


Рисунок 89 Окно Выбор категории

## Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

| на | формат магазина_магазин у дома | формат магазина_неизвестно | формат магазина_франчайзинг |
|----|--------------------------------|----------------------------|-----------------------------|
| 0  | 0                              | 0                          | 0                           |
| 0  | 0                              | 0                          | 0                           |
| 0  | 0                              | 0                          | 1                           |
| 0  | 0                              | 0                          | 0                           |
| 0  | 0                              | 0                          | 0                           |

Рисунок 90 Таблица с примером посчитанных трех фиктивных переменных

- Таблица с результатами кодирования.

| Результаты кодирования |   | Поиск...              | Настроить Таблицу |
|------------------------|---|-----------------------|-------------------|
| Переменная             | Закодированные категории                | Исключенные категории |                   |
| формат магазина        | магазин у дома, неизвестно, франчайзинг |                       | ТЦ                |
|                        |   |                       |                   |

**Рисунок 91 Пример таблицы с результатами кодирования**

В результате выполнения в выходных параметрах узла появятся столбцы-индикаторы с ролью первоначального столбца.

### 3.2.5.6.7. Узел «Заполнение пропусков»

**Узел «Заполнение пропусков»** обрабатывает пропущенные значения.

В зависимости от задачи Пользователь может использовать тот или иной метод заполнения отсутствующих элементов. Заменить пропуски можно на:

- **Моду** – наиболее часто встречающееся значение (подходит для категориальных переменных).
- **Константу** – выбранное Пользователем конкретное значение.
- **Среднее** – находится суммированием всех чисел в выборке и делением полученной суммы на количество чисел.
- **Медиану** – если взять все элементы множества и отсортировать, то это число делит множество пополам. Одна половина множества равна или больше этого числа, а другая меньше или равна этому числу.
- **Минимум.**
- **Максимум.**

**Список параметров узла** представлен в таблице ниже.

| Параметр                | Возможные значения и ограничения           | Описание  |
|-------------------------|--|---|
| Название                | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе  |
| Описание                | Ручной ввод<br>Ограничений на значение нет | Описание узла   |
| Создать общий индикатор | Чекбокс                                    | Выбор данного чекбокса указывает на необходимость расчета общего для всех переменных набора данных индикатора |

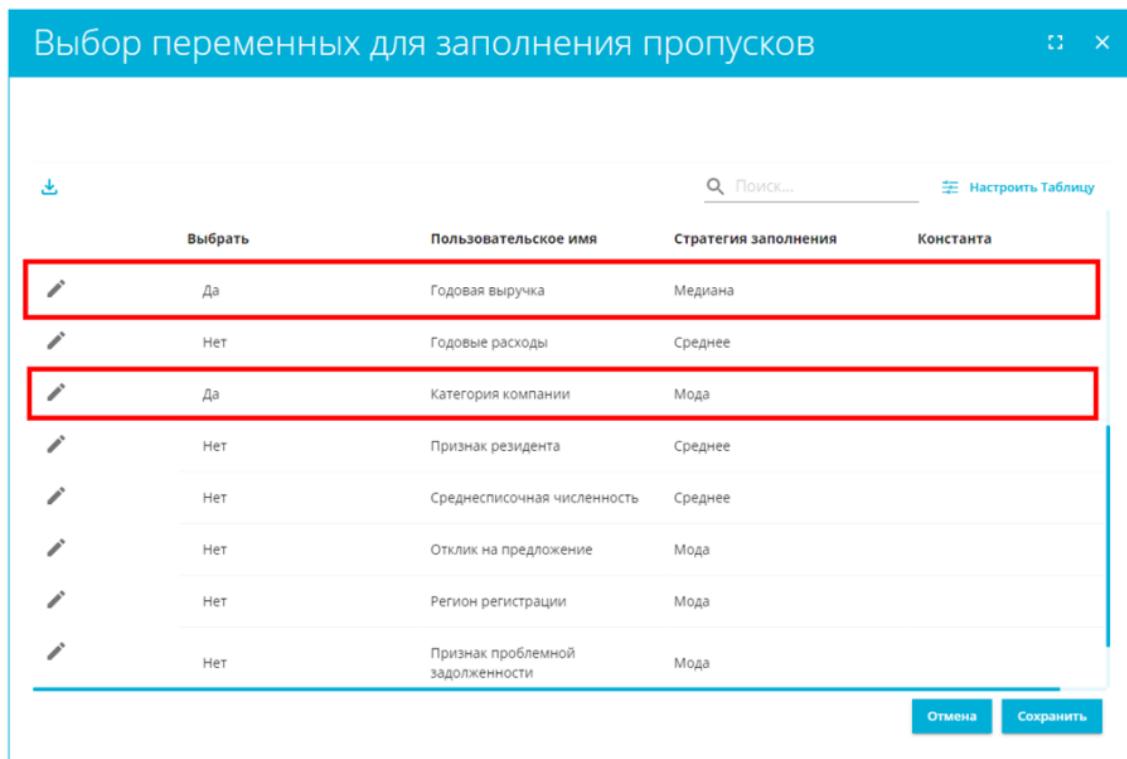
| Параметр   | Возможные значения и ограничения | Описание  |
|--|----------------------------------|---|
| пропущенных значений                                   |                                  | пропущенных значений (столбца, в котором будет указано наличие или отсутствие пропущенных значений в наблюдении). После успешного выполнения узла созданный индикатор можно найти по ссылке « <b>Выходные параметры</b> ». Ему будет назначена роль <b>Предиктора</b> с типом <b>Binary</b> .   |
| Создать индивидуальные индикаторы пропущенных значений | Чекбокс                          | Выбор данного чекбокса указывает на необходимость расчета индивидуальных индикаторов пропущенных значений для выбранной далее переменной (для каждой из выбранных переменных будет посчитан свой столбец).<br>После успешного выполнения узла созданный индикатор можно найти по ссылке « <b>Выходные параметры</b> ». Ему будет назначена роль <b>Предиктора</b> с типом <b>Binary</b> . |
| Максимальная доля пропусков                            | Ручной ввод                      | Данный параметр указывает максимальную долю пропусков   |
| Выбор переменных для заполнения пропусков              | Кнопка                           | При выборе кнопки откроется окно <b>Выбор переменных для заполнения пропусков</b> .   |

**Таблица 15 Параметры узла «Заполнение пропусков»**

#### **Окно Выбор переменных для заполнения пропусков**

В окне **Выбор переменных для заполнения пропусков** Пользователь может задать метод замены пропущенных данных (Рисунок 86). Для этого необходимо:

- Рядом с интересующей переменной выбрать иконку  .
- В столбце **Выбрать** нажать на чекбокс и при необходимости изменить **Стратегию заполнения**.
  - При выборе **Стратегии заполнения** константой нужно указать значение в соответствующем столбце **Константа**.
- Сохранить изменения, выбрав иконку  .



**Рисунок 92 Окно Выбор переменных для заполнения пропусков**

#### Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

| Пример данных                   |                                   |                          | IMP_Годовая выручка |   | IMP_Категория компании |  |
|---------------------------------|-----------------------------------|--------------------------|---------------------|---|------------------------|--|
| Среднесписочн...<br>численность | Организацион...<br>правовая форма | Отклик на<br>предложение |                     |   |                        |  |
| 168                             | Акционерное<br>общество           | 0                        | 12000               | 0 |                        |  |
| 177                             | Акционерное<br>общество           | 0                        | 5410                | 0 |                        |  |
| 165                             | Акционерное<br>общество           | 0                        | 4030                | 1 |                        |  |
| 165                             | Акционерное<br>общество           | 0                        | 4500                | 0 |                        |  |
| 225                             | Акционерное<br>общество           | 0                        | 3100                | 0 |                        |  |
| 159                             | Акционерное<br>общество           | 0                        | 4400                | 1 |                        |  |

**Рисунок 93 Окно Выбор переменных для заполнения пропусков**

В результате выполнения узла в наборе данных будут рассчитаны новые переменные с заполненными пропущенными значениями (с префиксом **IMP\_**).

- Таблица с количеством пропусков по каждой переменной набора данных.

| Количество пропусков | Поиск... | Настройка Таблицы |
|----------------------|----------|-------------------|
| Переменная           | Пропуски |                   |
| Годовая выручка      | 12011    |                   |
| Годовые расходы      | 0        |                   |
| Категория компании   | 1526     |                   |

**Рисунок 94 Пример таблицы с количеством пропущенных значений по каждой переменной**

- Таблица замененных значений.

| Заменённые значения | Поиск... | Настройка Таблицы |
|---------------------|----------|-------------------|
| Переменная          | Значение |                   |
| Годовая выручка     | 5410     |                   |
| Категория компании  | 0        |                   |

**Рисунок 95 Пример таблицы с замененными значениями**

### 3.2.5.6.8. Узел «Трансформация»

**Узел «Трансформация»** позволяет рассчитать новые переменные.

Пользователь может получить дополнительные сведения для моделирования, выполнив вычисления на основе имеющихся переменных. Для этого в узле «Трансформация» предусмотрено создание вычисляемых показателей при помощи расчетных формул.

**Список параметров узла** представлен в таблице ниже.

| Параметр                  | Возможные значения и ограничения           | Описание   |
|---------------------------|--|--|
| <b>Название</b>           | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе       |
| <b>Описание</b>           | Ручной ввод<br>Ограничений на значение нет | Описание узла  |
| <b>Расчетные атрибуты</b> | Кнопка                                     | При выборе кнопки откроется окно <b>Расчетные атрибуты</b> . |

**Таблица 16 Параметры узла «Трансформация»**

## Окно Расчетные атрибуты

В окне **Расчетные атрибуты** Пользователь имеет возможность создать новые переменные. Для этого необходимо:

- Выбрать иконку  .
- Задать следующие параметры:
  - Имя атрибута.
  - Формулу. Для этого нужно нажать на поле ввода, откроется окно **Формула** в котором и задается формула расчета.
  - Тип (подробнее в разделе Входные и выходные данные узла).  
Задается автоматически после задания расчетной формулы.
  - Роль (подробнее в разделе Входные и выходные данные узла).
- Сохранить изменения, выбрав иконку  .

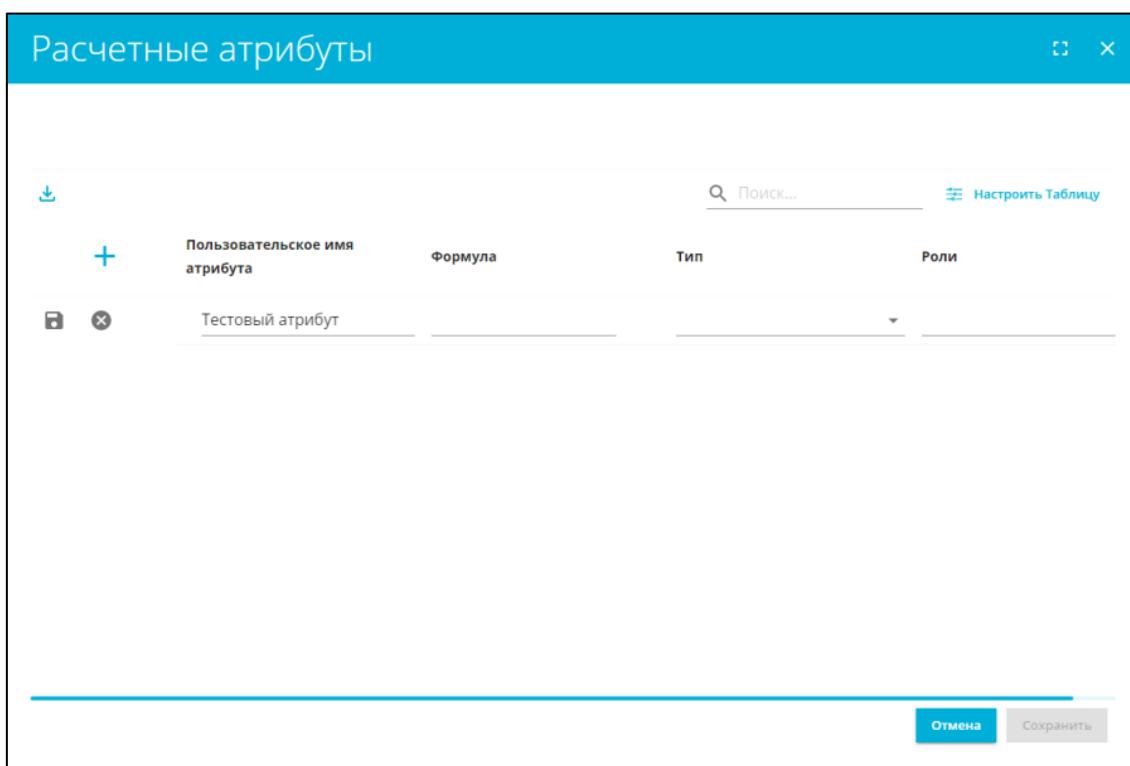
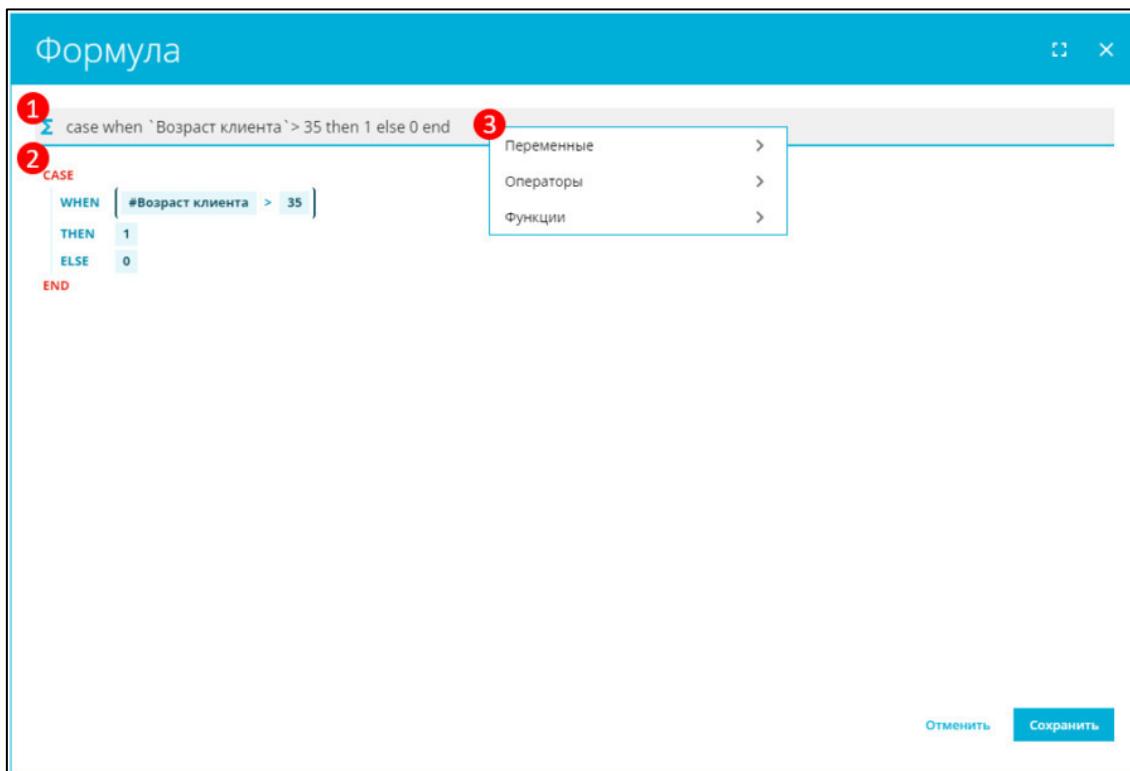


Рисунок 96 Окно Расчетные атрибуты

## Окно Формула

В окне **Формула** задается расчетная формула.

Основными элементами окна Формула являются:



**Рисунок 97 Окно Формула**

1. Стока ввода текстового представления условия.
2. Рабочее поле, в котором строится графическое представление условия.
3. Вкладка, которая открывается при щелчке правой кнопкой мыши в строке ввода, и включает в себя:
  - Переменные из набора данных.
  - Операторы и функции, представленные в таблице ниже.

Позволяет просматривать переменные набора данных и строить расчетные формулы.

Для **задания условия** необходимо:

- В строку текстового представления ввести формулу, используя операторы, функции и переменные набора данных.
  - Название переменной должно быть указано в обратных одинарных кавычках (` `). Пример: `Год` > 1996.
  - В качестве десятичного разделитель используется точка.
  - Стока должна указываться в двойных кавычках. Пример: `Пол` = "мужской".
- Посмотреть список операторов, функций и переменных, а также добавить их в условие можно и в панели, которая открывается при щелчке правой кнопкой мыши по строке ввода.

| Оператор/Функция | Описание                                      | Текстовое представление (формула)  |
|------------------|---|--|
| <b>And</b>       | Логическая операция «И»                       | <code>true and true</code><br>Вместо true нужно вставить условие<br><b>Пример:</b> `Год` > 1996 and `Пол` = "Женский"                                |
| <b>Or</b>        | Логическая операция «ИЛИ»                     | <code>true or true</code><br>Вместо true нужно вставить условие<br><b>Пример:</b> `Год` > 1996 or `Пол` = "Женский"                                  |
| <b>+</b>         | Операция сложения                             | <code>0 + 0</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` + ` Цена продукта M`         |
| <b>-</b>         | Операция вычитания                            | <code>0 - 0</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` - 1000                       |
| <b>*</b>         | Операция умножения                            | <code>0 * 0</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Цена продукта M` * `Скидка`                   |
| <b>/</b>         | Операция деления                              | <code>0 / 0</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> `Выручка за первый квартал` / `Выручка за год` |
| <b>&lt;</b>      | Меньше  | <code>"" &lt; ""</code><br>Вместо ` ` нужно вставить переменную/значение<br><b>Пример:</b> `Возраст` < 32  |
| <b>≤</b>         | Меньше или равно                              | <code>"" &lt;= ""</code><br>Вместо ` ` нужно вставить переменную/значение<br><b>Пример:</b> `Температура` <= 120                                     |
| <b>&gt;</b>      | Больше  | <code>"" &gt; ""</code><br>Вместо ` ` нужно вставить переменную/значение<br><b>Пример:</b> `Возраст` > "32"  |
| <b>≥</b>         | Больше или равно                              | <code>"" &gt;= ""</code><br>Вместо ` ` нужно вставить переменную/значение<br><b>Пример:</b> `Температура` >= 120                                     |
| <b>=</b>         | Равно   | <code>"" = ""</code><br>Вместо ` ` нужно вставить переменную/значение<br><b>Пример:</b> `Температура` = 120  |
| <b>≠</b>         | Не равно                                      | <code>"" != ""</code><br>Вместо ` ` нужно вставить переменную/значение<br><b>Пример:</b> `Температура` != 120  |
| <b>like</b>      | Проверяет, удовлетворяет ли символьная строка | <code>"" like ""</code><br>Вместо ` ` нужно вставить переменную/значение   |

| Оператор/Функция  | Описание  | Текстовое представление (формула)   |
|-------------------|---|---|
|                   | заданному образцу, который может содержать поисковые символы. Учитывает регистр                             | <b>Пример:</b> `Город` like "Москва"  |
| <b>ilike</b>      | Проверяет, удовлетворяет ли символьная строка заданному образцу, который может содержать поисковые символы. | <pre>"" ilike ""</pre> <p>Вместо "" нужно вставить переменную/значение</p> <b>Пример:</b> `Город` ilike "Москва"  |
| <b>startswith</b> | Проверяет, есть ли в начале одной текстовой строки другая текстовая строка                                  | <pre>"" startswith ""</pre> <p>Вместо "" нужно вставить переменную/значение</p> <b>Пример:</b> `Город` startswith "Сан"   |
| <b>endswith</b>   | Проверяет, есть ли в конце одной текстовой строки другая текстовая строка                                   | <pre>"" endswith ""</pre> <p>Вместо "" нужно вставить переменную/значение</p> <b>Пример:</b> `Город` endswith "бург"  |
| <b>contains</b>   | Проверяет, встречается ли указанная строка внутри другой строки   | <pre>"" contains ""</pre> <p>Вместо "" нужно вставить переменную/значение</p> <b>Пример:</b> `Название` contains "consulting"   |
| <b>between</b>    | Проверяет, входит ли значение в заданный диапазон   | <pre>0 between(0,0)</pre> <p>Вместо 0 нужно вставить числовую переменную/значение</p> <b>Пример:</b> `ID` between (1, 100500)   |
| <b>in</b>         | Проверяет наличие элемента в последовательности   | <pre>"" in("")</pre> <p>Вместо "" нужно вставить переменную/значение</p> <b>Пример:</b> `Регион` in ("Москва", "Московская область")  |
| <b>not in</b>     | Проверяет отсутствие элемента в последовательности  | <pre>"" not_in("")</pre> <p>Вместо "" нужно вставить переменную/значение</p> <b>Пример:</b> `Регион` not_in ("Москва", "Московская область")  |
| <b>CASE</b>       | Позволяет проверить несколько условий и выполнить разные операции на основе этих условий                    | <pre>case when true then "" else "" end</pre> <p>Вместо true нужно вставить условие, вместо "" – значение</p> <p>Для проверки дополнительных условий необходимо добавить дополнительные &lt;when true then ""&gt;</p> <b>Пример:</b> case when `Возраст` > 32 then 0 else 1 end |
| <b>coalesce</b>   | Возвращает данные из первого столбца, содержащего значение, отличное от NULL                                | <pre>coalesce("")</pre> <p>Вместо "" нужно вставить переменную/значение</p> <b>Пример:</b> coalesce(`ProductNumber`, `ProductName`)   |

| Оператор/Функция         | Описание   | Текстовое представление (формула)  |
|--------------------------|--|--|
| <b>not</b>               | Задает противоположное условие   | <code>not(true)</code><br>Вместо true нужно вставить условие<br><b>Пример:</b> not(`Год` > 1996)   |
| <b>create_date</b>       | Создает переменную даты из последовательно введенных года, месяца и дня                          | <code>create_date(0,0,0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> create_date(2016,10,25)                       |
| <b>current_date</b>      | Возвращает текущую дату  | <code>current_date()</code><br>Вводится без дополнительных параметров  |
| <b>current_timestamp</b> | Возвращает текущие дату и время  | <code>current_timestamp()</code><br>Вводится без дополнительных параметров   |
| <b>now</b>               | Возвращает текущие дату и время  | <code>now()</code><br>Вводится без дополнительных параметров   |
| <b>create_datetime</b>   | Создает переменную даты времени из последовательно введенных года месяца дня часа минут и секунд | <code>create_datetime(0,0,0,0,0,0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> create_datetime(2016,10,25,12,18,0) |
| <b>char_length</b>       | Возвращает длину строки  | <code>char_length("")</code><br>Вместо "" нужно вставить строку/строковую переменную<br><b>Пример:</b> char_length(`Код_продукта`)                               |
| <b>random</b>            | Возвращает случайное число   | <code>random()</code><br>Вводится без дополнительных параметров  |
| <b>In</b>                | Натуральный логарифм   | <code>In(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> In(`Числовая переменная`)                                  |
| <b>exp</b>               | Экспонента   | <code>exp(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> exp(`Числовая переменная`)                                |
| <b>power</b>             | Возведение в степень   | <code>power(0,0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> power(`Числовая переменная`,2)                        |
| <b>sqrt</b>              | Квадратный корень  | <code>sqrt(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> sqrt(`Числовая переменная`)                              |
| <b>abs</b>               | Абсолютное значение  | <code>abs(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> abs(`Числовая переменная`)                                |

| Оператор/Функция             | Описание   | Текстовое представление (формула)   |
|------------------------------|--|---|
| <b>ceil</b>                  | Возвращает наименьшее целое число, которое больше или равно текущему значению                | <code>ceil(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> <code>ceil(25.1)</code><br>Вернет значение 26   |
| <b>floor</b>                 | Возвращает наибольшее целое число, которое меньше или равно текущему значению                | <code>floor(0)</code><br>Вместо 0 нужно вставить числовое значение/числовую переменную<br><b>Пример:</b> <code>floor(25.1)</code><br>Вернет значение 25   |
| <b>extract_from_datetime</b> | Получает указанную часть (день, месяц, год, час, минута, секунда) из значения даты и времени | <code>extract_from_datetime("", "")</code><br>Вместо первых кавычек нужно указать необходимую часть для извлечения:<br><ul style="list-style-type: none"> <li>· SECOND</li> <li>· MINUTE</li> <li>· HOUR</li> <li>· DAY</li> <li>· MONTH</li> <li>· YEAR</li> </ul> Вместо вторых кавычек нужно указать переменную типа datetime<br><b>Пример:</b><br><code>extract_from_datetime("DAY", `datetime`)</code> |
| <b>extract_from_date</b>     | Получает указанную часть (день, месяц, год) из значения даты                                 | <code>extract_from_date("", "")</code><br>В первые кавычки нужно вписать необходимую часть для извлечения:<br><ul style="list-style-type: none"> <li>· DAY</li> <li>· MONTH</li> <li>· YEAR</li> </ul> Вместо вторых кавычек нужно указать переменную типа date<br><b>Пример:</b><br><code>extract_from_date("DAY", `date`)</code>  |
| <b>concat</b>                | Объединяет в единую строку указанные значения  | <code>concat("")</code><br>Вместо "" нужно вставить строки/строковые переменные<br><b>Пример:</b><br><code>concat(`Фамилия`, " ", `Имя`)</code>   |

Таблица 17 Операторы и функции

#### Результаты выполнения узла:

- Таблица с примером набора данных с посчитанными атрибутами. Отображаются первые 100 наблюдений.

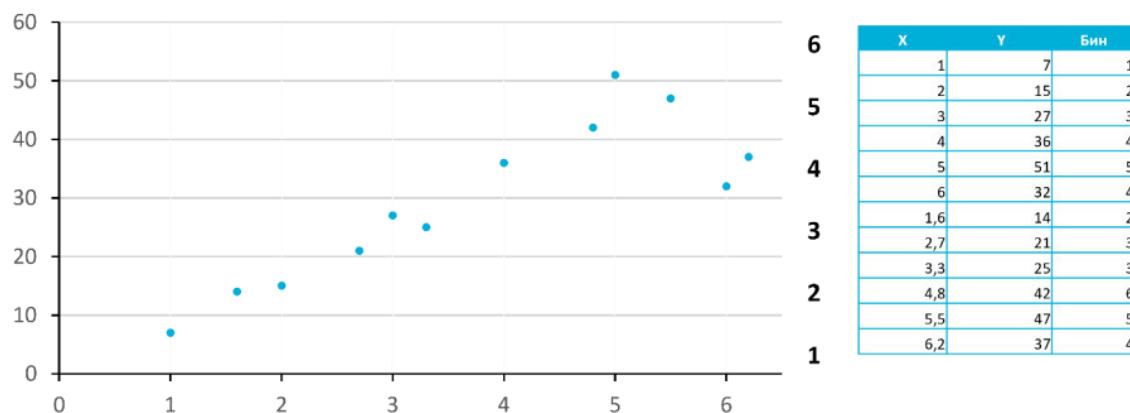
| Пример данных      |   |                               |                                       |                                 |                                   |                       |
|--------------------|---|-------------------------------|---------------------------------------|---------------------------------|-----------------------------------|-----------------------|
| Месяц коммуникации | Количество дней после регистрации клиента | Общее количество коммуникаций | Число дней с предыдущей коммуникацией | Наличие предыдущих коммуникаций | Результат предыдущей коммуникации | Отклик на предложение |
| май                | 1042                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     |
| май                | 1467                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     |
| май                | 1389                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     |
| май                | 579                                       | 1                             | -1                                    | 0                               | неизвестно                        | 1                     |
| май                | 673                                       | 2                             | -1                                    | 0                               | неизвестно                        | 1                     |
| май                | 562                                       | 2                             | -1                                    | 0                               | неизвестно                        | 1                     |
| май                | 1201                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     |
| май                | 1030                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     |

**Рисунок 98 Пример таблицы с набором данных с рассчитанным атрибутом Тестовый атрибут**

### 3.2.5.6.9. Узел «Биннинг/энкодинг»

**Узел «Биннинг/энкодинг»** включает в себя методы биннинга интервальных переменных и кодирования категориальных переменных.

**Биннинг** – процесс разделения диапазона непрерывной переменной на k интервалов. Может использоваться для сокращения размерности данных, что часто повышает точность модели.



**Рисунок 99 Пример биннинга по значениям переменной Y**

Предусмотрены следующие методы бинаризации переменной:

- **Однаковая ширина.** Данная стратегия разделяет диапазон значений переменной на указанное количество равных интервалов.
- **Квантильный.** Данная стратегия разделяет диапазон значений переменной таким образом, чтобы в каждый из них попало примерно одинаковое количество значений.
- **Дерево.** Данная стратегия разделяет диапазон значений переменной с помощью дерева решений, что предполагает связь с целевой переменной.

Разбиение начинается со всех наблюдений, которые представлены корневым узлом дерева. Алгоритм разбивает этот родительский узел на дочерние узлы (и листья) таким образом, чтобы значения (уровни) целевой переменной в пределах каждого дочернего региона были максимально похожи (критерий разбиения задается в параметрах **Критерий разбиения для регрессии** и **Критерий разбиения для классификации**).

Соответственно, параметр **Количество бинов** равен **Максимальной глубине дерева** в исходном алгоритме дерева решений, а параметр **Минимальное количество наблюдений – Минимальному количеству наблюдений в листе дерева** (подробнее про работу алгоритма в справке **узла «Дерево решений»**).

Если хотя бы для одной переменной используется бинаризация на основе дерева, то должна быть задана **целевая переменная** (задается **в узле «Метаданные»**).

Предусмотрены отдельные бины для пропущенных значений и значений вне диапазона (параметры **Обработка пропущенных значений** и **Обработка значений вне диапазона**).

Большинство алгоритмов машинного обучения не могут обрабатывать категориальные переменные. **Энкодинг** – процесс преобразования текстовых атрибутов в числовые значения.

Закодировать категориальные переменные в числовые можно четырьмя методами:

- **Количество** (Count Encoding) – для каждой категории ставится количество наблюдений с этой категорией.
- **Частота** (Freq Encoding) – для каждой категории ставится частота наблюдений с этой категорией.
- **Целевая** (Target Encoding) – для каждой категории ставится средневзвешенное среднего целевой на подвыборке соответствующей заданной категории и среднего целевой на всей обучающей выборке (train). Параметр веса рассчитывается через количество наблюдений в подвыборке, соответствующей заданной категории, и параметры **Сдвига** и **Масштаба**.

Для корректной работы необходимо указать целевую переменную (задается в узле **«Метаданные»**).

- **WOE** (Weight of Evidence) – для каждой категории ставится WOE, рассчитанный для подвыборки, соответствующей выбранной категории. Математически **WOE** определяется как логарифм отношения доли «хороших» наблюдений к доле «плохих» наблюдений.

Для корректной работы необходимо указать целевую переменную (задается в узле **«Метаданные»**).

Предусмотрены **Обработка пропущенных значений** и **Обработка неизвестных значений**, которые зависят от выбранного метода энкодинга (подробнее в описании соответствующих параметров).

Если целевая переменная категориальная и в ней больше 2-х категорий, то кодировщики, основанные на целевой переменной (**Target Encoding** и **WOE**),

исключают одну из категорий целевой переменной и для каждой из оставшихся рассматривают соответствующую задачу бинарной классификации (1 – выбранная категория, 0 – все остальные) и строятся по ней. Таким образом, если у целевой переменной К категорий,  $K > 2$ , то каждый кодировщик, основанный на целевой переменной, генерирует  $K - 1$  столбцов в результатах.

**Список параметров узла** представлен в таблице.

| Параметр                                | Возможные значения и ограничения   | Описание   |
|---|--|--|
| <b>Название</b>                         | Ручной ввод<br>Ограничений на значение нет   | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                         | Ручной ввод<br>Ограничений на значение нет   | Описание узла  |
| <b>Префикс выходных переменных</b>      | Ручной ввод<br>По умолчанию — BIN_   | Данный параметр задает префикс для выходных переменных узла<br>Используется для биннинга интервальных переменных   |
| <b>Обработка значений вне диапазона</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Включить в крайние бины (по умолчанию)</li> <li>• Отдельные бины</li> <li>• Игнорировать</li> </ul> | Данный параметр указывает что делать со значениями вне диапазона.<br>Используется для биннинга интервальных переменных   |
| <b>Бин для пропущенных значений</b>     | Чекбокс  | Выбор данного чекбокса указывает, что необходимо посчитать отдельный бин для пропущенных значений<br>Используется для биннинга интервальных переменных   |
| <b>Переменные</b>                       | Кнопка   | Используется для биннинга интервальных переменных<br>При выборе кнопки « <b>Переменные</b> » открывается окно <b>Переменные</b> , в котором необходимо выбрать переменные для биннинга и указать необходимые параметры.<br>Параметр <b>Метод</b> задает метод, который будет использоваться для бинаризации переменной. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• <b>Однаковая ширина</b> (по умолчанию)</li> <li>• <b>Квантильный</b></li> <li>• <b>Дерево</b></li> </ul> Параметр <b>Количество бинов</b> имеет значение по умолчанию, равное 10. Для метода <b>Дерево</b> значение должно быть больше или равно 2. Для остальных методов – больше или равно 1. |

| Параметр  | Возможные значения и ограничения  | Описание  |
|---|---|---|
| <b>Минимальное количество наблюдений в бине</b> | Ручной ввод целочисленного значения больше ли равно 1<br>По умолчанию — 1   | Является дополнительным параметром для <b>биннинга на основе дерева</b> и задает минимальное количество наблюдений в бине   |
| <b>Критерий разбиения для регрессии</b>         | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• squared error (по умолчанию)</li> <li>• friedman mse</li> <li>• absolute error</li> <li>• poisson</li> </ul> | Является дополнительным параметром для <b>биннинга на основе дерева</b> и задает критерий разбиения для регрессии.<br>Предусмотрены следующие критерии: <ul style="list-style-type: none"> <li>• <b>squared error</b> (среднеквадратичная ошибка)</li> <li>• <b>friedman mse</b> (среднеквадратичная ошибка с оценкой улучшения Фридмана)</li> <li>• <b>absolute error</b> (средняя абсолютная ошибка)</li> <li>• <b>poisson</b> (отклонение Пуассона)</li> </ul>                                 |
| <b>Критерий разбиения для классификации</b>     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• gini (по умолчанию)</li> <li>• entropy</li> </ul>  | Является дополнительным параметром для <b>биннинга на основе дерева</b> и задает критерий разбиения для классификации.<br>Предусмотрены следующие критерии: <ul style="list-style-type: none"> <li>• <b>gini</b> (коэффициент Джини)</li> <li>• <b>entropy</b> (критерий прироста информации, энтропия)</li> </ul>  |
| <b>Префикс выходных переменных</b>              | Ручной ввод<br>По умолчанию — ENC_  | Данный параметр используется для кодирования категориальных переменных и задает префикс для выходных переменных узла  |
| <b>Обработка пропущенных значений</b>           | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Значение</li> <li>• Пропуск</li> </ul>   | Данный параметр используется для кодирования категориальных переменных и задает метод обработки пропущенных значений. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• <b>Пропуск</b> – будет поставлено NaN.</li> <li>• <b>Значение</b> – зависит от выбранного метода энкодинга. Для Count и Freq Encoding – пропуск как отдельная категория, Target Encoding – среднее значение целевой переменной на обучающей выборке (train), WOE Encoding – значение 0.</li> </ul> |
| <b>Обработка неизвестных значений</b>           | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Значение</li> <li>• Пропуск</li> </ul>   | Данный параметр используется для кодирования категориальных переменных и задает метод обработки неизвестных значений. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• <b>Пропуск</b> – будет поставлено NaN.</li> <li>• <b>Значение</b> – зависит от выбранного метода энкодинга. Для Count, Freq и WOE Encoding – значение 0, Target Encoding – среднее значение целевой переменной на обучающей выборке (train).</li> </ul>  |
| <b>Минимальная частота значений</b>             | Ручной ввод<br>Число больше или   | Данный параметр используется для кодирования категориальных переменных и  |

| <b>Параметр</b>                            | <b>Возможные значения и ограничения</b>                        | <b>Описание</b>   |
|--|--|---|
|  | равно 0 и меньше или равно 1<br>По умолчанию — 0               | задаёт минимальный % значений с категорией для того, чтобы не считать её редкой. Все редкие категории собираются в одну категорию и при скоринге для всех редких категорий будет проставлено одно и то же значение. Пропущенные значения не входят в редкие категории.  |
| <b>Количество итераций кросс-валидации</b> | Ручной ввод<br>Число больше или равно 0<br>По умолчанию — 0    | Данный параметр используется для кодирования категориальных переменных и задает количество итераций кросс-валидации. Кросс-валидация используется для борьбы с переобучением при использовании кодировщиков, основанных на целевой переменной (методы <b>Target Encoding</b> и <b>WOE</b> ).  |
| <b>Переменные</b>                          | Кнопка   | Используется для кодирования категориальных переменных<br>При выборе кнопки « <b>Переменные</b> » открывается окно <b>Переменные</b> , в котором необходимо выбрать переменные для энкодинга и указать необходимые параметры.<br>Параметр <b>Метод</b> задает метод, который будет использоваться для энкодинга.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• <b>Количество</b> (Count Encoding)</li> <li>• <b>Частота</b> (Freq Encoding)</li> <li>• <b>Целевая</b> (Target Encoding)</li> <li>• <b>WOE</b> (Weight of Evidence)</li> </ul> |
| <b>Сдвиг</b>                               | Ручной ввод<br>По умолчанию — 1                                | Параметр для <b>Target Encoding</b><br>Данный параметр используется для кодирования категориальных переменных и задает сдвиг  |
| <b>Масштаб</b>                             | Ручной ввод<br>Число больше 0<br>По умолчанию — 0,001          | Параметр для <b>Target Encoding</b><br>Данный параметр используется для кодирования категориальных переменных и задает масштаб  |
| <b>Рандомизация</b>                        | Чекбокс  | Параметр для <b>WOE Encoding</b><br>Выбор чекбокса указывает на необходимость рандомизации  |
| <b>Стандартное отклонение</b>              | Ручной ввод<br>Число больше или равно 0<br>По умолчанию — 0,05 | Параметр для <b>WOE Encoding</b><br>Данный параметр используется для кодирования категориальных переменных и задает стандартное отклонение  |
| <b>Регуляризация</b>                       | Ручной ввод<br>Число больше или равно 0<br>По умолчанию — 1    | Параметр для <b>WOE encoding</b><br>Данный параметр используется для кодирования категориальных переменных и задает значение регуляризации  |

**Таблица 18 Параметры узла «Биннинг»**

## Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

| Количество дней после регистрации клиента | Общее количество коммуникаций | Число дней с предыдущей коммуникацией | Наличие предыдущих коммуникаций | Результат предыдущей коммуникации | Отклик на предложение | BIN_Возраст клиента | ENC_Формат магазина |
|---|-------------------------------|---------------------------------------|---------------------------------|-----------------------------------|-----------------------|---------------------|---------------------|
| 1042                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     | 5                   | 0.447               |
| 1467                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     | 5                   | 0.447               |
| 1389                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     | 3                   | 0.447               |
| 579                                       | 1                             | -1                                    | 0                               | неизвестно                        | 1                     | 5                   | 0.447               |
| 673                                       | 2                             | -1                                    | 0                               | неизвестно                        | 1                     | 5                   | 0.541               |
| 562                                       | 2                             | -1                                    | 0                               | неизвестно                        | 1                     | 3                   | 0.541               |
| 1201                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     | 5                   | 0.541               |
| 1030                                      | 1                             | -1                                    | 0                               | неизвестно                        | 1                     | 5                   | 0.447               |

Рисунок 100 Таблица с примером данных

В результате выполнения узла в наборе данных будут рассчитаны новые переменные с указанными в параметрах префиксами и ролью Предиктор. Переменная, полученная в ходе биннинга, будет иметь тип Nominal, в ходе кодирования – Interval.

- Таблица со статистиками биннинга.

| Статистика разбиения Возраст клиента |                        |           |
|--------------------------------------|------------------------|-----------|
| BIN_ID                               | BIN_INT                | Обучающая |
| -1                                   | NaN                    | 0         |
| 1                                    | (-inf, 31.000000]      | 2503      |
| 2                                    | (31.000000, 36.000000] | 2300      |
| 3                                    | (36.000000, 42.000000] | 1986      |
| 4                                    | (42.000000, 52.000000] | 2301      |
| 5                                    | (52.000000, +inf)      | 2072      |

Рисунок 101 Пример таблицы со статистиками бинаризации

- Таблица со статистиками кодирования.

| Формат магазина | Freq  | Rare level | ENC_Формат магазина |
|-----------------|-------|------------|---------------------|
| ТЦ              | 0.49  | 0          | 0.447               |
| магазин у дома  | 0.134 | 0          | 0.394               |
| неизвестно      | 0.044 | 0          | 0.507               |
| франчайзинг     | 0.33  | 0          | 0.541               |
| NaN             | 0     | 0          | 0.473               |
| Unknown levels  | 0     | 0          | 0.473               |

**Рисунок 102 Пример таблицы со статистиками кодирования**

### 3.2.5.6.10. Узел «Дисперсионный анализ»

**Узел «Дисперсионный анализ»** применяется для исследования влияния одной или нескольких качественных переменных (факторов) на одну зависимую количественную переменную.

В ходе дисперсионного анализа тестируется гипотеза о равенстве средних (целевой непрерывной переменной) в группах (которые задаются одной или несколькими категориальными переменными) с помощью сравнения (анализа) дисперсий.

Дисперсия случайной величины — мера разброса значений случайной величины относительно её математического ожидания. Разделение общей дисперсии на несколько источников, позволяет сравнить дисперсию, вызванную различием между группами, с дисперсией, вызванной внутригрупповой изменчивостью.

**ВАЖНО!** Для работы узла необходимо, чтобы целевая переменная была **интервального типа (Interval)**.

**Список параметров узла** представлен в таблице ниже.

| Параметр | Возможные значения и ограничения           | Описание   |
|----------|--|--|
| Название | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе |
| Описание | Ручной ввод<br>Ограничений на значение нет | Описание узла  |

**Таблица 19 Параметры узла «Дисперсионный анализ»**

### Результаты выполнения узла:

- Табличное представление результатов дисперсионного анализа. Основные понятия таблицы с результатами дисперсионного анализа:
  - **Intercept** — коэффициент β0.
  - **Grouping** — межгрупповая дисперсия.
  - **Residual** — внутригрупповая дисперсия.
  - **Сумма квадратов** — суммы квадратов отклонений.
  - **Степени свободы** — количество значений, которые могут свободно изменяться.
  - **F-статистика** — значение критерия Фишера
  - **p-value** — вероятность получить F-значение, равное или превышающее то значение, которое рассчитано по имеющимся выборочным данным (при условии, что нулевая гипотеза верна).

| Результаты дисперсионного анализа |                 |                 |              |         |
|-----------------------------------|-----------------|-----------------|--------------|---------|
|                                   | Сумма квадратов | Степени свободы | F-статистика | p-value |
| Intercept                         | 106.183         | 1               | 24.401       | 0       |
| Grouping                          | 92.922          | 3               | 7.117        | 0.009   |
| Residual                          | 39.164          | 9               | -            | -       |

Рисунок 103 Пример таблицы с результатами дисперсионного анализа

- Табличное представление средних значений зависимой переменной для каждой группы.

| Средние по группам |         |
|--------------------|---------|
| Type               | density |
| Газ                | 0.002   |
| Лантаноид          | 8.305   |
| Металл             | 8.944   |

Рисунок 104 Пример таблицы со средними значениями целевой переменной (density) по группам (Type)

### 3.2.5.6.11. Узел «Стандартизация»

**Узел «Стандартизация»** приводит признаки в разных единицах измерения и диапазонах значений к общей шкале.

**Стандартизация** – преобразование числовых наблюдений с целью приведения их к некоторой общей шкале. Необходимость стандартизации вызвана тем, что разные признаки из обучающего набора могут быть представлены в разных масштабах и изменяться в разных диапазонах, что влияет на выявление некорректных зависимостей моделью.

Предусмотрены следующие методы:

- **Стандартное отклонение (std)** – преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.
- **Диапазон (range)** – линейно преобразует значения переменных в диапазон [0, 1].

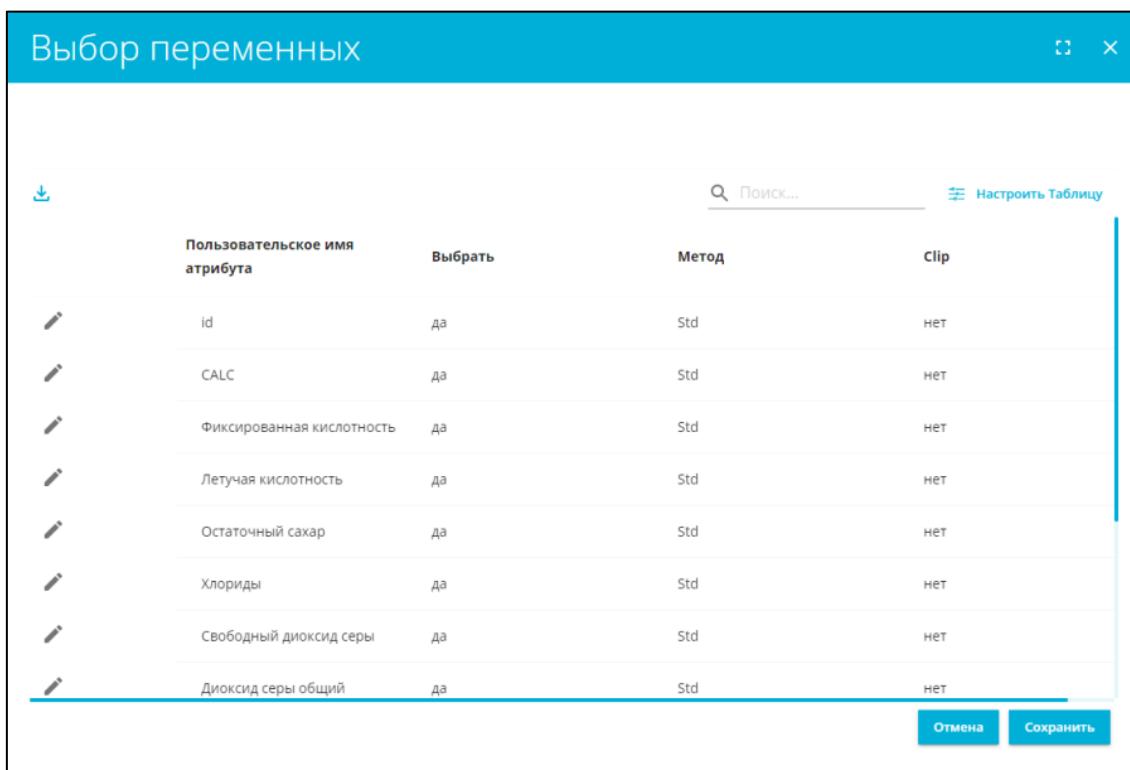
Список параметров узла представлен в таблице ниже.

| Параметр                    | Возможные значения и ограничения           | Описание   |
|-----------------------------|--|--|
| Название                    | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе   |
| Описание                    | Ручной ввод<br>Ограничений на значение нет | Описание узла  |
| Префикс выходных переменных | Ручной ввод                                | Данный параметр задает префикс стандартизованным выходным переменным   |
| Выбор переменных            | Кнопка                                     | При выборе кнопки <b>«Выбор переменных»</b> открывается окно <b>Выбор переменных</b> , в котором необходимо выбрать переменные для стандартизации. |

**Таблица 20 Параметры узла «Стандартизация»**

В окне **Выбор переменных** Пользователь имеет возможность выбрать атрибут для стандартизации и указать метод. Для этого необходимо:

- Рядом с интересующей переменной выбрать иконку
- В выпадающих меню выбрать **Метод** (std или range) и ограничить диапазон при стандартизации в [0, 1] на скоринге (столбец Clip). Если в столбце Clip стоит нет, то при стандартизации в [0, 1] на обучающей выборке всё точно будет в [0, 1], а на других выборках не обязательно.
- Сохранить изменения, выбрав иконку



**Рисунок 105 Окно Выбор переменных**

### Результаты выполнения узла:

- Таблица с примером полученных данных. Отображаются первые 100 наблюдений.

| STD_Диоксид серы общий | STD_Качество (баллы) | STD_Кислотность (pH) | STD_Летучая кислотность | STD_Остаточный сахар | STD_Плотность | STD_Свободный диоксид серы | STD_Содержание алкоголя |
|------------------------|----------------------|----------------------|-------------------------|----------------------|---------------|----------------------------|-------------------------|
| -0.38                  | -0.788               | 1.288                | 0.961                   | -0.454               | 0.558         | -0.467                     | -0.961                  |
| 0.624                  | -0.788               | -0.72                | 1.967                   | 0.043                | 0.028         | 0.872                      | -0.585                  |
| 0.229                  | -0.788               | -0.332               | 1.297                   | -0.17                | 0.134         | -0.084                     | -0.585                  |
| 0.411                  | 0.45                 | -0.98                | -1.385                  | -0.454               | 0.664         | 0.107                      | -0.585                  |
| -0.38                  | -0.788               | 1.288                | 0.961                   | -0.454               | 0.558         | -0.467                     | -0.961                  |
| -0.197                 | -0.788               | 1.288                | 0.738                   | -0.525               | 0.558         | -0.275                     | -0.961                  |
| 0.381                  | -0.788               | -0.073               | 0.403                   | -0.667               | -0.184        | -0.084                     | -0.961                  |
| -0.775                 | 1.689                | 0.511                | 0.682                   | -0.95                | -1.138        | -0.084                     | -0.398                  |
| -0.866                 | 1.689                | 0.316                | 0.291                   | -0.383               | 0.028         | -0.658                     | -0.867                  |
| 1.688                  | -0.788               | 0.251                | -0.156                  | 2.526                | 0.558         | 0.107                      | 0.072                   |
| 0.563                  | -0.788               | -0.202               | 0.291                   | -0.525               | -0.449        | -0.084                     | -1.148                  |
| 1.688                  | -0.788               | 0.251                | -0.156                  | 2.526                | 0.558         | 0.107                      | 0.072                   |

**Рисунок 106 Таблица с примером полученных данных**

В результате выполнения узла в наборе данных будут рассчитаны новые стандартизованные переменные (с префиксом **STD\_**).

### 3.2.5.6.12. Узел «Веса классов»

**Узел «Веса классов»** позволяет скорректировать дисбаланс классов при помощи задания весов.

В **задаче классификации** данные называются **несбалансированными**, когда в обучающей выборке доли объектов разных классов существенно различаются. Для такой ситуации существуют разные стратегии перебалансировки данных (например, замена большого класса подвыборкой по мощности равной малому классу, или **undersampling**, который предусмотрен в **узле «Sample»**).

В **узле «Веса классов»** Пользователь может задать веса для объектов каждого из классов вручную, либо сбалансировать автоматически.

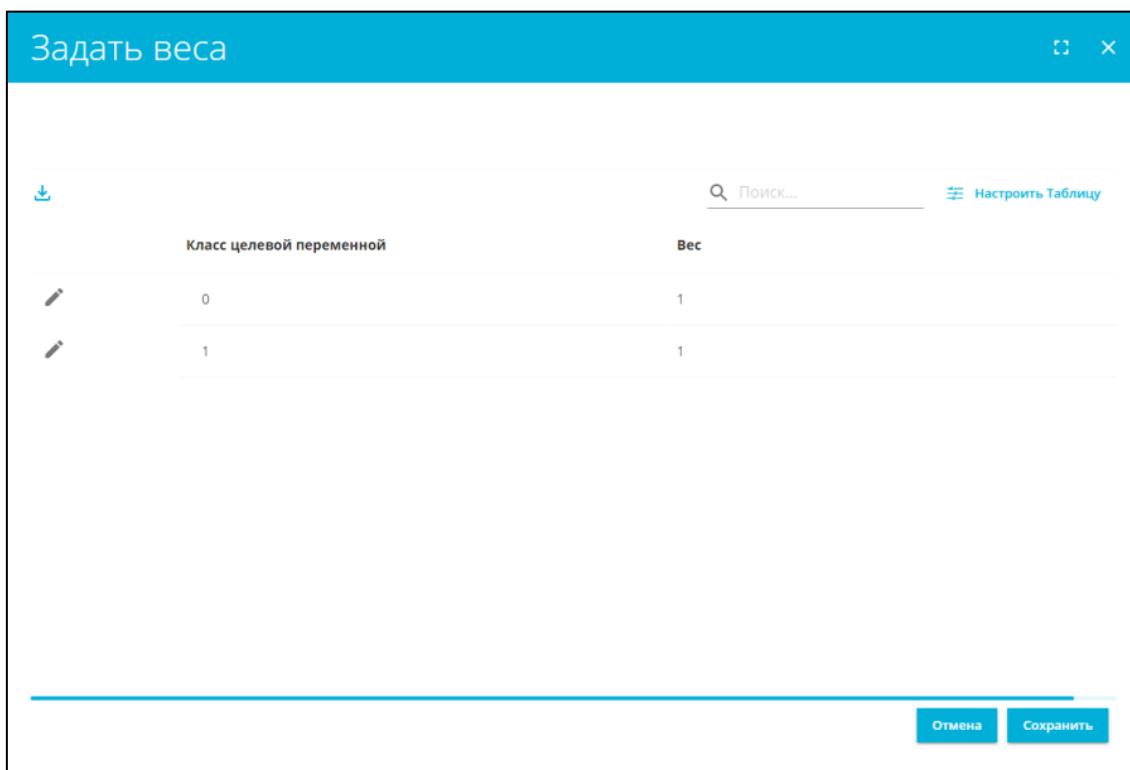
**Список параметров узла** представлен в таблице.

| Параметр                   | Возможные значения и ограничения           | Описание  |
|----------------------------|--|---|
| <b>Название</b>            | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>            | Ручной ввод<br>Ограничений на значение нет | Описание узла   |
| <b>Сбалансировать веса</b> | Чекбокс                                    | Выбор данного чекбокса указывает, что необходимо использовать автоматическую корректировку весов, которая считается обратно пропорционально частотам классов во входных данных. |
| <b>Задать веса</b>         | Кнопка                                     | При выборе кнопки открывается окно <b>Задать веса</b> , в котором необходимо задать вес каждому из классов.   |

**Таблица 21 Параметры узла «Веса классов»**

В окне **Задать веса** Пользователь может вручную задать веса каждому из классов. Для этого необходимо:

- Рядом с интересующей переменной выбрать иконку 
- В столбце **Вес** задать необходимое значение.
- Сохранить изменения, выбрав иконку 



**Рисунок 107 Окно Выбор переменных**

Для работы узла достаточно выбрать чекбокс **Сбалансировать веса** и запустить расчет узла. Весы классов посчитываются автоматически (без необходимости указывать их вручную в окне **Задать веса**).

#### Результаты выполнения узла:

- Таблица с примером полученных данных. Отображаются первые 100 наблюдений.

| Содержание алкоголя | Качество (0\1) | Уровень алкоголя | Уровень pH         | _class_weight_0 |
|---------------------|----------------|------------------|--------------------|-----------------|
| 9.4                 | 0              | Низкий           | Слабокислая среда  | 0.578           |
| 9.8                 | 0              | Ниже среднего    | Сильнокислая среда | 0.578           |
| 9.8                 | 0              | Ниже среднего    | Сильнокислая среда | 0.578           |
| 9.8                 | 0              | Ниже среднего    | Сильнокислая среда | 0.578           |
| 9.4                 | 0              | Низкий           | Слабокислая среда  | 0.578           |
| 9.4                 | 0              | Низкий           | Слабокислая среда  | 0.578           |

**Рисунок 108 Таблица с примером полученных данных**

В результате выполнения узла в наборе данных будет рассчитана новая переменная с весами и ролью **ClassWeight** (имя переменной **\_class\_weight\_0**).

- Таблица с весами классов.

| Веса классов |       | Поиск... | Настроить таблицу |
|--------------|-------|----------|-------------------|
| Класс        | Вес   |          |                   |
| 0            | 0.578 |          |                   |
| 1            | 3.684 |          |                   |

**Рисунок 109 Таблица с весами классов**

### 3.2.5.6.13. Узел «Автоэнкодер (PyTorch)»

**Узел "Автоэнкодер (PyTorch)"** позволяет создать сжатое представление данных (т.е. сократить размерность) в скрытом слое за счет специальной архитектуры нейронной сети.

**Список параметров узла** представлен в таблице ниже.

| Параметр                               | Возможные значения и ограничения           | Описание   | Группа параметров |
|--|--|--|-------------------|
| <b>Название</b>                        | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе   | Общие параметры   |
| <b>Описание</b>                        | Ручной ввод<br>Ограничений на значение нет | Описание узла  | Общие параметры   |
| <b>Конфигурация слоев автоэнкодера</b> | Кнопка                                     | <p>При выборе данной кнопки откроется окно <b>Конфигурация слоев нейросети</b>, где можно задавать слои, количество нейронов в слоях и функции активации для каждого узла. Для этого необходимо выбрать кнопку <b>Добавить</b> и в появившемся списке выбрать необходимую функцию активации и настроить параметры.</p> <p>Предусмотрены:</p> <ul style="list-style-type: none"> <li>• <b>Функция активации ReLU (Rectified Linear Unit)</b> – <math>f(x) = \max(0, x)</math></li> <li>• <b>Функция активации CELU</b> – <math>f(x) = \max(0, x) + \min(0, \alpha * (\exp(x/\alpha) - 1))</math></li> </ul> <p>Дополнительно необходимо задать <math>\alpha</math></p> <ul style="list-style-type: none"> <li>• <b>Функция активации ELU</b> – <math>f(x) = \begin{cases} \alpha * (\exp(x) - 1) &amp; \text{for } x \leq 0 \\ x &amp; \text{for } x &gt; 0 \end{cases}</math></li> </ul> | Общие параметры   |

| Параметр              | Возможные значения и ограничения  | Описание   | Группа параметров |
|-----------------------|---|--|-------------------|
|                       |   | <p>Дополнительно необходимо задать <math>\alpha</math></p> <ul style="list-style-type: none"> <li>• <b>Функция активации Sigmoid –</b><br/> <math display="block">\sigma(x) = \frac{1}{1 + e^{-x}}</math></li> <li>• <b>Функция активации Softmax –</b><br/> <math display="block">\sigma(x_i) = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}} \text{ for } i = 1, 2, \dots, K</math></li> <li>• <b>Линейный слой –</b><br/>         линейная функция, результат пропорционален переданному аргументу<br/> <math>f(x) = x</math></li> </ul> <p>Дополнительно необходимо задать <b>Количество выходных переменных</b> и при необходимости выбрать чекбокс <b>Добавить константу</b></p> <ul style="list-style-type: none"> <li>• <b>Функция активации Logsigmoid –</b><br/> <math display="block">\text{LogSigmoid}(x) = \log(\frac{1}{1 + e^{-x}})</math></li> <li>• <b>Исключение (Dropout) –</b><br/>         Во время обучения случайным образом обнуляет некоторые элементы входного тензора с вероятностью <math>p</math>, используя выборки из распределения Бернулли.</li> </ul> <p>Дополнительно необходимо задать <b>Вероятность исключения</b></p> <ul style="list-style-type: none"> <li>• <b>Tanh –</b> функция гиперболического тангенса<br/> <math display="block">f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}}</math></li> </ul> <p>После настройки необходимой конфигурации слоев нейросети выбрать кнопку <b>Сохранить</b>.</p> |                   |
| <b>Функция потерь</b> | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>• MAE</li> <li>• MSE</li> <li>• Poisson loss</li> <li>• Negative log likelihood</li> <li>• Cross entropy</li> </ul> | <p>Данный параметр задает функцию потерь. Предусмотрены:</p> <ul style="list-style-type: none"> <li>• <b>MAE</b> (средняя абсолютная ошибка)</li> <li>• <b>MSE</b> (среднеквадратическая ошибка)</li> <li>• <b>Poisson loss</b></li> <li>• <b>Negative log likelihood</b> (отрицательное логарифмическое правдоподобие)</li> <li>• <b>Cross entropy</b> (перекрестная энтропия)</li> </ul>   | Общие параметры   |

| Параметр                    | Возможные значения и ограничения  | Описание  | Группа параметров |
|-----------------------------|---|---|-------------------|
| <b>Алгоритм оптимизации</b> | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>• SGD</li> <li>• Adam (по умолчанию)</li> <li>• Adadelta</li> <li>• Adamax</li> <li>• LBFGS</li> <li>• AdamW</li> <li>• ASGD</li> <li>• NAdam</li> <li>• RAdam</li> <li>• Rprop</li> <li>• RMSprop</li> </ul> | <p>Данный параметр задает метод оптимизации, который будет использоваться для обновления весов нейронов скрытых слоев нейронной сети. Предусмотрены следующие методы:</p> <ul style="list-style-type: none"> <li>• <b>SGD (stochastic gradient descent)</b> – стохастический градиентный спуск. Данный метод делает шаг постоянной величины в направлении, указанном градиентом в текущей точке</li> <li>• <b>Adam (adaptive moment estimation)</b> – адаптивная оценка момента. Данный метод сочетает в себе идею метода Momentum о накоплении градиента и идею методов Adadelta и RMSProp об экспоненциальном сглаживании информации о предыдущих значениях квадратов градиентов.</li> <li>• <b>Adadelta (adaptive learning rate)</b> – метод адаптивной скорости обучения</li> <li>• <b>Adamax</b> – модификация метода Adam, основанная на бесконечной норме (max)</li> <li>• <b>LBFGS (limited-memory BFGS)</b> – BFGS с ограниченной памятью.</li> <li>• <b>AdamW</b>. Данный метод основан на адаптивной оценке моментов первого и второго порядка с добавленным методом уменьшения весов</li> <li>• <b>ASGD (average stochastic gradient descent)</b> – усредненный стохастический градиентный спуск. Данный метод усредняет веса, вычисляемые на каждой итерации.</li> <li>• <b>Nadam (Nesterov-accelerated adaptive momentum)</b>. Данный метод представляет собой модификацию оптимизатора Adam с добавлением</li> </ul> | Общие параметры   |

| Параметр                          | Возможные значения и ограничения    | Описание   | Группа параметров |
|-----------------------------------|-------------------------------------|--|-------------------|
|                                   |                                     | <p>момента Нестерова при вычислении градиентов.</p> <ul style="list-style-type: none"> <li><b>RAdam (Rectified Adam).</b> Данный метод является модификацией Adam, более устойчивой к изменению значений скорости обучения.</li> <li><b>Rprop (resilient backpropagation) –</b> устойчивый алгоритм обратного распространения. Данный метод использует только знаки частных производных для подстройки весовых коэффициентов. Также Rprop поддерживает отдельные дельты для каждого веса и смещения и адаптирует эти дельты во время обучения.</li> <li><b>RMSPROP (root mean square propagation) –</b> среднеквадратичное распространение корня. Данный метод использует усредненный по истории квадрат градиента.</li> </ul> |                   |
| <b>Скорость обучения</b>          | Ручной ввод<br>По умолчанию - 0,001 | Данный параметр задает скорость обучения, которая управляет размером шага при обновлении весов.  | Общие параметры   |
| <b>Количество эпох</b>            | Ручной ввод<br>По умолчанию - 10    | Данный параметр задает сколько раз алгоритм обучения будет обрабатывать весь набор обучающих данных.   | Общие параметры   |
| <b>Размер пакета</b>              | Ручной ввод<br>По умолчанию - 128   | Данный параметр определяет количество выборок, которые необходимо обработать перед обновлением параметров модели.  | Общие параметры   |
| <b>Seed</b>                       | Ручной ввод<br>По умолчанию - 42    | Начальное числовое значение для генератора случайных чисел. Используется для воспроизведения результатов при повторном запуске узла.   | Общие параметры   |
| <b>L2 регуляризация</b>           | Ручной ввод<br>По умолчанию - 0,1   | Данный параметр задает значение L2-регуляризации   | Общие параметры   |
| <b>Доля валидационной выборки</b> | Ручной ввод<br>По умолчанию - 0,1   | Данный параметр задает долю валидационной выборки, которая будет отобрана из исходной тестовой   | Общие параметры   |
| <b>Количество итераций без</b>    | Ручной ввод<br>По умолчанию - 5     | Данный параметр задает количество эпох без улучшения,  | Общие параметры   |

| Параметр                                  | Возможные значения и ограничения   | Описание  | Группа параметров  |
|---|--|---|--|
| <b>существенного улучшения</b>            |  | после которых скорость обучения будет снижена.<br>Пример: если значение параметра = 2, то первые 2 эпохи без улучшений loss будут проигнорированы, и только после 3-й эпохи скорость обучения уменьшится.   |  |
| <b>Режим динамического расчета порога</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"><li>• <b>rel</b></li><li>• <b>abs</b></li></ul>                                      | Данный параметр задает режим расчета порога<br>Предусмотрены: <ul style="list-style-type: none"><li>• <b>rel</b> – <math>dynamic\_threshold = best * (1 - threshold)</math></li><li>• <b>abs</b> – <math>dynamic\_threshold = best - threshold</math></li></ul>   | Общие параметры  |
| <b>Порог</b>                              | Ручной ввод<br>По умолчанию - 0,0001   | Данный параметр задает порог расчета нового оптимума, чтобы сосредоточиться только на значительных изменениях   | Общие параметры  |
| <b>Стандартизация</b>                     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"><li>• <b>no</b></li><li>• <b>std</b> (по умолчанию)</li><li>• <b>range</b></li></ul> | Данный параметр отвечает за выбор метода стандартизации числовых переменных.<br><b>Стандартизация</b> – преобразование числовых наблюдений с целью приведения их к некоторой общей шкале.<br>Необходимость стандартизации вызвана тем, что разные признаки из обучающего набора могут быть представлены в разных масштабах и изменяться в разных диапазонах, что влияет на выявление некорректных зависимостей моделью.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"><li>• <b>no</b> — стандартизация не нужна</li><li>• <b>std</b> — стандартное отклонение - преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li><li>• <b>range</b> — диапазон - линейно преобразует значения переменных в диапазон [0, 1].</li></ul> | Общие параметры  |
| <b>Beta1</b>                              | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,9   | Данный параметр задает коэффициент, используемый для управления скоростью затухания скользящих средних значений градиента (первого момента)   | Параметры алгоритма оптимизации <b>Adam</b> , <b>Adamax</b> , <b>AdamW</b> , |

| <b>Параметр</b>                      | <b>Возможные значения и ограничения</b>                      | <b>Описание</b>   | <b>Группа параметров</b>   |
|--------------------------------------|--|---|--|
|                                      |  |   | <b>RAdam, Nadam</b>  |
| <b>Beta2</b>                         | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,999 | Данный параметр задает коэффициент, используемый для управления скоростью затухания средних значений вторых моментов градиентов (некентрированной дисперсии)  | Параметры алгоритма оптимизации <b>Adam, Adamax, AdamW, RAdam, Nadam</b> |
| <b>Epsilon</b>                       | Ручной ввод числа с плавающей точкой<br>По умолчанию - 1e-8  | Данный параметр задает значение, добавляемое к знаменателю для улучшения числовой стабильности Minimal decay applied to lr. If the difference between new and old lr is smaller than eps, the update is ignored. Default: 1e-8.                       | Параметры алгоритма оптимизации <b>Adam, Adamax, AdamW, RAdam, Nadam</b> |
| <b>Использовать алгоритм AMSGrad</b> | Чекбокс  | Выбор данного чекбокса указывает, что необходимо использовать вариант AMSGrad этого алгоритма.<br>Разница между AMSgrad и Adam заключается в рассчитанном векторе второго момента, который используется для обновления параметров.                    | Параметры алгоритма оптимизации <b>Adam, AdamW</b>                       |
| <b>Импульс (Momentum)</b>            | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0     | Данный параметр задает коэффициент импульса (запоминает скорость на предыдущем шаге и добавляет в указанное число раз меньшую величину на следующем шаге)   | Параметры алгоритма оптимизации <b>SGD, RMSProp</b>                      |
| <b>Dampening</b>                     | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0     | Данный параметр задает демпфирование импульса. Демпфирование гарантирует, что оптимизатор не сделает слишком больших шагов, что может произойти, если использовать только импульс. Чем выше градиент, тем больше демпфирование уменьшает размер шага. | Параметры алгоритма оптимизации <b>SGD</b>                               |
| <b>Момент Нестерова</b>              | Чекбокс  | Выбор данного чекбокса включает импульс Нестерова (использует производную не в текущей точке, а в следующей, если бы мы продолжали двигаться в этом же направлении без изменений)   | Параметры алгоритма оптимизации <b>SGD</b>                               |
| <b>Rho</b>                           | Ручной ввод<br>По умолчанию - 0,9                            | Данный параметр задает коэффициент, используемый для вычисления скользящего среднего квадратов градиентов   | Параметры алгоритма оптимизации <b>Adadelta</b>                          |

| <b>Параметр</b>                                  | <b>Возможные значения и ограничения</b>  | <b>Описание</b>  | <b>Группа параметров</b>                                 |
|--|--|--|--|
| <b>Epsilon</b>                                   | По умолчанию - 1e-8  | Данный сглаживающий параметр задает значение, предотвращающее деление на 0.                                      | Параметры алгоритма оптимизации <b>Adadelta, RMSProp</b> |
| <b>Максимум итераций за шаг оптимизации</b>      | Ручной ввод целочисленного значения<br>По умолчанию - 20   | Данный параметр задает максимальное число итераций за шаг оптимизации  | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Максимум вычислений за шаг оптимизации</b>    | Ручной ввод целочисленного значения<br>По умолчанию - 1  | Данный параметр задает максимальное число вычислений функции за шаг оптимизации                                  | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Tolerance grad</b>                            | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,00001                                     | Данный параметр задает допуск завершения при оптимальности первого порядка                                       | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Tolerance change</b>                          | Ручной ввод числа с плавающей точкой<br>По умолчанию - 1e-9  | Данный параметр задает допуск завершения при изменении значения/параметра функции                                | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Количество запоминаемых шагов оптимизации</b> | Ручной ввод целочисленного значения<br>По умолчанию - 100  | Данный параметр задает количество запоминаемых шагов оптимизации   | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Line search</b>                               | Список: <ul style="list-style-type: none"><li>• no (по умолчанию)</li><li>• strong Wolfe</li></ul> | Данный параметр задает метод линейного поиска  | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Lambda</b>                                    | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,0001                                      | Данный параметр задает затухание   | Параметры алгоритма оптимизации <b>ASGD</b>              |
| <b>Alpha</b>                                     | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,75  | Данный параметр задает мощность для обновления скорости обучения   | Параметры алгоритма оптимизации <b>ASGD</b>              |
| <b>t0</b>  | Ручной ввод числа с плавающей точкой   | Данный параметр задает точку, с которой начинается усреднение. Если требуемое количество итераций меньше данного | Параметры алгоритма оптимизации <b>ASGD</b>              |

| Параметр                                  | Возможные значения и ограничения                               | Описание  | Группа параметров                              |
|---|--|---|--|
|   | По умолчанию - 100000  | значения, то усреднение не произойдет.  |  |
| <b>Сокращение импульса</b>                | Ручной ввод<br>По умолчанию - 0,004                            | Данный параметр задает значение сокращения импульса   | Параметры алгоритма оптимизации <b>Nadam</b>   |
| <b>Коэффициент уменьшения (eta minus)</b> | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию - 0,5 | Данный параметр задает мультипликативный коэффициент уменьшения.<br>Если на текущем шаге частная производная по соответствующему весу поменяла свой знак, значит последнее изменение было большим, и алгоритм проскочил локальный минимум.<br>Следовательно, величину коррекции необходимо уменьшить на значение данного параметра и вернуть предыдущее значение весового коэффициента. | Параметры алгоритма оптимизации <b>Rprop</b>   |
| <b>Коэффициент увеличения (eta plus)</b>  | Ручной ввод<br>Число больше 1<br>По умолчанию - 1,2            | Данный параметр задает мультипликативный коэффициент увеличения.<br>Если на текущем шаге частная производная по соответствующему весу не поменяла свой знак, значит нужно увеличить величину коррекции на значение данного параметра для достижения более быстрой сходимости.   | Параметры алгоритма оптимизации <b>Rprop</b>   |
| <b>Минимальный размер шага</b>            | Ручной ввод<br>По умолчанию - 0,000001                         | Данный параметр задает минимальный размер шага.<br>Он необходим, чтобы не допустить слишком маленьких значений весов, ограничивает величину коррекции снизу.  | Параметры алгоритма оптимизации <b>Rprop</b>   |
| <b>Максимальный размер шага</b>           | Ручной ввод<br>По умолчанию - 50                               | Данный параметр задает максимальный размер шага.<br>Он необходим, чтобы не допустить слишком больших значений весов, ограничивает величину коррекции сверху.  | Параметры алгоритма оптимизации <b>Rprop</b>   |
| <b>Alpha</b>                              | Ручной ввод<br>По умолчанию - 0,99                             | Данный параметр задает константу сглаживания  | Параметры алгоритма оптимизации <b>RMSProp</b> |
| <b>Центрировать</b>                       | Чекбокс  | Выбор данного чекбокса указывает, что необходимо вычислить центрированный RMSProp, градиент которого нормализуется по оценке его дисперсии  | Параметры алгоритма оптимизации <b>RMSProp</b> |

Таблица 22 Параметры узла «Автоэнкодер (PyTorch)»

## Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

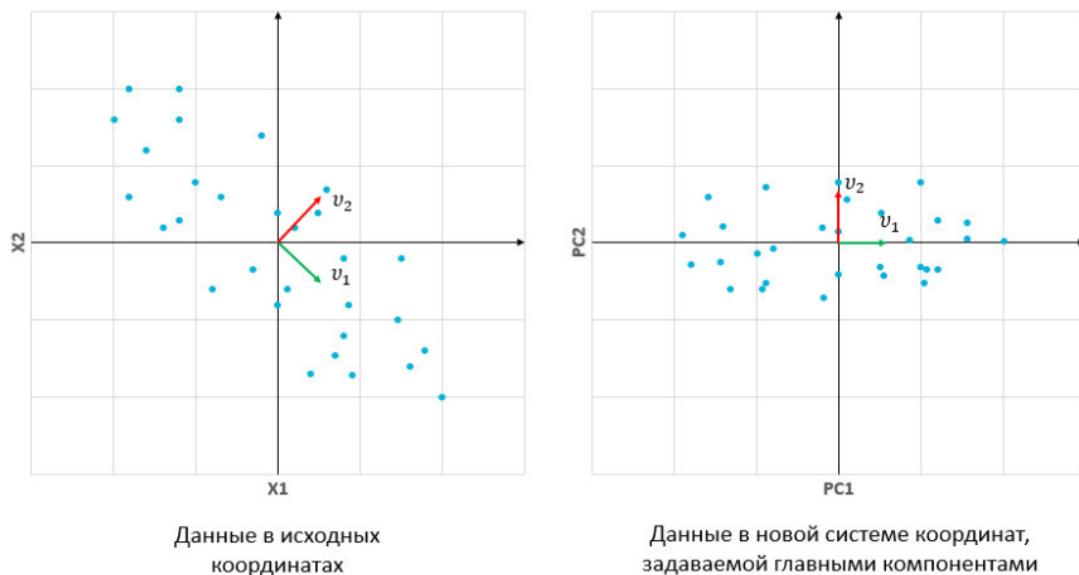
| Пример данных         |           |           |           |           |           |           |        |    |       |      |
|-----------------------|-----------|-----------|-----------|-----------|-----------|-----------|--------|----|-------|------|
| AUTOENC_0             | AUTOENC_1 | AUTOENC_2 | AUTOENC_3 | AUTOENC_4 | AUTOENC_5 | AUTOENC_6 | crim   | zn | indus | chas |
| 0.09308073162942929   | -0.042    | -0.193    | -0.055    | -0.117    | -0.027    | 0.078     | 0.158  | 0  | 10.81 | 0    |
| 0.07212171119944259   | -0.04     | -0.088    | -0.049    | -0.073    | -0.009    | 0.008     | 0.103  | 25 | 5.13  | 0    |
| -0.08740271555571276  | 0.094     | 0.041     | -0.06     | -0.039    | -0.091    | -0.047    | 0.349  | 0  | 9.9   | 0    |
| -0.56233789733710664  | 0.415     | 0.023     | -0.093    | -0.18     | -0.365    | 0.077     | 2.733  | 0  | 19.58 | 0    |
| -0.058821201727990206 | 0.067     | 0.003     | -0.177    | -0.076    | -0.093    | 0.022     | 0.043  | 21 | 5.64  | 0    |
| 0.05821181111786634   | 0.136     | 0.043     | -0.072    | -0.052    | -0.051    | 0.024     | 0.083  | 45 | 3.44  | 0    |
| 0.12470699752766681   | -0.049    | -0.148    | 0.035     | -0.055    | 0.006     | 0.028     | 0.19   | 22 | 5.86  | 0    |
| -0.09537271385771723  | 0.09      | 0.05      | -0.067    | -0.014    | -0.075    | -0.049    | 0.269  | 0  | 9.9   | 0    |
| -0.050994383251411576 | 0.048     | -0.059    | -0.099    | -0.108    | -0.045    | 0.036     | 10.062 | 0  | 18.1  | 0    |
| -0.36687899121522116  | 0.352     | -0.016    | -0.128    | -0.104    | -0.297    | 0.082     | 1.413  | 0  | 19.58 | 1    |
| -0.0553996995188863   | 0.197     | -0.128    | -0.147    | -0.128    | -0.114    | 0.016     | 25.84  | 0  | 18.1  | 0    |
| 0.03798152314908942   | 0.075     | -0.01     | -0.083    | -0.081    | -0.047    | 0.053     | 0.092  | 30 | 4.93  | 0    |

Рисунок 110 Таблица с примером данных

### 3.2.5.6.14. Узел «PCA»

**Анализ главных компонент (PCA – Principal component analysis)** – это метод, который преобразует большой набор переменных в меньший с минимальными потерями информативности. Использование PCA позволяет ускорить расчет модели за счет уменьшения количества входных переменных.

PCA используется для разложения многомерного набора данных на набор последовательных ортогональных компонентов, которые объясняют максимальную величину дисперсии (дисперсия – степень разброса данных). Снижение размерности достигается с помощью **SVD** (Singular Value Decomposition, сингулярное разложение).



В узле предусмотрены следующие **Алгоритмы главных компонент**:

- **auto** – алгоритм подбирается автоматически исходя из числа переменных и количества наблюдений в наборе данных: если входные данные больше  $500 \times 500$ , а количество извлекаемых компонентов меньше 80% наименьшего размера данных, то будет выбран метод **randomized**. В остальных случаях будет выбран **full**.
- **full** – точный полный SVD (стандартный решатель LAPACK) и отбор компонент с помощью постобработки (**параметр Методы определения количества компонент**)
- **arpack** – SVD, усеченный количеством компонент в пределе строго больше 0 и строго меньше  **$\min(\text{количество наблюдений}, \text{количество переменных})$**  (решатель ARPACK).
- **randomized** – алгоритм находит приближенное усеченное SVD, используя рандомизацию для ускорения вычислений (по методу Halko).

Рекомендуется нормализовать данные перед использованием PCA, иначе можно получить вводящие в заблуждение компоненты (**чекбокс Нормализация компонент**).

PCA работает с числовыми данными. Поэтому если в наборе данных присутствует категориальная переменная, то она кодируется методом One-hot encoding и количество рассчитанных главных компонент будет равно (**исходное кол-во некатегориальных атрибутов**) + (**столбцы One-hot encoding**).

**Пример:** изначально было 14 атрибутов один из которых категориальный (с 3 категориями), то в результатах получится 15 (подробнее про One-hot encoding в описании одноименного узла).

Для отбора компонент предусмотрены следующие **Методы определения количества компонент**:

- **full** – отбирается количество компонент в соответствии с формулой  **$\min(\text{количество наблюдений}, \text{количество переменных}) - 1$** .
- **number** – отбирамое количество компонент задается вручную.
- **variance** – будут отобраны компоненты, которые равны или превышают заданный процент объясненной дисперсии.
- **mle** – для отбора компонентов используется MLE Минка.

**Список параметров узла** представлен в таблице ниже.

| <b>Параметр</b>                                     | <b>Возможные значения и ограничения</b>   | <b>Описание</b>  |
|---|---|--|
| <b>Название</b>                                     | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                                     | Ручной ввод<br>Ограничений на значение нет  | Описание узла  |
| <b>Метод определения количества компонент</b>       | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• full (по умолчанию)</li> <li>• number</li> <li>• variance</li> <li>• mle</li> </ul>    | Данный параметр задает метод определения количества компонент. Предусмотрены: <ul style="list-style-type: none"> <li>• full</li> <li>• number</li> <li>• variance</li> <li>• mle</li> </ul>  |
| <b>Количество компонент</b>                         | Ручной ввод целочисленного значения<br>По умолчанию - 1   | Актуален при выборе <b>Метода определения количества компонент = number</b> .<br>Задает ограничение на количество компонент, которые в конечном итоге будут отобраны алгоритмом.   |
| <b>Процент объясненной дисперсии</b>                | Ручной ввод<br>Значение больше 0 и меньше или равно 1<br>По умолчанию – 0,95  | Актуален при выборе <b>Метода определения количества компонент = variance</b> .<br>Задает ограничение на процент объясненной дисперсии компонентами, которые в конечном итоге будут отобраны алгоритмом.   |
| <b>Нормализация компонент</b>                       | Чекбокс   | Данный чекбокс отвечает за преобразование числовых наблюдений с целью приведения их к общей шкале. Необходимость нормализации вызвана тем, что разные признаки из исходного набора могут быть представлены в разных масштабах и изменяться в разных диапазонах |
| <b>Алгоритм главных компонент</b>                   | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• auto (по умолчанию)</li> <li>• full</li> <li>• arpack</li> <li>• randomized</li> </ul> | Данный параметр задает алгоритм главных компонент. Предусмотрены следующие варианты: <ul style="list-style-type: none"> <li>• auto</li> <li>• full</li> <li>• arpack</li> <li>• randomized</li> </ul>  |
| <b>Погрешность</b>                                  | Ручной ввод числового значения<br>Значение должно быть в диапазоне [0.0, inf)<br>По умолчанию – 0   | Актуален при выборе <b>Алгоритма главных компонент = arpack</b><br>Задает допустимую погрешность для сингулярных значений.   |
| <b>Задать количество итераций степенного метода</b> | Чекбокс   | Актуален при выборе <b>Алгоритма главных компонент = randomized</b><br>Выбор данного чекбокса указывает, что   |

| Параметр  | Возможные значения и ограничения  | Описание  |
|---|---|---|
|   |   | необходимо задать количество итераций степенного метода   |
| <b>Количество итераций степенного метода</b>        | Ручной ввод числового значения<br>Значение должно быть в диапазоне [0, inf)<br>По умолчанию – 10  | Актуален при выборе <b>Алгоритма главных компонент = randomized</b><br>Перед увеличением <b>Количества итераций степенного метода</b> следует увеличивать <b>Количество дополнительных случайных векторов</b> , поскольку принцип рандомизированного метода заключается в том, чтобы избежать использования этих более дорогостоящих шагов итерации.  |
| <b>Количество дополнительных случайных векторов</b> | По умолчанию – 10   | Актуален при выборе <b>Алгоритма главных компонент = randomized</b><br>Задает дополнительное количество случайных векторов, что обеспечивает лучшую аппроксимацию сингулярных векторов и сингулярных значений.  |
| <b>Нормализация итераций</b>                        | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• auto (по умолчанию)</li> <li>• QR</li> <li>• LU</li> <li>• randomized</li> </ul> | Актуален при выборе <b>Алгоритма главных компонент = randomized</b><br>Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• auto – не применяет нормализацию, если <b>Количество итераций степенного метода &lt;=2</b>, и переключается на LU разложение в противном случае</li> <li>• QR – пошаговое разложение матрицы. Самое медленное, но наиболее точное</li> <li>• LU – LU разложение матрицы. Численно стабильное, но может терять в точности</li> <li>• randomized – наиболее быстрый метод, но нестабильный, если <b>Количество итераций степенного метода</b> большое (5 и больше)</li> </ul> |
| <b>Seed</b>   | Ручной ввод числового значения<br>Значение больше 0<br>По умолчанию – 42  | Начальное числовое значение для генератора случайных чисел.<br>Используется для воспроизведения результатов при повторном запуске   |
| <b>Исключить оригинальные предикторы</b>            | Чекбокс   | Выбор данного чекбокса указывает, что необходимо удалить оригинальные предикторы  |

Таблица 23 Параметры узла «PCA»

## Результаты выполнения узла:

- Таблица с примером данных. Отображаются первые 100 наблюдений.

| PC_1     | PC_2    | PC_3   | PC_4   | PC_5   | PC_6   | PC_7   | PC_8   | PC_9   | PC_10  | PC_11  |
|----------|---------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| -797.035 | 0.967   | 0.551  | -1.272 | -0.019 | 0.253  | 0.008  | 0.859  | 0.353  | -0.757 | -0.26  |
| -796.085 | 7.234   | 1.13   | 1.841  | -1.43  | -0.389 | 1.06   | 0.258  | -0.244 | -0.163 | -0.166 |
| -794.875 | -19.399 | 1.195  | -2.079 | -0.195 | 0.864  | -0.552 | -0.225 | 0.265  | 0.317  | 0.073  |
| -793.921 | -13.097 | 1.731  | -2.076 | -0.288 | 0.777  | -0.549 | -0.229 | 0.255  | 0.326  | 0.075  |
| -793.078 | 5.867   | -0.608 | -1.416 | -0.811 | 0.51   | -0.59  | -0.338 | 0.21   | -0.502 | 0.053  |
| -791.756 | -31.163 | 7.956  | -2.245 | -0.459 | -0.419 | -0.52  | 1.112  | -1.306 | 0.096  | -0.157 |
| -790.744 | -35.417 | 2.773  | -1.609 | 0.104  | 0.358  | -1.057 | -0.04  | -1.476 | 0.411  | 0.121  |
| -790.436 | 48.27   | -6.235 | -0.315 | 3.598  | 0.183  | 0.353  | 1.019  | 0.413  | 0.215  | 0.053  |
| -789.128 | 11.757  | -1.96  | -2.442 | -0.343 | 0.732  | -0.58  | -0.335 | 0.055  | -0.696 | -0.04  |
| -788.436 | 48.285  | -8.243 | -0.313 | 3.597  | 0.184  | 0.355  | 1.018  | 0.413  | 0.215  | 0.053  |
| -787.073 | 6.15    | 0.387  | -3.686 | 0.073  | 0.043  | 0.073  | 0.949  | 0.274  | -0.146 | -0.217 |

Рисунок 111 Таблица с примером данных

В данной таблице будут отображены вычисленные компоненты (переменные PC\_1, PC\_2 и т.д.).

- Столбчатая диаграмма, на которой отображены компоненты (их номера) и значения объясненной ими дисперсии (какая доля общего разброса точек приходится на каждую из новых координат).

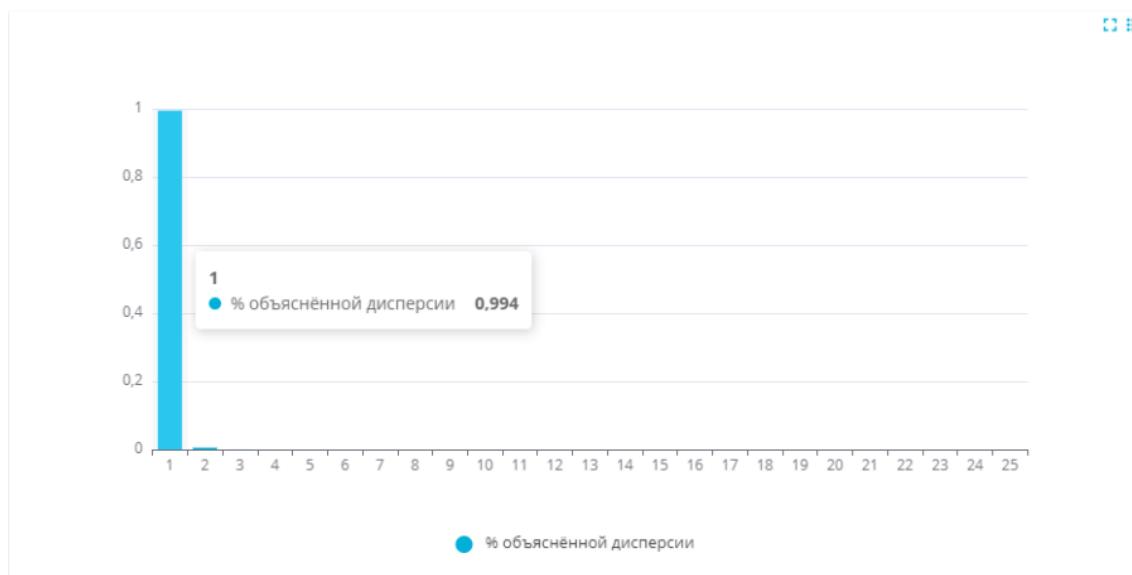
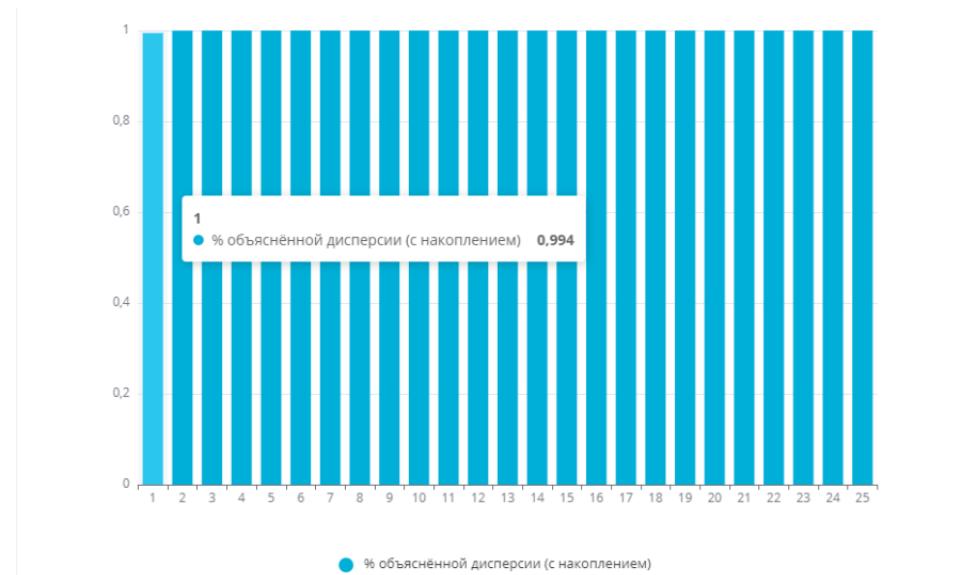


Рисунок 112 Пример столбчатой диаграммы

- Столбчатая диаграмма аналогичная первой, но с кумулятивной суммой.



**Рисунок 113 Пример столбчатой диаграммы с кумулятивной суммой**

### 3.2.5.6.15. Узел «Профилирование»

**Узел «Профилирование»** позволяет исследовать данные с целью выяснения статистических характеристик переменных.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения           | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе |
| <b>Описание</b> | Ручной ввод<br>Ограничений на значение нет | Описание узла  |

**Таблица 24 Параметры узла «Веса классов»**

**Результаты выполнения узла:**

- Таблица с примером данных.** Отображаются первые 100 наблюдений.

| Пример данных  |         |                  |                    |                           |                     |                  |                   |
|----------------|---------|------------------|--------------------|---------------------------|---------------------|------------------|-------------------|
| Уровень сахара | Хлориды | Уровень алкоголя | Уровень pH         | Фиксированная кислотность | Летучая кислотность | Лимонная кислота | Уровень плотности |
| Сухое          | 0.076   | Низкий           | Слабокислая среда  | 7.4                       | 0.7                 | 0                | Обычная           |
| Сухое          | 0.098   | Ниже среднего    | Сильнокислая среда | 7.8                       | 0.88                | 0                | Обычная           |
| Сухое          | 0.092   | Ниже среднего    | Сильнокислая среда | 7.8                       | 0.76                | 0                | Обычная           |
| Сухое          | 0.075   | Ниже среднего    | Сильнокислая среда | 11.2                      | 0.28                | 0                | Обычная           |
| Сухое          | 0.076   | Низкий           | Слабокислая среда  | 7.4                       | 0.7                 | 0                | Обычная           |
| Сухое          | 0.075   | Низкий           | Слабокислая среда  | 7.4                       | 0.66                | 0                | Обычная           |
| Сухое          | 0.069   | Низкий           | Сильнокислая среда | 7.9                       | 0.6                 | 0                | Обычная           |
| Сухое          | 0.065   | Ниже среднего    | Слабокислая среда  | 7.3                       | 0.65                | 0                | Обычная           |
| Сухое          | 0.073   | Низкий           | Слабокислая среда  | 7.8                       | 0.58                | 0                | Обычная           |

**Рисунок 114 Таблица с примером данных**

- Профиль данных.**

Профиль данных рассчитывается по всему набору данных и по каждой из переменных и зависит от ее типа (для категориальных и количественных переменных разный набор статистик).

Для того, чтобы отобразить статистики по нужной переменной, необходимо выбрать ее в списке.



**Рисунок 115 Пример профиля данных**

### Общие статистики

Для всех наблюдений набора данных считается следующий набор статистик (**пункт Статистика по набору данных**):

- Количество значений.
- Количество уникальных значений.
- Процент уникальных значений (в процентах).
- Количество дублирующих строк.
- Процент дублирующих строк (в процентах).

### Категориальные переменные

Для категориальных переменных рассчитывается набор статистик в соответствии с таблицей ниже.

| Статистика                             | Описание   |
|--|--|
| <b>Количество уникальных значений</b>  | Количество неповторяющихся значений                        |
| <b>Процент уникальных значений</b>     | Процент уникальных значений от общего количества значений  |
| <b>Количество пропущенных значений</b> | Количество пропущенных значений                            |
| <b>Процент пропущенных значений</b>    | Процент пропущенных значений от общего количества значений |
| <b>Максимальная длина</b>              | Количество знаков максимального по длине значения          |
| <b>Минимальная длина</b>               | Количество знаков минимального по длине значения           |
| <b>Top</b>                             | З значения с максимальной частотой                         |
| <b>Bottom</b>                          | З значения с минимальной частотой                          |

**Таблица 25 Набор статистик, рассчитываемых для категориальных переменных**

Для удобства интерпретации и анализа также предусмотрено построение **Облака слов**.

**Облако слов** (или облако тегов) визуализирует частоту появления значения переменной. Размер облака отражает частоту появления значения. Цветовая гамма не несет в себе смысла и выполняет исключительно эстетическую функцию.

Посмотреть **Облако слов** можно в том же контейнере с профилированием, выбрав в правом верхнем углу иконку  .



**Рисунок 116 Пример Облака слов**

### Количественные переменные

Для количественных переменных посчитывается набор статистик в соответствии с таблицей ниже.

| Статистика                             | Описание                                       |
|--|--|
| <b>Количество уникальных значений</b>  | количество неповторяющихся значений переменной |
| <b>Процент уникальных значений</b>     | процент неповторяющихся значений переменной    |
| <b>Количество пропущенных значений</b> | количество пропущенных значений переменной     |
| <b>Процент пропущенных значений</b>    | процент пропущенных значений                   |
| <b>Минимальное значение</b>            | наименьшее значение переменной                 |

| <b>Статистика</b>                      | <b>Описание</b>  |
|--|--|
| <b>Максимальное значение</b>           | наибольшее значение переменной   |
| <b>Среднее значение</b>                | сумма всех значений переменной, разделенная на число этих значений   |
| <b>5-я персентиль</b>                  | это некоторое значение $X$ из данного ряда, которое делит все имеющиеся в нем значения на две группы: 5% значений, которые меньше $X$ , и оставшиеся значения (то есть 95%), которые превышают $X$ .   |
| <b>95-я персентиль</b>                 | это некоторое значение $X$ из данного ряда, которое делит все имеющиеся в нем значения на две группы: 95% значений, которые меньше $X$ , и оставшиеся значения (то есть 5%), которые превышают $X$ .   |
| <b>1-я квартиль</b>                    | 25-я персентиль  |
| <b>3-я квартиль</b>                    | 75-я персентиль  |
| <b>Медиана</b>                         | значение, которое делит распределение пополам (его площадь в т.ч.): половина значений больше медианы, половина — не больше.  |
| <b>Межквартильный размах</b>           | разница между 3-м и 1-м квартилями   |
| <b>Количество выбросов</b>             | Выбросами считаются наблюдения, которые отклоняются от своего математического ожидания более чем на три среднеквадратических отклонения (правило трех сигм).   |
| <b>Коэффициент вариации</b>            | величина, равная отношению стандартного (среднеквадратичного) отклонения случайной величины к ее математическому ожиданию. Он применяется для сравнения вариативности одного и того же признака в нескольких совокупностях с различным средним арифметическим. Если значение коэффициента вариации не превышает 33%, то совокупность считается однородной, а если больше 33%, то — неоднородной.       |
| <b>Коэффициент эксцесса</b>            | характеризует меру высоты графика. Если коэффициент больше нуля, то распределение является более высоким («островершинным») относительно «эталонного» нормального распределения. Если коэффициент ниже нуля, то более низким и пологим.  |
| <b>Медианное абсолютное отклонение</b> | вычисляется как медиана абсолютного значения для каждого значения минус медианное значение группы. Является статистикой, более устойчивой к выбросам в наборе данных, чем стандартное отклонение   |
| <b>Асимметрия</b>                      | характеризует меру скошенности графика влево/вправо. Если коэффициент асимметрии отрицателен, то скос левосторонний. Если коэффициент положителен, то скос правосторонний. И чем коэффициент больше по модулю, тем сильнее скос распределения.   |
| <b>Стандартное отклонение</b>          | статистическая характеристика распределения случайной величины, показывающая среднюю степень разброса значений величины относительно математического ожидания. Большее значение среднеквадратического отклонения показывает больший разброс наблюдаемых значений признака относительно среднего; меньшее значение, соответственно, показывает, что величины в множестве сгруппированы вокруг среднего. |
| <b>Дисперсия</b>                       | величина, которая характеризует меру разброса значений случайной величины относительно ее математического ожидания.  |

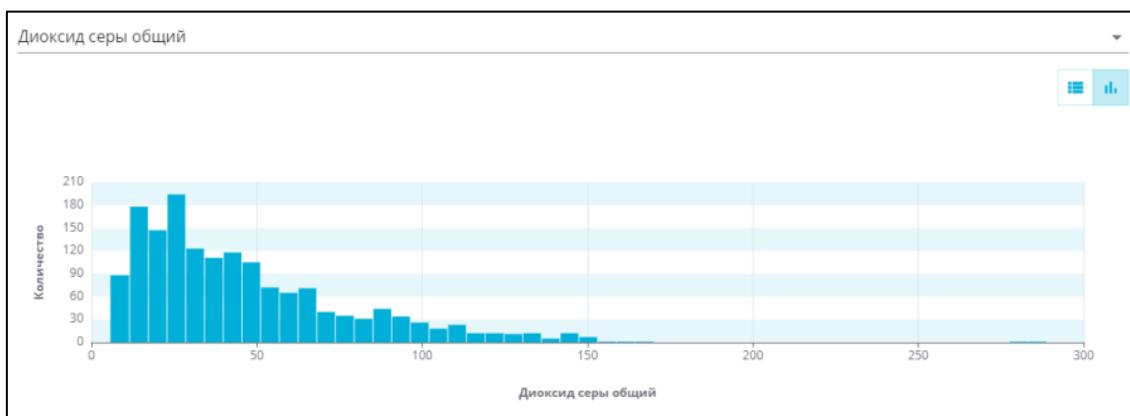
| Статистика          | Описание  |
|---------------------|---|
| <b>T-статистика</b> | T-статистика — это разница между средним по выборке и гипотетическим средним (предполагаемым равным нулю), деленная на расчетную стандартную ошибку среднего.   |
| <b>Пи-значение</b>  | Уровень значимости — вероятность получить Т-значение, равное или превышающее то значение, которое мы в действительности рассчитали по имеющимся выборочным данным (при условии, что нулевая гипотеза верна) |

**Таблица 26 Набор статистик, рассчитываемых для количественных переменных**

Для удобства интерпретации и анализа предусмотрено построение **Гистограммы**.

**Гистограмма** визуализирует распределение данных в рамках непрерывного интервала. На горизонтальной оси отмечаются интервалы (бины), а на вертикальной оси отмечается частота попаданий наблюдений в каждый интервал. Количество бинов не изменяется и по умолчанию равно 50.

Посмотреть гистограмму можно в том же контейнере, выбрав в правом верхнем углу иконку .



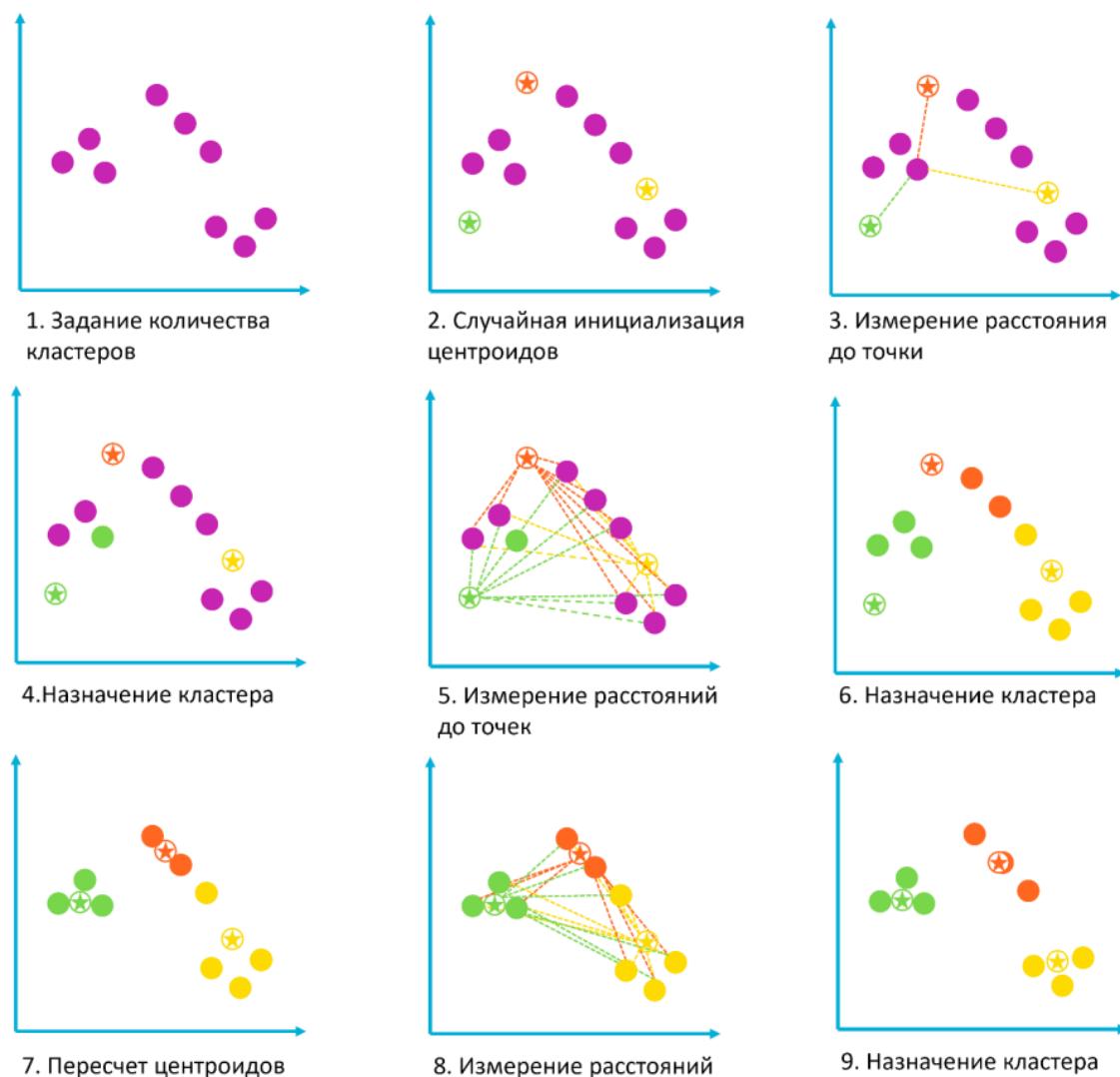
**Рисунок 117 Пример Гистограммы**

### 3.2.5.7. Группа узлов «Обучение без учителя»

#### 3.2.5.7.1. Узел «Кластерный анализ (k-means)»

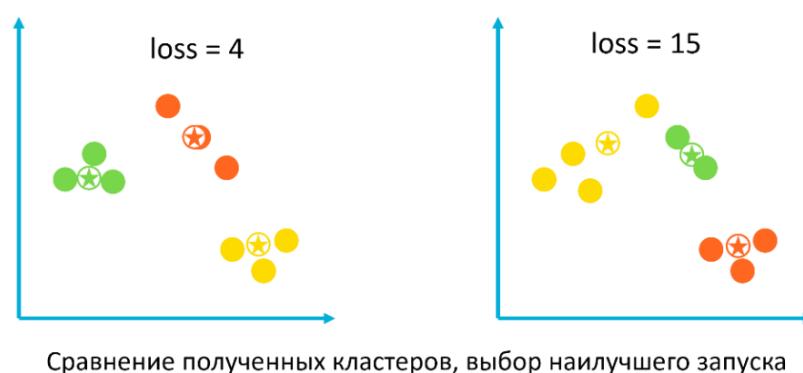
**Узел «Кластерный анализ (k-means)»** используется для кластеризации набора данных в отдельные группы (кластеры) исходя из выявленных шаблонов во входном наборе данных. Наблюдения группируются таким образом, чтобы они были похожи друг на друга внутри кластера, но различались с наблюдениями из других кластеров.

**Алгоритм работы:** Модель k-средних определяет начальный набор центроидов для кластеров (исходя из параметров **Количество кластеров** и **Метод инициализации кластеров**). Затем каждое наблюдение определяется в кластер с наиболее близким центроидом. Центроиды кластеров обновляются в соответствии с набором наблюдений, назначенным в каждый кластер. Далее итерационно проверяется необходимо ли переназначить наблюдение в другой кластер. Данный процесс продолжается до момента достижения максимального числа итераций (параметр **Максимальное количество итераций**).



**Рисунок 118 Принцип работы узла «Кластерный анализ (k-means)»**

Запуск алгоритма с определением начального набора центроидов происходит ограниченное количество раз (в соответствии с заданным параметром **Количество запусков**). После выполнения всех запусков выбирается запуск с минимальным критерием **инерции** (суммой квадратов расстояний между точками и центроидом внутри кластеров).



**Рисунок 119 Принцип работы узла «Кластерный анализ (k-means)»**

**Список параметров узла** представлен в таблице ниже.

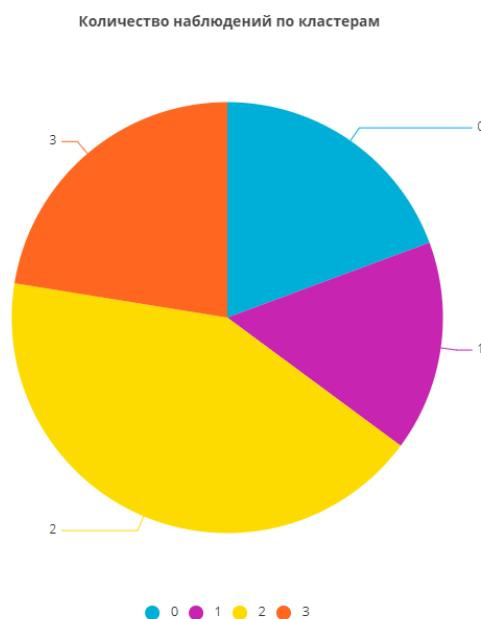
| Параметр                             | Возможные значения и ограничения   | Описание   |
|--------------------------------------|--|--|
| <b>Название</b>                      | Ручной ввод<br>Ограничений на значение нет   | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                      | Ручной ввод<br>Ограничений на значение нет   | Описание узла  |
| <b>Количество кластеров</b>          | Ручной ввод целочисленного значения<br>По умолчанию — 5  | Задание числа кластеров, на которые будет делиться пространство признаков.<br>Для определения количества кластеров можно воспользоваться априорной информацией об исходных данных в разделе Исследования данных при Кластеризации исходного набора данных.   |
| <b>Seed</b>                          | Ручной ввод целочисленного значения<br>По умолчанию — 42   | Начальное числовое значение для генератора случайных чисел.<br>Используется для воспроизведения результатов при повторном запуске узла   |
| <b>Метод инициализации кластеров</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• k-means++ (по умолчанию)</li> <li>• Forgy</li> </ul>                                | Данный параметр отвечает за выбор метода инициализации начальных точек кластеров. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• <b>k-means++</b><br/>Идея метода k-means++ состоит в том, чтобы случайным образом выбрать начальные точки, которые находятся как можно дальше друг от друга.</li> <li>• <b>Forgy</b><br/>Метод Forgy случайным образом выбирает <math>k</math> наблюдений (по числу заданных кластеров) из набора данных и использует их в качестве начальных значений.</li> </ul>  |
| <b>Стандартизация</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Нет</li> <li>• Стандартное отклонение (по умолчанию)</li> <li>• Диапазон</li> </ul> | Данный параметр отвечает за выбор метода стандартизации числовых переменных. <b>Стандартизация</b> – преобразование числовых наблюдений с целью приведения их к некоторой общей шкале. Необходимость стандартизации вызвана тем, что разные признаки из обучающего набора могут быть представлены в разных масштабах и изменяться в разных диапазонах, что влияет на выявление некорректных зависимостей моделью. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• <b>Нет.</b></li> <li>• <b>Стандартное отклонение</b> – преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li> <li>• <b>Диапазон</b> – линейно преобразует значения переменных в диапазон [0, 1].</li> </ul> |

| <b>Параметр</b>                          | <b>Возможные значения и ограничения</b>   | <b>Описание</b>  |
|--|---|--|
| <b>Количество запусков</b>               | Ручной ввод целочисленного значения<br>По умолчанию — 10  | Данный параметр задает число запусков алгоритма с разными начальными центроидами   |
| <b>Максимальное количество итераций</b>  | Ручной ввод целочисленного значения<br>По умолчанию — 300   | Данный параметр задает максимальное количество итераций в рамках одного запуска  |
| <b>Алгоритм K-средних</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• elkan</li> <li>• full</li> </ul> | Данный параметр отвечает за выбор алгоритма k-средних. <ul style="list-style-type: none"> <li>• <b>elkan</b> – может быть более эффективным для некоторых наборов данных с четко определенными кластерами за счет использования неравенства треугольника. Однако требует больше памяти</li> <li>• <b>full</b> – классический алгоритм</li> </ul> |
| <b>Размер выборки</b>                    | Ручной ввод целочисленного значения<br>По умолчанию — 1000  | Данный параметр задает размер выборки, которая будет отображена на графике <b>Силуэт</b> в результатах узла  |
| <b>Расстояние между кластерами</b>       | Ручной ввод целочисленного значения<br>По умолчанию — 25  | Данный параметр задает расстояние между кластерами на графике <b>Силуэт</b> в результатах узла   |
| <b>Количество бинов</b>                  | Ручной ввод целочисленного значения<br>По умолчанию — 10  | Данный параметр задает количество бинов, на которое будет делиться количественная переменная на <b>графике с параллельными осями</b> в результатах узла  |
| <b>Переменные, по которым делать оси</b> | Раскрывающийся список с выбором нескольких переменных   | Данный параметр задает переменные, которые будут отражены на <b>графике с параллельными осями</b> в результатах узла   |
| <b>Максимальное количество линий</b>     | Ручной ввод<br>По умолчанию — 50  | Данный параметр задает максимальное количество линий, которые будут отображаться на <b>графике с параллельными осями</b> в результатах узла  |

Таблица 27 Параметры узла «Кластерный анализ (k-means)»

### Результаты выполнения узла:

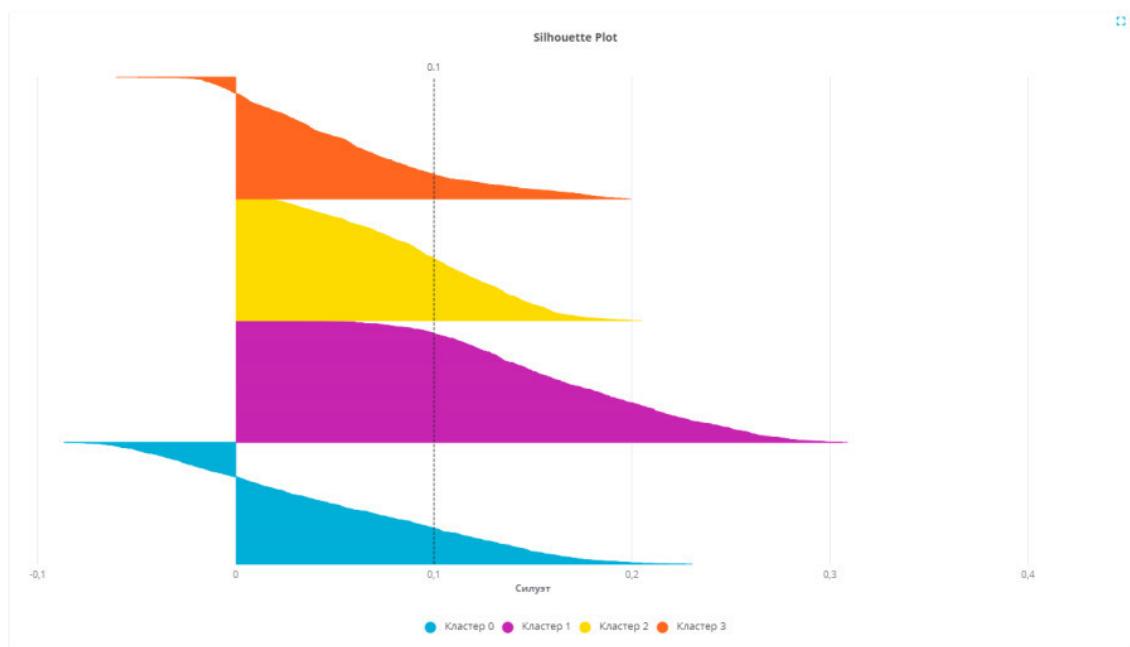
- Круговая диаграмма с количеством наблюдений по кластерам.



**Рисунок 120 Пример Круговой диаграммы с результатами кластеризации**

При наведении курсора мыши на сектор кластера можно узнать количество наблюдений в нем.

- Силуэт – Silhouette Plot



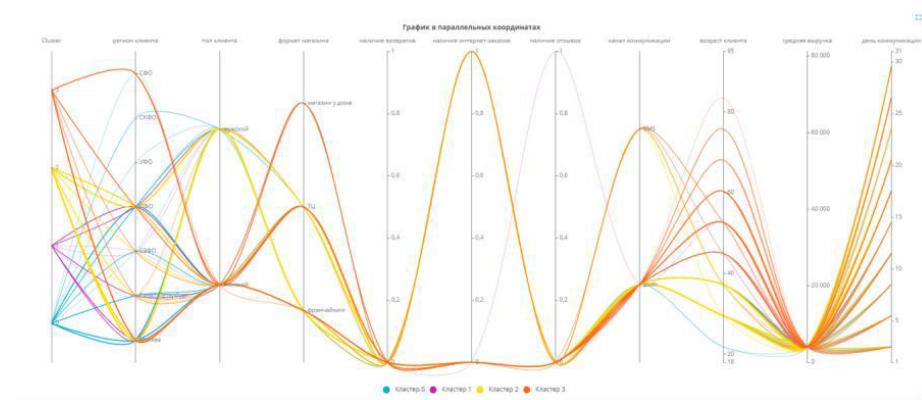
**Рисунок 121 Пример Silhouette Plot**

**Значение Silhouette** для каждого наблюдения (на графике отображается указанное в параметре **Размер выборки число наблюдений**) является мерой того, насколько это наблюдение похоже на наблюдения в собственном кластере по сравнению с наблюдениями в других кластерах.

**Значение Silhouette** находится в диапазоне от -1 до 1. Высокое значение указывает на то, что наблюдение хорошо соответствует собственному кластеру и плохо соответствует другим кластерам.

Если большинство наблюдений имеют низкое или отрицательное значение Silhouette, тогда пользователь должен перестроить кластеризацию с большим или меньшим количеством кластеров.

- График в параллельных координатах.



**Рисунок 122 Пример графика в параллельных координатах**

График в параллельных координатах позволяет интерпретировать построенные кластеры.

На данном графике каждой переменной присваивается собственная ось (согласно параметру **Переменные, по которым делать оси**). Оси располагаются параллельно друг другу, и каждая имеет свою собственную шкалу. Начальная ось отражает кластер, к которому модель отнесла наблюдение. Каждое наблюдение наносится на график в виде линии (параметр **Максимальное количество линий**), пересекающейся с каждой из осей. Таким образом, пользователь может выявить паттерны и корреляции между разными переменными.

- Таблица с примером данных. Отображаются первые 100 наблюдений.

Пример таблицы с данными, содержащей 100 строк и 8 столбцов. Столбец Cluster\_ID0 выделен красным квадратом.

| отклик на предложение | название акции механика          | тип акционной механики | время действия акции | минимальная сумма чека | применение на повторную покупку | Индикатор, возраста до 21 года | Cluster_ID0 |
|-----------------------|----------------------------------|------------------------|----------------------|------------------------|---------------------------------|--------------------------------|-------------|
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 1           |
| 0                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 1           |
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 1           |
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 1           |

**Рисунок 123 Таблица с примером данных**

В результате выполнения узла будет рассчитана новая переменная с результатами кластеризации (переменная **Cluster\_ID0**).

- Таблица с координатами центроидов, где в качестве строк выступают номера кластеров и значения переменных, в которых находятся центроиды этих кластеров.

| Координаты центроидов |                 |                 |                   |   |                               |                                       |                                 |                     |                       |                    |
|-----------------------|-----------------|-----------------|-------------------|---|-------------------------------|---------------------------------------|---------------------------------|---------------------|-----------------------|--------------------|
| Номер кластера        | возраст клиента | средняя выручка | день коммуникации | количество дней после регистрации клиента | общее количество коммуникаций | число дней с предыдущей коммуникацией | наличие предыдущих коммуникаций | регион клиента_ДВФО | регион клиента_Москва | регион клиента_ПФО |
| 0                     | -0.151          | 0.011           | -0.217            | -0.184                                    | -0.216                        | 1.651                                 | 1.252                           | 0.015               | 0.261                 | 0.246              |
| 1                     | -0.072          | -0.01           | 0.041             | 1.804                                     | 0.038                         | -0.338                                | -0.274                          | 0.022               | 0.214                 | 0.255              |
| 2                     | -0.554          | -0.161          | 0.074             | -0.407                                    | 0.1                           | -0.452                                | -0.34                           | 0.016               | 0.23                  | 0.286              |
| 3                     | 1.226           | 0.3             | 0.016             | -0.344                                    | -0.032                        | -0.333                                | -0.245                          | 0.047               | 0.212                 | 0.162              |

**Рисунок 124 Пример таблицы с координатами центроидов кластеров**

- Таблица со статистиками по кластерам
  - Номер кластера.
  - Количество наблюдений.
  - Среднеквадратичное расстояние между наблюдениями внутри кластера.
  - Сумма расстояний между наблюдениями внутри кластера.
  - Расстояние между центроидом и ближайшим наблюдением.
  - Расстояние между центроидом и наиболее удаленным наблюдением.
  - Расстояние между центроидом и вторым по удаленности наблюдением.
  - Расстояние между центроидом и третьим по удаленности наблюдением.
  - Ближайший кластер.
  - Расстояние до ближайшего центроида.
  - Среднее расстояние между центроидом и наблюдениями в кластере.
  - Сумма расстояний между наблюдениями и центроидом.

| Статистика по кластерам |                       |  |   |   |  |   |  |                   |                                    |  |
|-------------------------|-----------------------|--|---|---|--|---|--|-------------------|------------------------------------|--|
| Номер кластера          | Количество наблюдений | Среднеквадратичное расстояние между наблюдениями внутри кластера | Сумма расстояний между наблюдениями внутри кластера | Расстояние между центроидом и ближайшим наблюдением | Расстояние между центроидом и наиболее удаленным наблюдением | Расстояние между центроидом и вторым по удаленности наблюдением | Расстояние между центроидом и третьим по удаленности наблюдением | Ближайший кластер | Расстояние до ближайшего центроида | Среднее расстояние между центроидом и наблюдением в кластере |
| 0                       | 2159                  | 0.452  | 31329.011   | 2.58  | 23.921   | 23.115  | 17.376   | 2                 | 2.936                              | 3.628  |
| 1                       | 1763                  | 0.369  | 17094.239   | 1.919   | 14.004   | 11.143  | 10.311   | 2                 | 2.937                              | 3.01   |
| 2                       | 4734                  | 0.382  | 49193.842   | 2.208   | 22.497   | 15.117  | 15.105   | 3                 | 2.031                              | 3.124  |
| 3                       | 2506                  | 0.416  | 30826.603   | 2.271   | 25.005   | 24.893  | 20.09  | 2                 | 2.031                              | 3.298  |

**Рисунок 125 Пример со статистиками по кластерам**

- Таблица со статистиками по переменным кластера. По каждому кластеру отражены среднее и стандартное отклонение для каждой переменной.

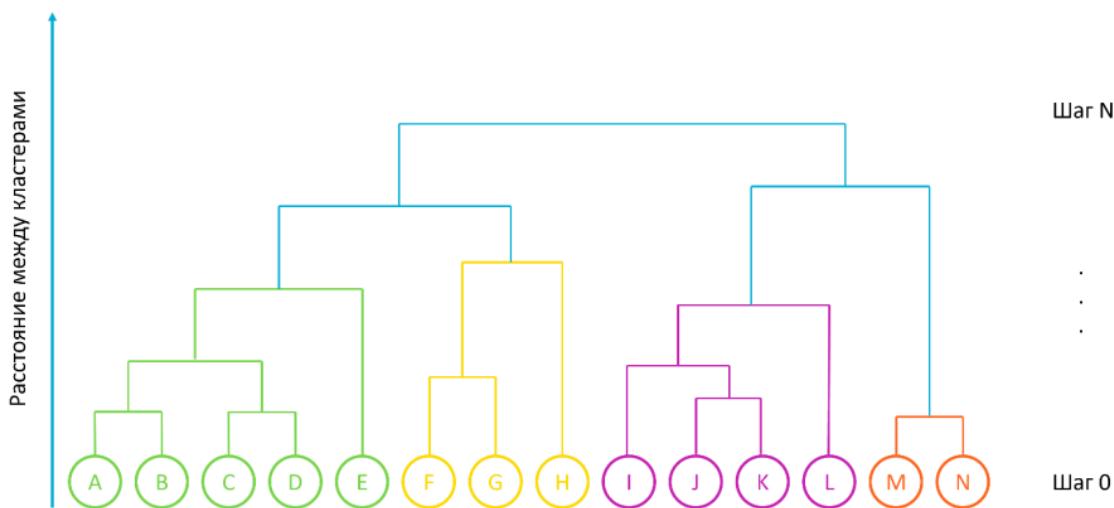
| Имя параметра                             | Номер кластера | Среднее  | Стандартное отклонение |
|---|----------------|----------|------------------------|
| возраст клиента                           | 0              | 39.438   | 10.331                 |
| средняя выручка                           | 0              | 1584.324 | 2679.87                |
| день коммуникации                         | 0              | 13.832   | 7.904                  |
| количество дней после регистрации клиента | 0              | 308.341  | 219.287                |
| общее количество коммуникаций             | 0              | 1.922    | 1.334                  |
| число дней с предыдущей коммуникацией     | 0              | 230.919  | 123.72                 |

**Рисунок 126 Пример таблицы со статистиками по переменным кластера**

### 3.2.5.7.2. Узел «Иерархическая кластеризация»

В основе **узла «Иерархическая кластеризация»** лежит алгоритм кластеризации, направленный на создание иерархии вложенных кластеров.

**Алгоритм работы:** Каждое наблюдение начинается в своем собственном кластере (Шаг 0), далее кластеры последовательно объединяются. Так, первоначально рассчитываются расстояния (расчет расстояния задает параметр **Метрика**) между наблюдениями, ближайшие из них объединяются в один кластер. Параметр **Критерий объединения** определяет стратегию слияния кластеров. Затем вычисляется расстояние между кластерами и ближайшие объединяются в один большой кластер. Слияние кластеров происходит до тех пор, пока не будет синтезирован один единый кластер (Шаг N).



**Рисунок 127 Схема работы алгоритма Иерархической кластеризации**

Для остановки алгоритма необходимо указать в параметре **Критерий остановки** требуемый вариант – по достижению заданного количества кластеров, либо по минимальному расстоянию между кластерами.

Для определения значения количества кластеров или расстояния между кластерами рекомендуется воспользоваться дендрограммой в результатах узла.

**Список параметров узла** представлен в таблице ниже.

| <b>Параметр</b>               | <b>Возможные значения и ограничения</b>   | <b>Описание</b>   |
|-------------------------------|---|---|
| <b>Название</b>               | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>               | Ручной ввод<br>Ограничений на значение нет  | Описание узла   |
| <b>Критерий остановки</b>     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Количество кластеров (по умолчанию)</li> <li>• Расстояние</li> </ul> | Данный параметр отвечает за выбор критерия остановки алгоритма.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• <b>Количество кластеров</b><br/>Достигнуто указанное в параметре <b>Количество кластеров</b> значение</li> <li>• <b>Расстояние</b><br/>Достигнуто указанное в параметре <b>Минимальное расстояние</b> значение.</li> </ul>  |
| <b>Количество кластеров</b>   | Ручной ввод целочисленного значения<br>Число больше или равно 1<br>По умолчанию — 5   | Данный параметр задает число кластеров, на которые будет делиться пространство признаков. Действителен при выбранном <b>Критерии остановки Количество кластеров</b>   |
| <b>Минимальное расстояние</b> | Ручной ввод<br>Число больше или равно 0<br>По умолчанию — 0   | Данный параметр задает минимальное расстояние между кластерами для остановки алгоритма. Действителен при выбранном <b>Критерии остановки Расстояние</b><br>Для определения значения минимального расстояния можно воспользоваться <b>Дендрограммой</b> в результатах узла.  |
| <b>Стандартизация</b>         | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Нет</li> <li>• Стандартное отклонение</li> <li>• Диапазон</li> </ul> | Данный параметр отвечает за выбор метода стандартизации числовых переменных.<br><b>Стандартизация</b> – преобразование числовых наблюдений с целью приведения их к некоторой общей шкале. Необходимость стандартизации вызвана тем, что разные признаки из обучающего набора могут быть представлены в разных масштабах и изменяться в разных диапазонах, что влияет на выявление некорректных зависимостей моделью.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• Нет.</li> <li>• Стандартное отклонение – преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li> <li>• Диапазон – линейно преобразует значения переменных в диапазон [0, 1].</li> </ul> |

| <b>Параметр</b>                          | <b>Возможные значения и ограничения</b>   | <b>Описание</b>   |
|--|---|---|
| <b>Метрика</b>                           | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Евклидова метрика</li> <li>• Манхэттенская метрика</li> <li>• Косинус</li> </ul>   | Данный параметр отвечает за выбор метрики, которая задает расчет расстояния между наблюдениями. Выбор метрики влияет на форму кластеров, поскольку некоторые элементы могут быть относительно ближе друг к другу по одной метрике, чем по другой.<br>Предусмотрены следующие метрики: <ul style="list-style-type: none"> <li>• Евклидова метрика</li> <li>• Манхэттенская метрика</li> <li>• Косинус</li> </ul>   |
| <b>Критерий объединения</b>              | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Ward (можно использовать только если в качестве метрики близости наблюдений выбрана Евклидова метрика)</li> <li>• Average</li> <li>• Maximum</li> <li>• Minimum</li> </ul> | Метрика, используемая для объединения кластеров.<br>Предусмотрены следующие метрики близости кластеров: <ul style="list-style-type: none"> <li>• Ward минимизирует сумму квадратов разностей во всех кластерах</li> <li>• Average минимизирует среднее расстояние между всеми наблюдениями пар кластеров.</li> <li>• Maximum сводит к минимуму максимальное расстояние между наблюдениями пар кластеров</li> <li>• Minimum минимизирует расстояние между ближайшими наблюдениями пар кластеров</li> </ul> |
| <b>Seed</b>                              | Ручной ввод целочисленного значения<br>По умолчанию — 42  | Начальное числовое значение для генератора случайных чисел. Используется для воспроизведения результатов при повторном запуске узла   |
| <b>Размер выборки</b>                    | Ручной ввод целочисленного значения<br>Значение больше или равно 2<br>По умолчанию — 1000   | Данный параметр задает размер выборки, которая будет отображена на графике <b>Силуэт</b> в результатах узла   |
| <b>Расстояние между кластерами</b>       | Ручной ввод<br>Значение больше или равно 0<br>По умолчанию — 25   | Данный параметр задает расстояние между кластерами на графике <b>Силуэт</b> в результатах узла  |
| <b>Количество бинов</b>                  | Ручной ввод целочисленного значения<br>Значение больше или равно 1<br>По умолчанию — 10   | Данный параметр задает количество бинов, на которое будет делиться количественная переменная на <b>графике с параллельными осями</b> в результатах узла   |
| <b>Переменные, по которым делать оси</b> | Раскрывающийся список с выбором нескольких переменных   | Данный параметр задает переменные, которые будут отражены на <b>графике с параллельными осями</b> в результатах узла  |

| Параметр                             | Возможные значения и ограничения  | Описание  |
|--------------------------------------|---|---|
| <b>Максимальное количество линий</b> | Ручной ввод целочисленного значения<br>Значение больше или равно 1<br>По умолчанию — 50 | Данный параметр задает максимальное количество линий, которые будут отражены на <b>графике с параллельными осями</b> в результатах узла |

Таблица 28 Список параметров узла "Иерархическая кластеризация"

### Параметры узла «Иерархическая кластеризация»

#### Результаты выполнения узла:

- Круговая диаграмма с количеством наблюдений по кластерам.

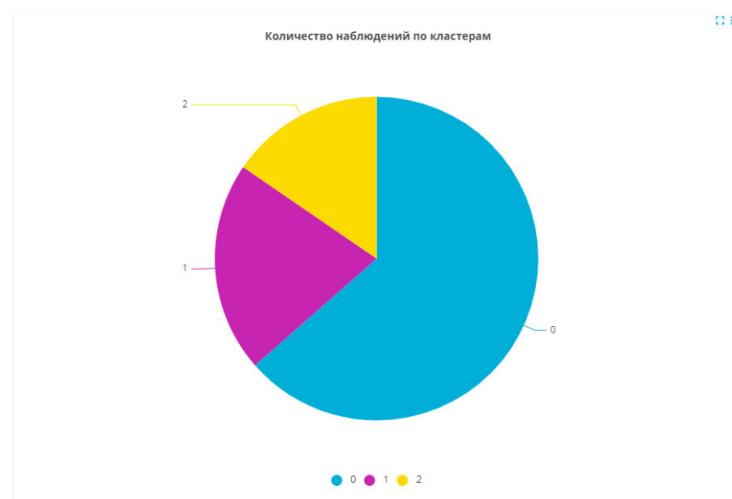


Рисунок 128 Пример круговой диаграммы

При наведении курсора мыши на сектор кластера можно узнать количество наблюдений в нем.

- Силуэт – Silhouette Plot.

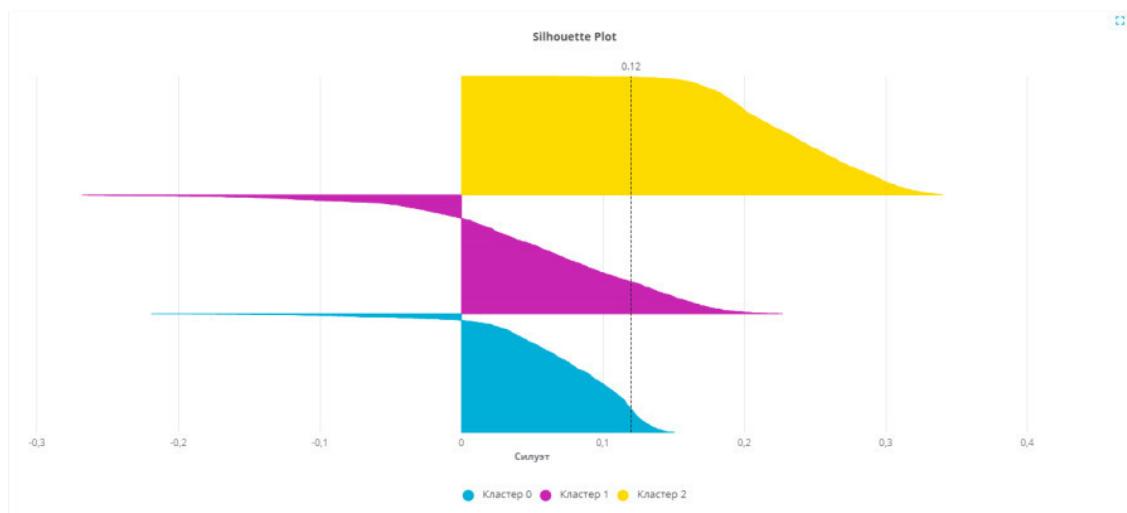


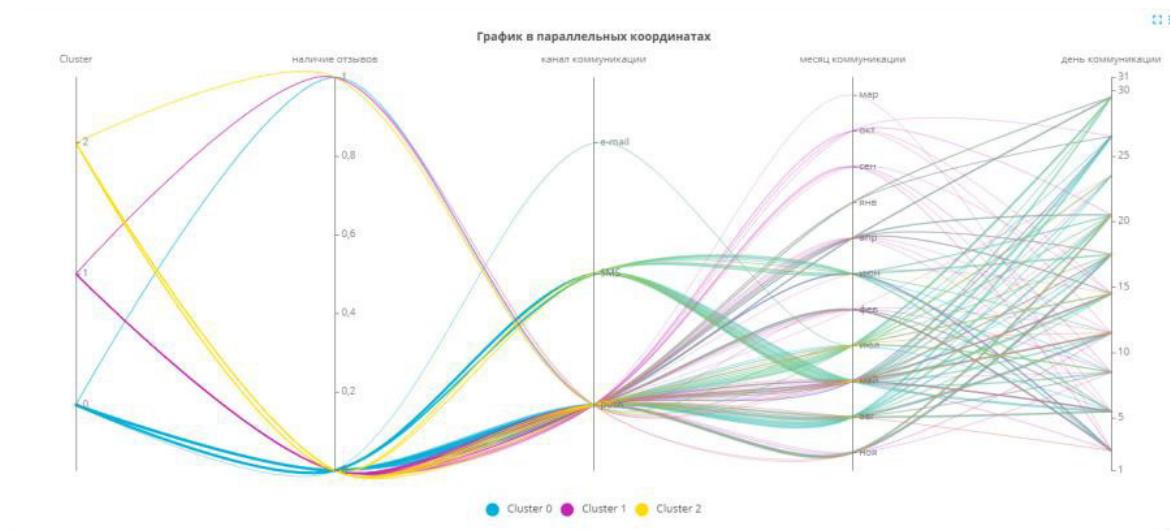
Рисунок 129 Пример Silhouette Plot

**Значение Silhouette** для каждого наблюдения (на графике отображается указанное в параметре **Размер выборки число наблюдений**) является мерой того, насколько это наблюдение похоже на наблюдения в собственном кластере по сравнению с наблюдениями в других кластерах.

**Значение Silhouette** находится в диапазоне от -1 до 1. Высокое значение указывает на то, что наблюдение хорошо соответствует собственному кластеру и плохо соответствует другим кластерам.

Если большинство наблюдений имеют низкое или отрицательное значение Silhouette, тогда пользователь должен перестроить кластеризацию с большим или меньшим количеством кластеров.

- График в параллельных координатах.

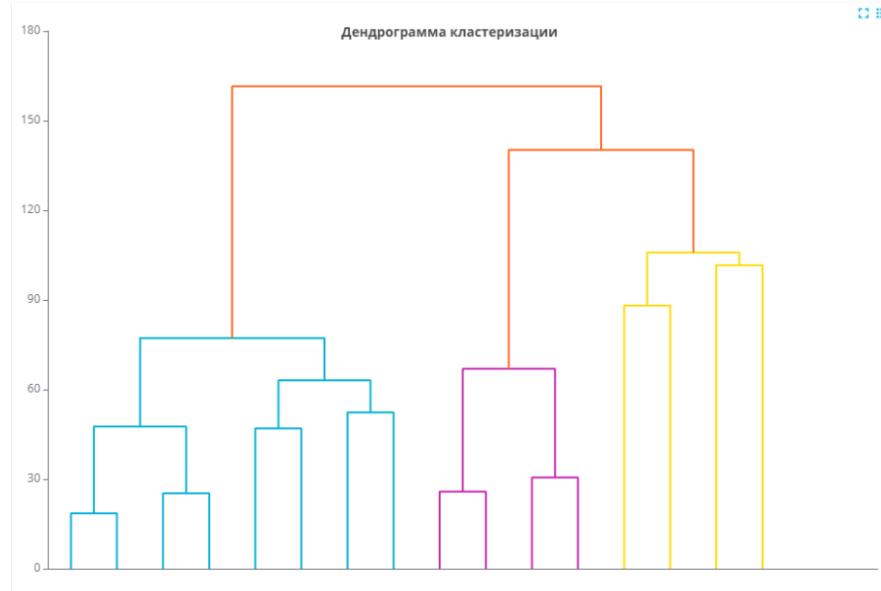


**Рисунок 130 Пример графика в параллельных координатах**

График в параллельных координатах позволяет интерпретировать построенные кластеры.

На данном графике каждой переменной присваивается собственная ось (согласно параметру **Переменные, по которым делать оси**). Оси располагаются параллельно друг другу, и каждая имеет свою собственную шкалу. Начальная ось отражает кластер, к которому модель отнесла наблюдение. Каждое наблюдение наносится на график в виде линии (параметр **Максимальное количество линий**), пересекающейся с каждой из осей. Таким образом, пользователь может выявить паттерны и корреляции между разными переменными.

- Дендрограмма кластеризации.



**Рисунок 131 Пример дендограммы**

Дендрограмма показывает близость отдельных наблюдений и кластеров, а также последовательность их объединения. Количество уровней соответствует количеству слияний кластеров. По оси Y расположена шкала, на которой откладывается расстояние между объектами в пространстве признаков.

- Таблица с примером данных. Отображаются первые 100 наблюдений.

| Пример данных         |                                  |                        |                      |                        |                                 |                                |             | Поиск... | Настроить таблицу |
|-----------------------|----------------------------------|------------------------|----------------------|------------------------|---------------------------------|--------------------------------|-------------|----------|-------------------|
| отклик на предложение | название акции механика          | тип акционной механики | время действия акции | минимальная сумма чека | применение на повторную покупку | Индикатор, возраста до 21 года | Cluster_ID0 |          |                   |
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 2           |          |                   |
| 0                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 2           |          |                   |
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 2           |          |                   |
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 2           |          |                   |
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 2           |          |                   |
| 1                     | Скидка 5% на покупку от 3000 руб | Прямая скидка          | 7 дней               | 3000 руб               | 0                               | 0                              | 2           |          |                   |

**Рисунок 132 Пример таблицы**

В результате выполнения узла будет рассчитана новая переменная с результатами кластеризации (переменная **Cluster\_ID0**).

- Таблица со статистиками по переменным кластера. По каждому кластеру отражены среднее и стандартное отклонение для каждой переменной.

| Имя параметра                             | Номер кластера | Среднее  | Стандартное отклонение |
|---|----------------|----------|------------------------|
| возраст клиента                           | 0              | 41.947   | 12.469                 |
| средняя выручка                           | 0              | 1721.314 | 3774.803               |
| день коммуникации                         | 0              | 16.065   | 8.508                  |
| количество дней после регистрации клиента | 0              | 254.918  | 193.309                |
| общее количество коммуникаций             | 0              | 2.753    | 3.178                  |
| число дней с предыдущей коммуникации      | 0              | 9.553    | 42.849                 |
| наличие предыдущих коммуникаций           | 0              | 0.188    | 0.855                  |

**Рисунок 133 Пример таблицы со статистиками по переменным кластера**

### 3.2.5.8. Группа узлов «Обучение с учителем»

#### 3.2.5.8.1. Узел «Дерево решений»

В основе **узла «Дерево решений»** лежит алгоритм, обобщающий наблюдения правилами вида «Если..., то...» в иерархическую, последовательную структуру в виде дерева. Правила генерируются в процессе обучения.

**Алгоритм работы:** Процесс построения деревьев решений представляет собой последовательное, рекурсивное разбиение множества наблюдений на подмножества с применением решающих правил в **узлах**. Разбиение продолжается до момента, пока не будет достигнуто условие остановки алгоритма. Последний узел, который не осуществляет проверку и разбиение, становится **листом**.

В основе алгоритма построения дерева решений лежит принцип жадной максимизации прироста информации – на каждом шаге выбирается тот признак, при разделении по которому прирост информации оказывается наибольшим. Дальше процедура повторяется рекурсивно, пока энтропия не окажется равной нулю или какой-то малой величине.

#### Борьба с переобучением:

- Ограничить максимальную глубину дерева (параметр **Максимальная глубина**)

- Ограничить минимальное число объектов листе (параметр **Минимальное количество наблюдений в листе**)
- Ограничить максимальное количества листьев в дереве (параметр **Максимальное количество листов**)

**Список параметров узла** представлен в таблице ниже.

| Параметр   | Возможные значения и ограничения  | Описание  |
|--|---|---|
| <b>Название</b>  | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>  | Ручной ввод<br>Ограничений на значение нет  | Описание узла   |
| <b>Критерий разбиения для классификации</b>            | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• gini (по умолчанию)</li> <li>• entropy</li> </ul>  | Данный параметр задает критерий разбиения на узлы для классификации. Предусмотрены следующие критерии: <ul style="list-style-type: none"> <li>• gini (неопределенность Джини) – направлен на максимизацию количества пар объектов, одного класса, оказавшихся в одном поддереве.</li> <li>• entropy (критерий прироста информации, энтропия) – направлен на максимизацию прироста информации</li> </ul> |
| <b>Критерий разбиения для регрессии</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• squared error (по умолчанию)</li> <li>• friedman mse</li> <li>• absolute error</li> <li>• poisson</li> </ul> | Данный параметр задает критерий разбиения для регрессионной задачи. Предусмотрены следующие критерии: <ul style="list-style-type: none"> <li>• squared error (среднеквадратичная ошибка)</li> <li>• friedman mse (среднеквадратичная ошибка с оценкой улучшения Фридмана)</li> <li>• absolute error (средняя абсолютная ошибка)</li> <li>• poisson (отклонение Пуассона)</li> </ul>                     |
| <b>Стратегия разбиения</b>                             | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• best (по умолчанию)</li> <li>• random</li> </ul>   | Данный параметр задает стратегию разделения на каждом узле. Предусмотрены следующие стратегии: <ul style="list-style-type: none"> <li>• best – выбор наилучшей функции сегментации и точки сегментации</li> <li>• random – случайное разделение</li> </ul>  |
| <b>Максимальная глубина</b>                            | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 5  | Данный параметр задает максимальную глубину дерева, после достижения которой алгоритм останавливает работу  |
| <b>Минимальное количество наблюдений для разбиения</b> | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 2  | Данный параметр задает минимальное количество наблюдений, которое должно быть в разбиении   |
| <b>Минимальное количество</b>                          | Ручной ввод<br>Неотрицательное  | Данный параметр задает минимальное количество наблюдений, которое может быть в листе  |

| <b>Параметр</b>  | <b>Возможные значения и ограничения</b>  | <b>Описание</b>  |
|--|--|--|
| <b>наблюдений в листе</b>  | число<br>По умолчанию — 5  |  |
| <b>Минимальная доля веса наблюдений в листе</b>  | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0   | Данный параметр определяет минимальный весовой коэффициент выборки в листовом узле. По умолчанию наблюдения имеют одинаковый вес   |
| <b>Максимальное количество признаков</b>   | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• all (по умолчанию)</li> <li>• sqrt</li> <li>• log2</li> <li>• number</li> <li>• frac</li> </ul> | Данный параметр определяет максимальное количество признаков, которое будет учитываться при поиске лучшего разделения. Предусмотрены следующие варианты: <ul style="list-style-type: none"> <li>• all – учитывать все доступные признаки</li> <li>• sqrt – учитывать sqrt(число всех доступных признаков)</li> <li>• log2 – учитывать log2(число всех доступных признаков)</li> <li>• number – учитывать указанное число признаков</li> <li>• frac – учитывать int(указанное число * число всех доступных признаков)</li> </ul> При выборе number или frac появится дополнительный параметр Число (вводится int) и Frac (вводится float) соответственно. |
| <b>Seed</b>  | Ручной ввод<br>По умолчанию — 12345  | Начальное числовое значение для генератора случайных чисел   |
| <b>Максимальное количество листов</b>  | Ручной ввод<br>Неотрицательное число   | Данный параметр определяет максимальное количество листов в дереве   |
| <b>Минимальное снижение неоднородности</b>   | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0   | Данный параметр определяет минимальное снижение неоднородности Узел будет разделен, если это разделение вызовет уменьшение неоднородности большее или равное указанному значению   |
| <b>ccp_alpha (отсечение с минимизацией стоимости-сложности)<br/>Метод отсечения дерева</b> | Ручной ввод<br>По умолчанию — 0  | Данный параметр регулирует количество отсекаемых узлов. Чем больше значение ccp_alpha, тем большее количество узлов удаляется из дерева  |

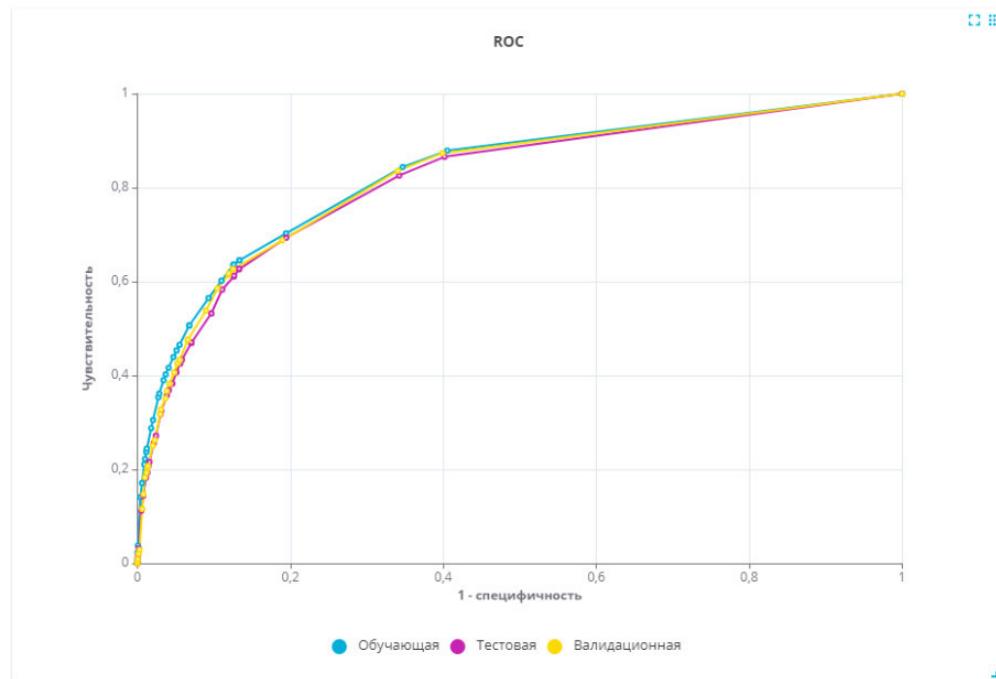
Таблица 29 Параметры узла «Дерево решений»

#### Результаты выполнения узла:

Узел «Дерево решений» имеет разные результаты в зависимости от решаемой задачи.

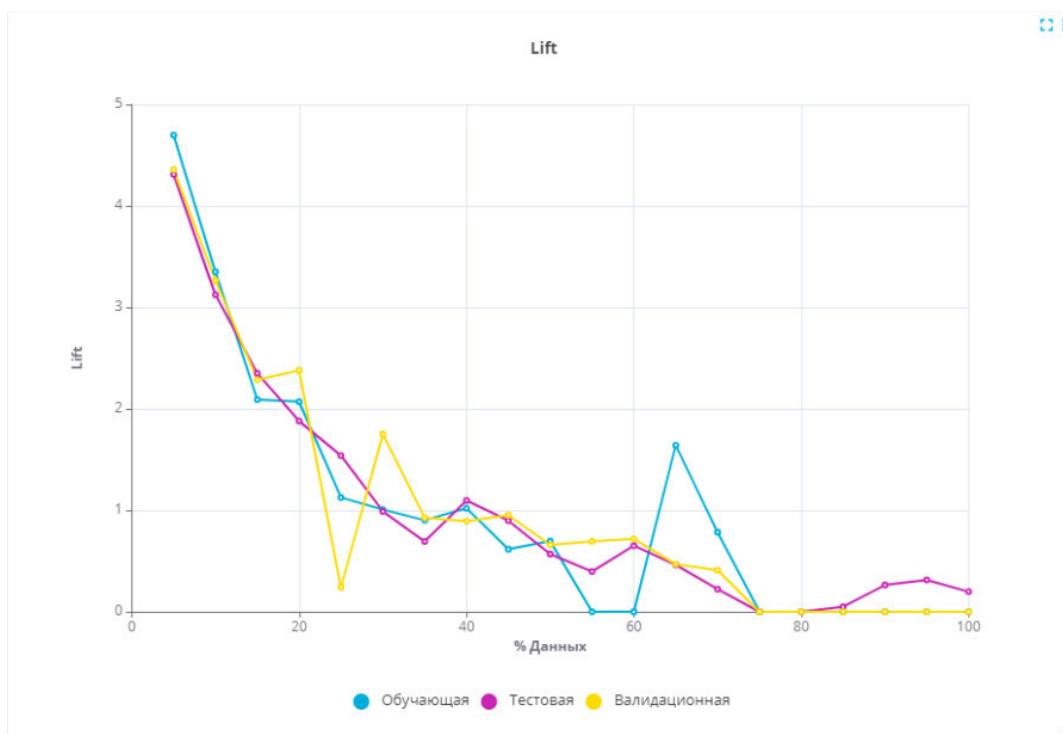
Результаты бинарной классификации представлены следующими объектами:

- График ROC .



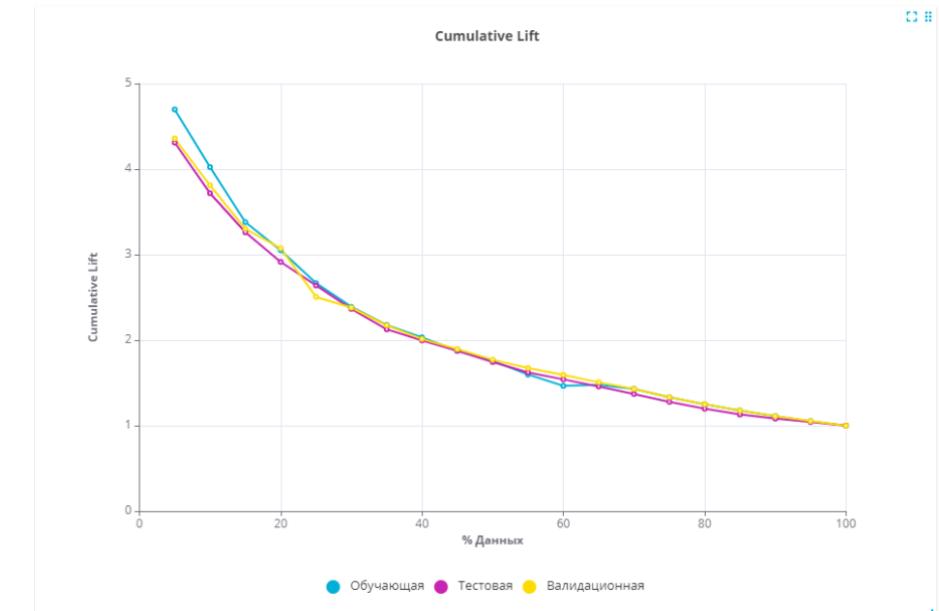
**Рисунок 134 Пример графика ROC**

- График Lift.



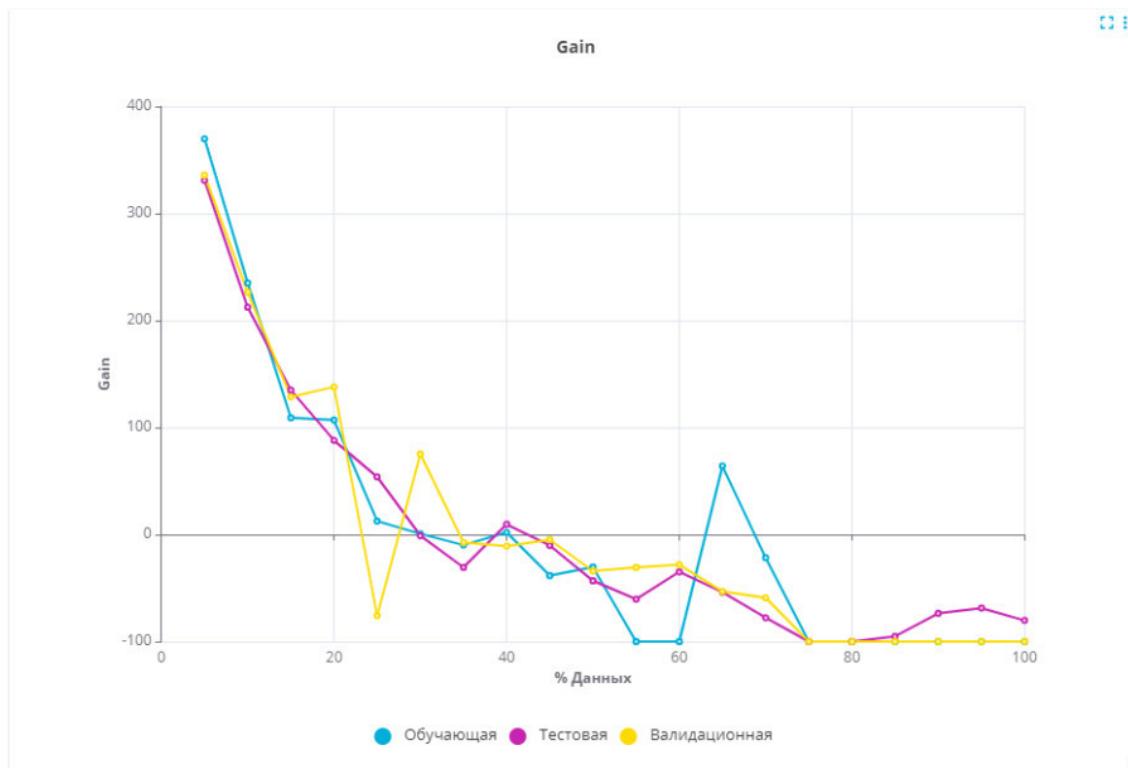
**Рисунок 135 Пример графика Lift**

- График Cumulative Lift.



**Рисунок 136 Пример графика Cumulative Lift**

- График Gain.



**Рисунок 137 Пример графика Gain**

- График Cumulative Gain.

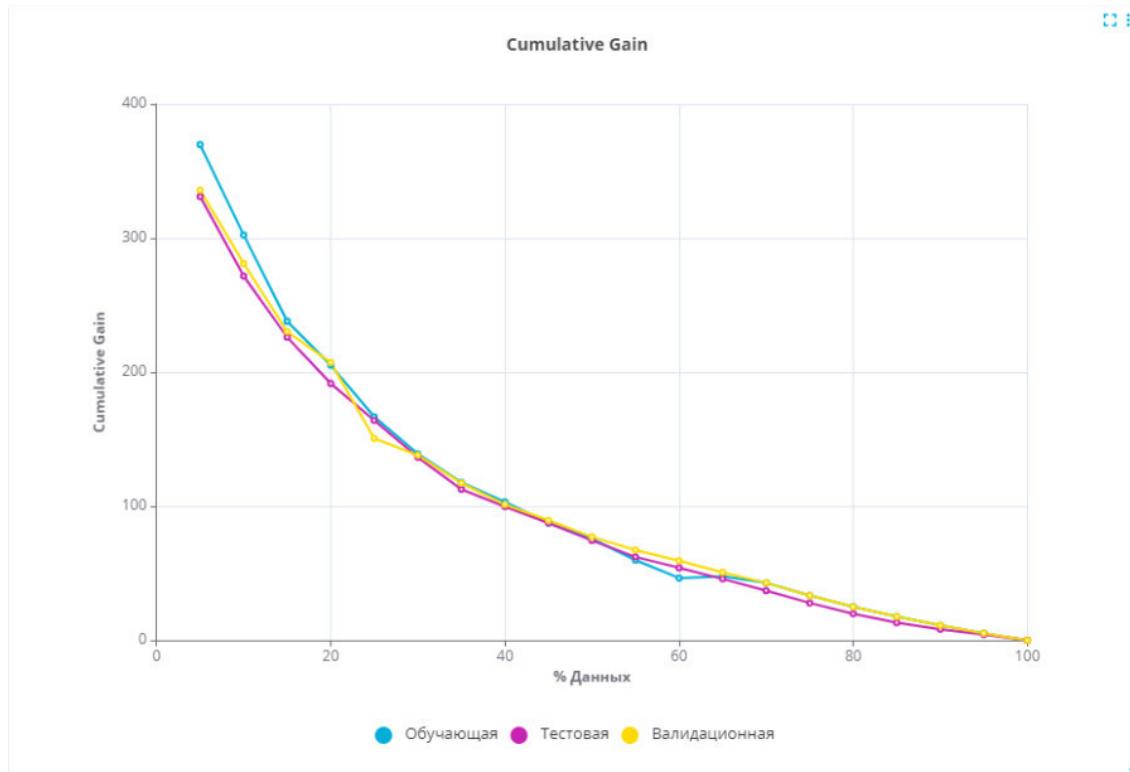


Рисунок 138 Пример графика Cumulative Gain

- Диаграмма дерева решений.

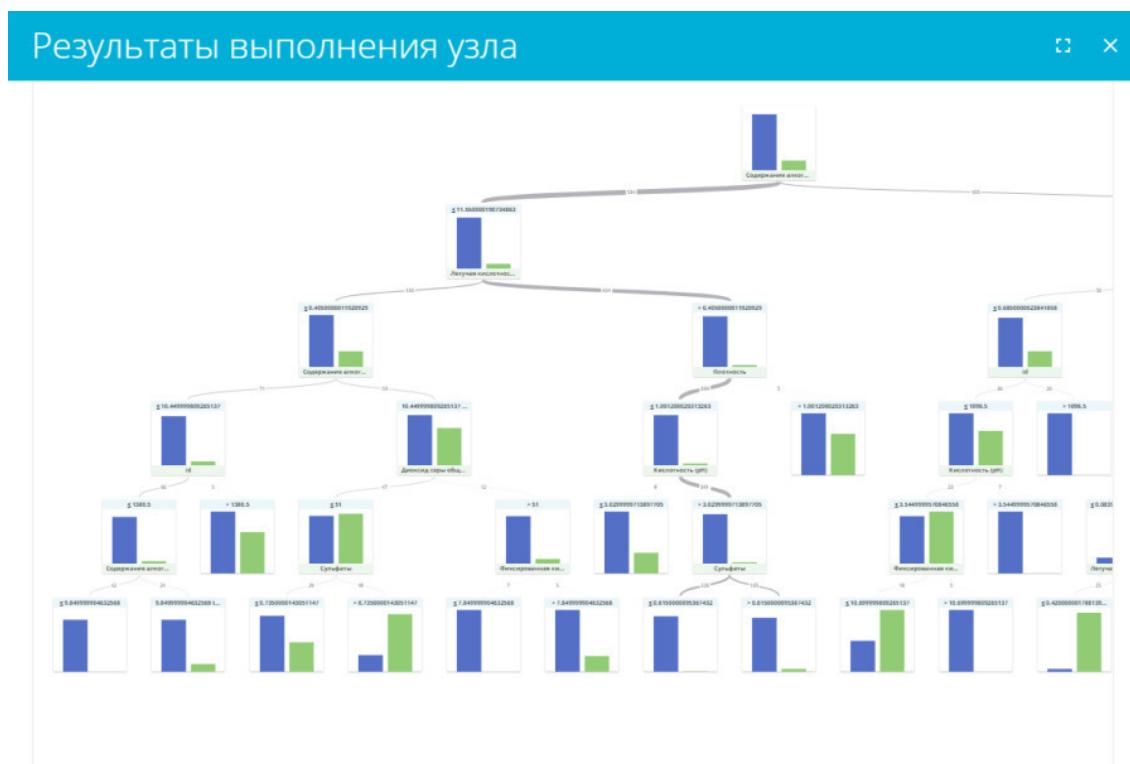


Рисунок 139 Пример графика дерева решений

- Таблица с метриками качества модели.

| ↑             | AUC ROC ↑ | gini ↑    | log loss ↑ | nobs ↑    |
|---------------|-----------|-----------|------------|-----------|
| Filter...     | Filter... | Filter... | Filter...  | Filter... |
| Валидационная | 0.823     | 0.646     | 0.347      | 14003     |
| Обучающая     | 0.83      | 0.66      | 0.334      | 18591     |
| Тестовая      | 0.817     | 0.635     | 0.351      | 14084     |

**Рисунок 140 Пример таблицы с метриками качества модели**

- Таблица с метриками качества модели для задачи классификации.

| ↑             | misclassification ↑ | mcc ↑     | nobs ↑    |
|---------------|---------------------|-----------|-----------|
| Filter...     | Filter...           | Filter... | Filter... |
| Валидационная | 0.141               | 0.413     | 14003     |
| Обучающая     | 0.133               | 0.457     | 18591     |
| Тестовая      | 0.142               | 0.413     | 14084     |

**Рисунок 141 Пример таблицы с метриками качества модели для задачи классификации**

- Таблица со списком переменных, сортированных по важности.

| Переменные                                       | Важность ↓ |
|--|------------|
| Filter...  | Filter...  |
| Количество просрочек 90+ дней                    | 0.584      |
| Количество просрочек 30-59 дней                  | 0.195      |
| Использование необеспеченного кредитного баланса | 0.087      |
| Количество просрочек 60-89 дней                  | 0.057      |
| Род деятельности_Прочее                          | 0.023      |

**Рисунок 142 Пример таблицы со списком переменных, отсортированных по важности**

Результаты **многоклассовой классификации** представлены следующими объектами:

- Диаграмма дерева решений.
- Таблица с метриками качества модели.

| <a href="#">↑</a>         | <a href="#">log loss ↑</a> | <a href="#">nobs ↑</a>    |
|---------------------------|----------------------------|---------------------------|
| <a href="#">Filter...</a> | <a href="#">Filter...</a>  | <a href="#">Filter...</a> |
| Валидационная             | 2.883                      | 27                        |
| Обучающая                 | 0.556                      | 36                        |
| Тестовая                  | 5.524                      | 27                        |

**Рисунок 143 Пример таблицы с метриками качества модели**

- Таблица с метриками качества модели для задачи классификации.

| <a href="#">↑</a>         | <a href="#">misclassi... ↑</a> | <a href="#">mcc ↑</a>     | <a href="#">nobs ↑</a>    |
|---------------------------|--------------------------------|---------------------------|---------------------------|
| <a href="#">Filter...</a> | <a href="#">Filter...</a>      | <a href="#">Filter...</a> | <a href="#">Filter...</a> |
| Валидационна              | 0.111                          | 0.8                       | 27                        |
| Обучающая                 | 0.25                           | 0.601                     | 36                        |
| Тестовая                  | 0.259                          | 0.564                     | 27                        |

**Рисунок 144 Пример таблицы с метриками качества модели для задачи классификации**

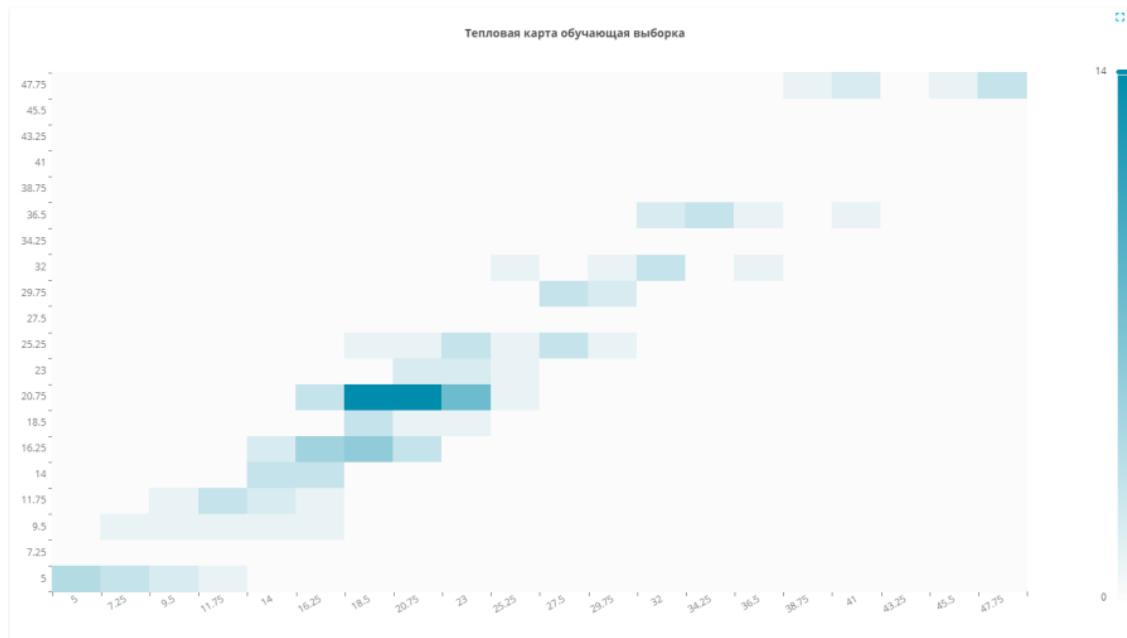
- Таблица со списком переменных, сортированных по важности.

| Переменные          | Важность ↓ |
|---------------------|------------|
| Filter...           | Filter...  |
| electronegativity   | 0.789      |
| boiling_temperarure | 0.185      |
| density             | 0.024      |
| name_Сера           | 0          |
| name_Неодим         | 0          |

**Рисунок 145 Пример таблицы со списком переменных, сортированных по важности**

Результаты **регрессии** представлены следующими объектами:

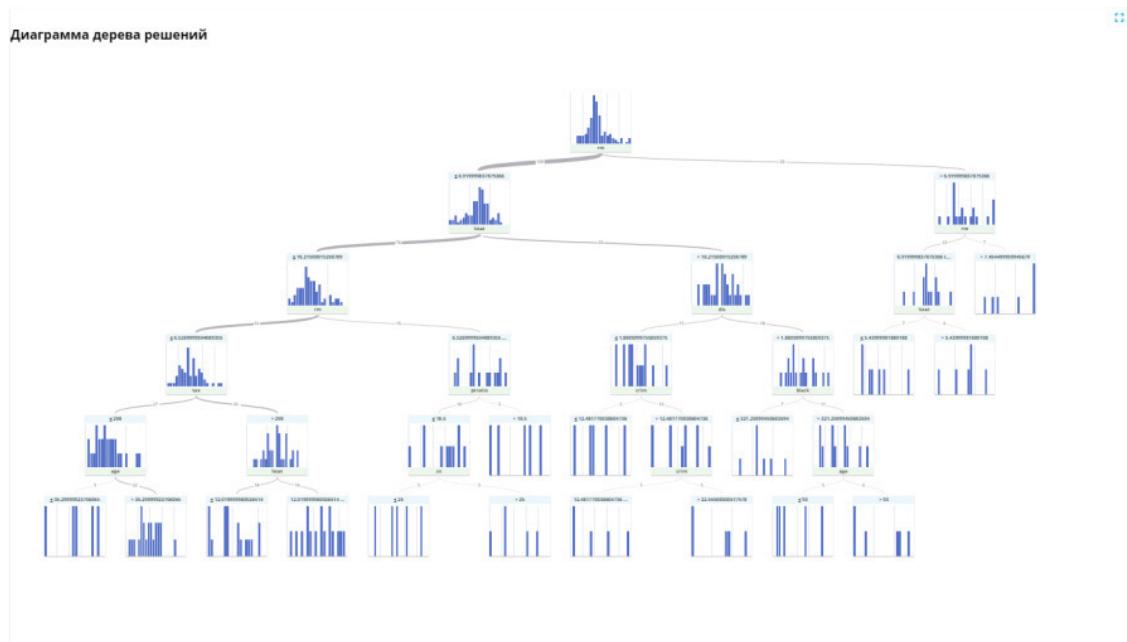
- Тепловые карты обучающей, тестовой и валидационной выборок.



**Рисунок 146 Пример тепловой карты на данных обучающей выборки**

«Тепловые» карты, отражающие корреляцию или ассоциацию типа хи-квадрат первоначальных значений с целевыми признаками по сегментам.

- Диаграмма дерева решений.



**Рисунок 147 Пример диаграммы дерева решений для задачи регрессии**

- Таблица с метриками качества модели.

|               | mse ↑     | rmse ↑    | mae ↑     | mape ↑    | r2 ↑      | nobs ↑    |
|---------------|-----------|-----------|-----------|-----------|-----------|-----------|
| Валидационная | Filter... | Filter... | Filter... | Filter... | Filter... | Filter... |
| Валидационная | 21.978    | 4.688     | 3.222     | 0.178     | 0.758     | 96        |
| Обучающая     | 4.933     | 2.221     | 1.784     | 0.092     | 0.937     | 129       |
| Тестовая      | 26.483    | 5.146     | 3.367     | 0.162     | 0.523     | 98        |

**Рисунок 148 Пример таблицы с метриками качества модели**

- Таблица со списком переменных, сортированных по важности (рисунок выше).

### 3.2.5.8.2. Узел «Случайный лес»

В основе **узла «Случайный лес»** лежит алгоритм машинного обучения, который представляет собой ансамбль деревьев решений.

**Алгоритм работы: Ансамблирование** — это тип обучения, при котором объединяются различные типы алгоритмов или тот же алгоритм несколько раз, что позволяет сформировать более мощную прогнозную модель. Так, в алгоритме **Случайного леса** каждое дерево предсказывает класс (в случае задачи классификации) или значение (в случае задачи регрессии) на основании своего разбиения, и выбирается то предсказание, которое получило наибольшее количество голосов (в случае задачи классификации) или среднее значение всех предсказанных значений (в случае задачи регрессии).

**Список параметров узла** представлен в таблице ниже.

| Параметр   | Возможные значения и ограничения   | Описание  |
|--|--|---|
| <b>Название</b>  | Ручной ввод<br>Ограничений на значение нет   | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>  | Ручной ввод<br>Ограничений на значение нет   | Описание узла   |
| <b>Количество деревьев</b>                             | Ручной ввод<br>Неотрицательное число   | Данный параметр определяет количество деревьев в случайном лесу   |
| <b>Критерий разбиения для классификации</b>            | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Gini</li> <li>• entropy</li> </ul>  | Данный параметр задает критерий разбиения на узлы для классификации. Предусмотрены следующие критерии: <ul style="list-style-type: none"> <li>• gini (коэффициент Джини)</li> <li>• entropy (критерий прироста информации, энтропия)</li> </ul>   |
| <b>Критерий разбиения для регрессии</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• squared error</li> <li>• friedman mse</li> <li>• absolute error</li> <li>• poisson</li> </ul> | Данный параметр задает критерий разбиения для регрессионной задачи. Предусмотрены следующие критерии: <ul style="list-style-type: none"> <li>• squared error (среднеквадратичная ошибка)</li> <li>• friedman mse (среднеквадратичная ошибка с оценкой улучшения Фридмана)</li> <li>• absolute error (средняя абсолютная ошибка)</li> <li>• poisson (отклонение Пуассона)</li> </ul> |
| <b>Максимальная глубина</b>                            | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 5   | Данный параметр задает максимальную глубину дерева, после достижения которой алгоритм останавливает работу.   |
| <b>Минимальное количество наблюдений для разбиения</b> | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 2   | Данный параметр задает минимальное количество наблюдений, которое должно быть в разбиении   |
| <b>Минимальное количество наблюдений в листе</b>       | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 5   | Данный параметр задает минимальное количество наблюдений, которое может быть в листе  |

| <b>Параметр</b>                                 | <b>Возможные значения и ограничения</b>   | <b>Описание</b>  |
|---|---|--|
| <b>Максимальное количество признаков</b>        | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• all</li> <li>• sqrt</li> <li>• log2</li> <li>• number</li> <li>• frac</li> </ul> | Данный параметр определяет максимальное количество признаков, которое будет учитываться при поиске лучшего разделения. Предусмотрены следующие варианты: <ul style="list-style-type: none"> <li>• all – учитывать все доступные признаки</li> <li>• sqrt – учитывать <math>\sqrt{}</math>(число всех доступных признаков)</li> <li>• log2 – учитывать <math>\log_2</math>(число всех доступных признаков)</li> <li>• number – учитывать указанное число признаков</li> <li>• frac – учитывать <math>\frac{1}{\text{число}} * \text{число}</math> всех доступных признаков)</li> </ul> При выборе number или frac появится дополнительный параметр Число (вводится int) и Frac (вводится float) соответственно. |
| <b>Seed</b>                                     | Ручной ввод числового значения<br>По умолчанию — 12345  | Начальное числовое значение для генератора случайных чисел   |
| <b>Максимальное количество листов</b>           | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0  | Данный параметр определяет максимальное количество листов в дереве   |
| <b>Минимальное снижение неоднородности</b>      | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0  | Данный параметр определяет минимальное снижение неоднородности Узел будет разделен, если это разделение вызовет уменьшение неоднородности большее или равное указанному значению   |
| <b>Использовать бутстреп</b>                    | Чекбокс   | Данный чекбокс указывает на необходимость использования метода повторной выборки наблюдений  |
| <b>Размер бутстреп-выборок</b>                  | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 1  | Данный параметр задает размер бутстреп-выборок   |
| <b>Минимальная доля веса наблюдений в листе</b> | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0  | Данный параметр определяет минимальный весовой коэффициент выборки в листовом узле. По умолчанию наблюдения имеют одинаковый вес   |
| <b>Количество параллельных сессий</b>           | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0  | Данный параметр задает количество параллельных сессий  |
| <b>ccp_alpha</b>                                | Ручной ввод<br>По умолчанию — 0   | Данный параметр регулирует количество отсекаемых узлов. Чем больше значение ccp_alpha, тем большее количество узлов удаляется из дерева  |

**Таблица 30 Параметры узла «Случайный лес»**

#### **Результаты выполнения узла:**

Узел «Случайный лес» имеет разные результаты в зависимости от решаемой задачи.

Результаты **бинарной классификации** представлены следующими объектами:

- График ROC (аналогично узлу «Дерево решений»).
- График Lift (аналогично узлу «Дерево решений»).
- График Cumulative Lift (аналогично узлу «Дерево решений»).
- График Gain (аналогично узлу «Дерево решений»).
- График Cumulative Gain (аналогично узлу «Дерево решений»).
- Таблица с метриками качества модели (аналогично узлу «Дерево решений»).
- Таблица с метриками качества модели для задачи классификации (аналогично узлу «Дерево решений»).
- Таблица со списком переменных, сортированных по важности (аналогично узлу «Дерево решений»).

Результаты **многоклассовой классификации** представлены следующими объектами:

- Таблица с метриками качества модели (аналогично узлу «Дерево решений»).
- Таблица с метриками качества модели для задачи классификации (аналогично узлу «Дерево решений»).
- Таблица со списком переменных, сортированных по важности (аналогично узлу «Дерево решений»).

Результаты **регрессии** представлены следующими объектами:

- Тепловые карты обучающей, тестовой и валидационной выборок (аналогично узлу «Дерево решений»).
- Таблица с метриками качества модели (аналогично узлу «Дерево решений»).
- Таблица со списком переменных, сортированных по важности (аналогично узлу «Дерево решений»).

### 3.2.5.8.3. Узел «Байесовская регрессия»

**Узел «Байесовская регрессия»** представляет собой линейную регрессию с применением распределения вероятностей параметров, а не точечных оценок.

**Алгоритм работы:** В основе узла лежит Байесовская гребневая регрессия, где оптимизируются параметры регуляризации *lambda* (точность весов) и *alpha* (точность шума).

**Список параметров узла** представлен в таблице ниже.

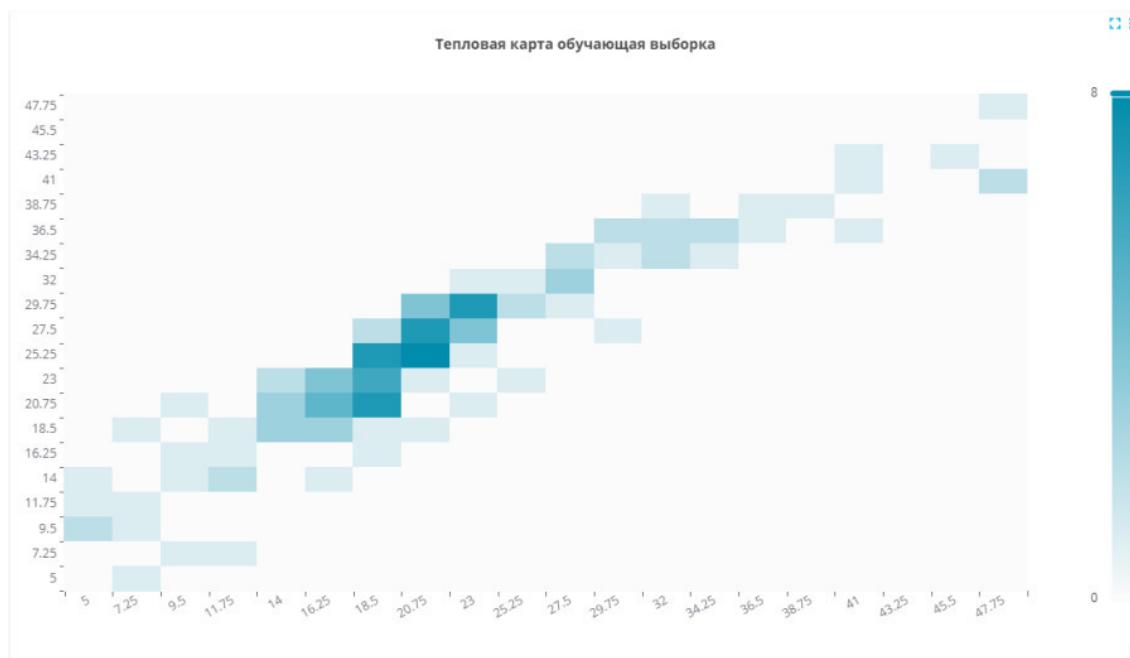
| Параметр        | Возможные значения и ограничения           | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе |

| <b>Параметр</b>                    | <b>Возможные значения и ограничения</b>   | <b>Описание</b>  |
|------------------------------------|---|--|
| <b>Описание</b>                    | Ручной ввод<br>Ограничений на значение нет  | Описание узла  |
| <b>Количество итераций</b>         | Ручной ввод<br>Неотрицательное число, больше или равно 1<br>По умолчанию — 300  | Данный параметр задает количество итераций, после достижения которого алгоритм останавливается   |
| <b>Добавить константу в модель</b> | Чекбокс   | Выбор данного чекбокса добавит константу в модель  |
| <b>Допустимая погрешность</b>      | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,0001   | Данный параметр задает допустимую погрешность, после достижения которой алгоритм останавливается   |
| <b>Стандартизация</b>              | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• no (по умолчанию)</li> <li>• std</li> <li>• range</li> </ul> | Данный параметр отвечает за выбор метода стандартизации данных.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• no — нет</li> <li>• std — стандартное отклонение - преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li> <li>• range — диапазон - линейно преобразует значения переменных в диапазон [0, 1].</li> </ul> |
| <b>Alpha 1</b>                     | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,000001   | Данный параметр представляет собой параметр формы для предварительного распределения Гамма над параметром альфа  |
| <b>Alpha 2</b>                     | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,000001   | Данный параметр представляет собой обратный параметр масштаба (параметр скорости) для предварительного распределения Гамма над параметром альфа  |
| <b>Alpha init</b>                  | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0  | Данный параметр представляет собой начальное значение для alpha (точность шума)  |
| <b>Lambda 1</b>                    | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,000001   | Данный параметр представляет собой параметр формы для предварительного распределения Гамма над параметром лямбда   |
| <b>Lambda 2</b>                    | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,000001   | Данный параметр представляет собой обратный параметр масштаба (параметр скорости) для предварительного распределения Гамма над параметром лямбда   |
| <b>Lambda init</b>                 | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 1  | Данный параметр представляет собой начальное значение для лямбды (точность весов)  |

**Таблица 31 Параметры узла «Байесовская регрессия»**

## Результаты выполнения узла:

- Тепловые карты для обучающей, валидационной и тестовой выборок.



**Рисунок 149 Пример тепловой карты на данных обучающей выборки**

- Таблица с метриками качества модели.

| Метрики модели |           |           |           |           |           |           |
|----------------|-----------|-----------|-----------|-----------|-----------|-----------|
| ↑              | mse ↑     | rmse ↑    | mae ↑     | mape ↑    | r2 ↑      | nobs ↑    |
| Filter...      | Filter... | Filter... | Filter... | Filter... | Filter... | Filter... |
| Валидационная  | 26.783    | 5.175     | 3.513     | 0.186     | 0.705     | 96        |
| Обучающая      | 12.267    | 3.502     | 2.622     | 0.14      | 0.844     | 129       |
| Тестовая       | 21.466    | 4.633     | 3.009     | 0.153     | 0.613     | 98        |

**Рисунок 150 Пример таблицы с метриками качества модели**

- Таблица с коэффициентами переменных.

| Коэффициенты Переменных |               |
|-------------------------|---------------|
| Переменные ↑            | Коэффициент ↑ |
| Filter...               | Filter...     |
| Intercept               | 19.877        |
| age                     | -4.549        |
| black                   | 5.565         |
| chas_1.0                | -0.024        |
| crim                    | -3.823        |

**Рисунок 151 Пример таблицы с коэффициентами переменных**

#### 3.2.5.8.4. Узел «Линейная регрессия»

В основе **узла «Линейная регрессия»** лежит модель зависимости между входными и выходными переменными с линейной функцией связи.

**Список параметров узла** представлен в таблице ниже.

| Параметр                           | Возможные значения и ограничения                           | Описание   |
|------------------------------------|--|--|
| <b>Название</b>                    | Ручной ввод<br>Ограничений на значение нет                 | Название узла, которое будет отображаться в интерфейсе     |
| <b>Описание</b>                    | Ручной ввод<br>Ограничений на значение нет                 | Описание узла  |
| <b>Добавить константу в модель</b> | Чекбокс  | Выбор данного чекбокса добавит константу в модель          |
| <b>L1</b>                          | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,5 | Данный параметр задает значение L1-регуляризации           |
| <b>L2</b>                          | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,5 | Данный параметр задает значение L2-регуляризации           |
| <b>Seed</b>                        | Ручной ввод числового значения<br>По умолчанию — 42        | Начальное числовое значение для генератора случайных чисел |

| <b>Параметр</b>                                  | <b>Возможные значения и ограничения</b>  | <b>Описание</b>  |
|--|--|--|
| <b>Количество итераций</b>                       | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 100   | Данный параметр задает количество итераций, после достижения алгоритм останавливается  |
| <b>Метод оптимизации для гребневой регрессии</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• svd</li> <li>• cholesky</li> <li>• lsqr (по умолчанию)</li> <li>• sparse_cg</li> <li>• sag</li> <li>• saga</li> </ul> | Данный параметр задает метод оптимизации для гребневой регрессии. Предусмотрены следующие варианты: <ul style="list-style-type: none"> <li>• svd - использует сингулярное разложение</li> <li>• cholesky</li> <li>• lsqr - использует специальную процедуру упорядоченных наименьших квадратов</li> <li>• sparse_cg - использует решатель сопряженных градиентов</li> <li>• sag - использует stochastic average gradient descent</li> <li>• saga - использует улучшенную версию stochastic average gradient descent</li> </ul> |
| <b>Допустимая погрешность</b>                    | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,0001  | Данный параметр задает допустимую погрешность, после достижения которой алгоритм останавливается   |
| <b>Правило обновления коэффициентов модели</b>   | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• cyclic (по умолчанию)</li> <li>• random</li> </ul>  | Данный параметр задает правило обновления коэффициентов модели. Предусмотрены следующие варианты: <ul style="list-style-type: none"> <li>• cyclic - последовательный перебор</li> <li>• random - случайные коэффициенты обновляются каждую итерацию</li> </ul>   |

**Таблица 32 Параметры узла «Линейная регрессия»**

#### **Результаты выполнения узла:**

- Тепловые карты для обучающей, валидационной и тестовой выборок (Аналогично узлу «Байесовская регрессия»).
- Таблица с метриками качества модели (Аналогично узлу «Байесовская регрессия»).
- Таблица с коэффициентами переменных (Аналогично узлу «Байесовская регрессия»).

#### 3.2.5.8.5. Узел «Логистическая регрессия»

В основе **узла «Логистическая регрессия»** лежит метод построения линейного классификатора, позволяющий оценивать апостериорные вероятности принадлежности объектов классам.

**Список параметров узла** представлен в таблице ниже.

| <b>Параметр</b>                         | <b>Возможные значения и ограничения</b>   | <b>Описание</b>  |
|---|---|--|
| <b>Название</b>                         | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                         | Ручной ввод<br>Ограничений на значение нет  | Описание узла  |
| <b>Тип многоклассовой классификации</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Бинарная классификация (по умолчанию)</li> <li>• Многоклассовая классификация</li> </ul>                 | Данный параметр задает тип многоклассовой классификации. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• Бинарная классификация – метод «один против остальных»</li> <li>• Многоклассовая классификация – метод «многие против многих» Если решается бинарная логистическая регрессия, разницы между рассматриваемыми методами нет.</li> </ul>  |
| <b>Добавить константу в модель</b>      | Чекбокс   | Выбор данного чекбокса добавит константу в модель  |
| <b>L1</b>                               | Ручной ввод<br>Число больше или равное 0<br>По умолчанию — 0,5  | Данный параметр задает значение L1-регуляризации   |
| <b>L2</b>                               | Ручной ввод<br>Число больше или равное 0<br>По умолчанию — 0,5  | Данный параметр задает значение L2-регуляризации   |
| <b>Seed</b>                             | Ручной ввод числового значения<br>По умолчанию — 42   | Начальное числовое значение для генератора случайных чисел   |
| <b>Количество итераций</b>              | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 100  | Данный параметр задает количество итераций, после достижения алгоритм останавливается  |
| <b>Метод оптимизации</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• lbfgs (по умолчанию)</li> <li>• newton-cg</li> <li>• liblinear</li> <li>• sag</li> <li>• saga</li> </ul> | Данный параметр задает метод оптимизации. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• lbfgs – использует вторую производную матрицу функции потерь для итерационной оптимизации функции потерь</li> <li>• newton-cg – использует алгоритм Newton-Conjugate-Gradient</li> <li>• liblinear – использует метод спуска оси координат для итеративной оптимизации функции потерь</li> <li>• sag – стохастический средний градиентный спуск</li> <li>• saga – использует улучшенную версию стохастического среднего градиентного спуска</li> </ul> Метод liblinear доступен только при выборе бинарной классификации. |

| Параметр                      | Возможные значения и ограничения                              | Описание  |
|-------------------------------|---|---|
|                               |   | Методы newton-cg, lbfgs, sag и saga обрабатывают L2 или без штрафа. Методы liblinear и saga также обрабатывают штраф L1. Метод saga также поддерживает штраф elasticnet (L1 + L2). Для небольших наборов данных метод liblinear - хороший выбор, тогда как sag и saga быстрее для больших. В случае неправильно выставленных пользователем методов оптимизации, автоматически будет произведена замена. |
| <b>Допустимая погрешность</b> | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,0001 | Данный параметр задает допустимую погрешность для оптимизации, после достижения которой алгоритм останавливается  |

**Таблица 33 Параметры узла «Логистическая регрессия»**

**Результаты выполнения узла:**

- Таблица с метриками качества модели.

| Метрики модели |            |           |
|----------------|------------|-----------|
| ↑              | log loss ↑ | nobs ↑    |
| Filter...      | Filter...  | Filter... |
| Валидационная  | 1.08       | 27        |
| Обучающая      | 1.224      | 36        |
| Тестовая       | 1.099      | 27        |

**Рисунок 152 Пример таблицы с метриками качества модели**

- Таблица с метриками качества модели задачи классификации.

**Метрики модели (классификация)**

| ↑             | misclassification ↑ | mcc ↑     | nobs ↑    |
|---------------|---------------------|-----------|-----------|
| Filter...     | Filter...           | Filter... | Filter... |
| Валидационная | 0.222               | 0.581     | 27        |
| Обучающая     | 0.333               | 0.408     | 36        |
| Тестовая      | 0.333               | 0.427     | 27        |

**Рисунок 153 Пример таблицы с метриками качества модели задачи классификации**

- Таблица с коэффициентами переменных.

**Коэффициенты Переменных**

| Переменны... ↑ | Актиноид ↑ | Газ ↑     | Лантаноид ↑ | Металл ↑  | Неметалл ↑ | Полуметалл ↑ |
|----------------|------------|-----------|-------------|-----------|------------|--------------|
| Filter...      | Filter...  | Filter... | Filter...   | Filter... | Filter...  | Filter...    |
| Intercept      | -0.001     | -0.001    | -0.001      | 0         | -0.001     | -0.001       |
| Level 1_2      | -0.001     | -0.001    | -0.001      | 0         | -0.001     | -0.001       |
| Level 2        | -0.001     | -0.001    | -0.001      | 0         | -0.001     | -0.001       |
| Level 3        | -0.001     | -0.001    | -0.001      | 0         | -0.001     | -0.001       |
| Level 4        | -0.001     | -0.001    | 0           | 0         | -0.001     | -0.001       |

**Рисунок 154 Пример таблицы с коэффициентами переменных**

### 3.2.5.8.6. Узел «Линейные модели»

**Узел «Линейные модели»** объединяет в себе линейные классификаторы и регрессоры с обучением методом стохастического градиентного спуска и поддерживает различные функции потерь и штрафы.

**Список параметров узла** представлен в таблице ниже

| Параметр                                     | Возможные значения и ограничения  | Описание  |
|--|---|---|
| <b>Название</b>                              | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>                              | Ручной ввод<br>Ограничений на значение нет  | Описание узла   |
| <b>Функция потерь для классификации</b>      | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• hinge (по умолчанию)</li> <li>• log</li> <li>• modified_huber</li> <li>• squared_hinge</li> <li>• perceptron</li> <li>• squared_loss</li> <li>• huber</li> <li>• epsilon_insensitive</li> <li>• squared_epsilon_insensitive</li> </ul> | Данный параметр задает функцию потерь для классификационной задачи. Предусмотрены следующие функции: <ul style="list-style-type: none"> <li>• hinge – средняя потеря петель</li> <li>• log – логистическая регрессия</li> <li>• modified_huber – сглаженная потеря петли</li> <li>• squared_hinge – похож на hinge, но его штраф введен в квадрат</li> <li>• perceptron – перцептрон</li> <li>• squared_loss – метод наименьших квадратов</li> <li>• huber – потеря Хубера</li> <li>• epsilon_insensitive</li> <li>• squared_epsilon_insensitive</li> </ul> |
| <b>Функция потерь для регрессии</b>          | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• squared_error (по умолчанию)</li> <li>• huber</li> <li>• epsilon_insensitive</li> <li>• squared_epsilon_insensitive</li> </ul>   | Данный параметр задает функцию потерь для регрессионной задачи. Предусмотрены следующие функции: <ul style="list-style-type: none"> <li>• squared_error – метод наименьших квадратов</li> <li>• huber – потеря Хубера</li> <li>• epsilon_insensitive</li> <li>• squared_epsilon_insensitive</li> </ul>  |
| <b>Epsilon</b>                               | Ручной ввод<br>Неотрицательное число<br>По умолчанию – 0,1  | Данный параметр задает Epsilon в функцию потерь   |
| <b>L1</b>                                    | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 0,5   | Данный параметр задает значение L1-регуляризации  |
| <b>L2</b>                                    | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 0,5   | Данный параметр задает значение L2-регуляризации  |
| <b>Правило определения скорости обучения</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• constant</li> <li>• optimal (по умолчанию)</li> <li>• invscaling</li> <li>• adaptive</li> </ul>  | Данный параметр задает правило определения скорости обучения. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• constant – постоянная скорость обучения</li> <li>• optimal – определяется на основе эвристики, предложенной Леоном Ботту</li> <li>• invscaling – обратное масштабирование</li> <li>• adaptive – адаптивное уменьшение скорости обучения</li> </ul>  |

| <b>Параметр</b>  | <b>Возможные значения и ограничения</b>   | <b>Описание</b>   |
|--|---|---|
| <b>Начальная скорость обучения</b>                         | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,001  | Данный параметр задает начальную скорость обучения  |
| <b>Стандартизация</b>                                      | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• no (по умолчанию)</li> <li>• std</li> <li>• range</li> </ul> | Данный параметр отвечает за выбор метода стандартизации данных. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• no — нет</li> <li>• std — стандартное отклонение - преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li> <li>• range — диапазон - линейно преобразует значения переменных в диапазон [0, 1].</li> </ul> |
| <b>Метод построения вероятностей</b>                       | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• linear (по умолчанию)</li> <li>• logistic</li> </ul>         | Данный параметр задает метод построения вероятностей  |
| <b>Показатель степени для изменения скорости обучения</b>  | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,5  | Данный параметр задает показатель степени для изменения скорости обучения   |
| <b>Добавить константу в модель</b>                         | Чекбокс   | Выбор данного чекбокса добавит константу в модель   |
| <b>Перемешать наблюдения</b>                               | Чекбокс   | Выбор данного чекбокса указывает на необходимость перемешать наблюдения   |
| <b>Seed</b>  | Ручной ввод числового значения<br>По умолчанию — 42   | Начальное числовое значение для генератора случайных чисел  |
| <b>Усреднение коэффициентов</b>                            | Ручной ввод<br>Число больше или равно 0<br>По умолчанию — 0   | Данный параметр задает вычисление усредненных коэффициентов в результирующей линейной модели  |
| <b>Количество итераций</b>                                 | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 1000   | Данный параметр задает количество итераций, после достижения алгоритм останавливается   |
| <b>Ранняя остановка</b>                                    | Чекбокс   | Выбор данного чекбокса указывает на необходимость ранней остановки алгоритма, если валидационная оценка не улучшается   |
| <b>Размер % валидационной выборки для ранней остановки</b> | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию — 0,1  | Данный параметр задает долю обучающих данных, которые нужно отложить в качестве валидационного набора для ранней остановки  |

| Параметр                      | Возможные значения и ограничения                              | Описание  |
|-------------------------------|---|---|
| <b>Допустимая погрешность</b> | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,0001 | Данный параметр задает допустимую погрешность для оптимизации |

**Таблица 34 Параметры узла «Линейные модели»**

#### Результаты выполнения узла:

**Узел «Линейные модели»** имеет разные результаты в зависимости от решаемой задачи.

Результаты **регрессии** представлены следующими объектами:

- Тепловые карты для обучающей, валидационной и тестовой выборок (Аналогично узлу «Байесовская регрессия»).
- Таблица с метриками качества модели (Аналогично узлу «Байесовская регрессия»).
- Таблица с коэффициентами переменных (Аналогично узлу «Байесовская регрессия»).

Результаты **бинарной классификации** представлены следующими объектами:

- График ROC.
- График Lift.
- График Cumulative Lift.
- График Gain.
- График Cumulative Gain.
- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица с коэффициентами переменных.

Результаты **многоклассовой классификации** представлены следующими объектами:

- Таблица с метриками качества модели (Аналогично узлу «Логистическая регрессия»).
- Таблица с метриками качества модели задачи классификации (Аналогично узлу «Логистическая регрессия»).
- Таблица с коэффициентами переменных (Аналогично узлу «Логистическая регрессия»).

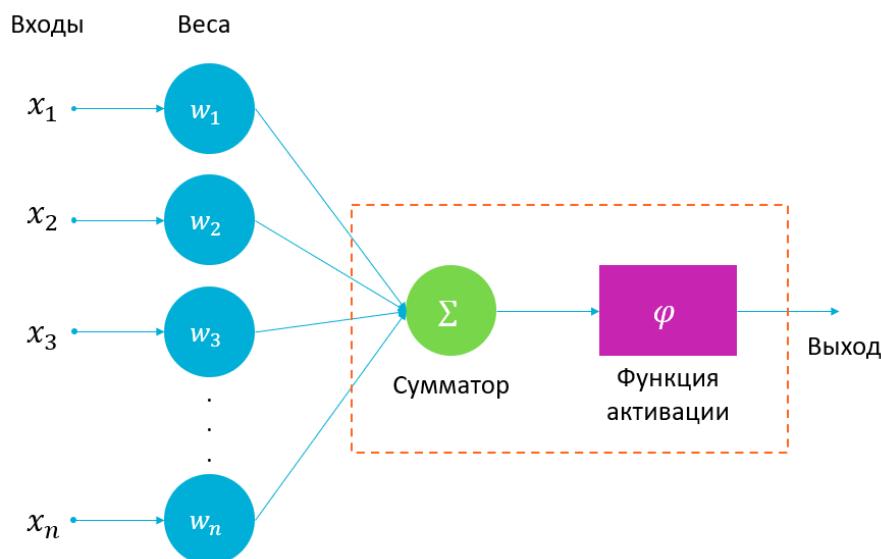
#### 3.2.5.8.7. Узел «Нейронная сеть»

**Узел «Нейронная сеть»** позволяет строить нейросеть типа **MLP (multilayer perceptron, многослойный перцептрон)**, в которой входной сигнал преобразуется в выходной, проходя последовательно через скрытые слои.

**Искусственная нейронная сеть** представляет собой сложную дифференцируемую функцию, которая в процессе обучения способна выявлять сложные зависимости между входными и выходными данными и выполнять обобщение.

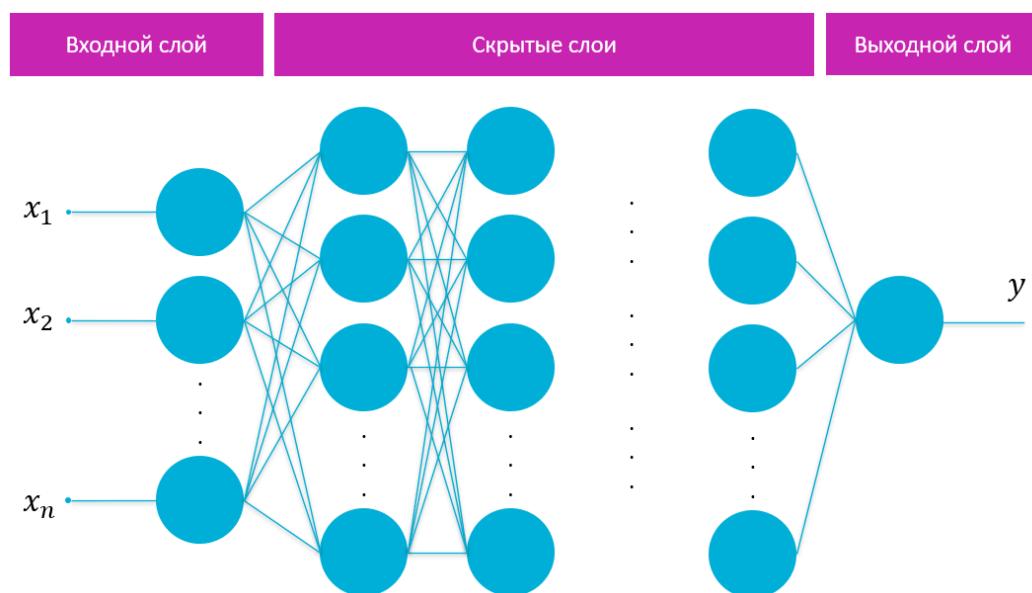
### 3.2.5.8.7.1. Архитектура нейронной сети

**Нейрон** – базовый элемент для построения искусственной нейронной сети (далее просто нейронной сети). Представляет собой вычислительный элемент, который имеет несколько входов и один выход. Каждый вход имеет некоторый **вес**, на который умножается поступившее на данный вход значение. Далее в нейроне происходит суммирование взвешенных входных данных, и полученная сумма преобразуется с помощью **функции активации** в значение, которое и будет являться выходом нейрона.



### **Рисунок 155 Строение искусственного нейрона**

В многослойных сетях нейроны сгруппированы в **слои**, причем каждый нейрон предыдущего слоя связан со всеми нейронами следующего слоя, а внутри слоев связь между нейронами отсутствует.



### **Рисунок 156 Слои нейронной сети**

Первый слой сети называется **входным**: его нейроны принимают элементы вектора признаков и передают их нейронам следующего слоя без какой-либо обработки.

Далее располагаются один или несколько **скрытых слоев** (число скрытых слоев и количество нейронов в них задается **параметром Количество нейронов в скрытом слое**, подробнее в таблице с описанием параметров ниже). Каждый нейрон в скрытом слое преобразует значения из предыдущего слоя с взвешенным линейным суммированием функцией активации, которая порождает более информативные признаковые описания, преобразую данные нелинейным образом (функция активации для скрытых слоев задается **параметром Функция активации**).

Последний слой – **выходной**, получает значения из последнего скрытого слоя и преобразует их в выходные значения сети исходя из решаемой задачи.

### 3.2.5.8.7.2. Обучение нейронной сети

---

Обучение нейронной сети является итерационным процессом (задается параметром **Количество итераций**), которое заключается в нахождении параметров модели (весов и смещений). На каждой итерации происходит два прохода сети:

1. Настройка нейронной сети начинается с **Прямого распространения**. Данные (задается параметром **Количество наблюдений для одной итерации**) подаются на вход сети, через все слои распространяются в направлении выходов, в ходе чего случайным образом генерируются значения весовых коэффициентов и смещений модели. Таким образом, формируются выходные значения. Далее *фактические выходные значения* оцениваются относительно *ожидаемых выходных значений* (целевых значений, *target*) с помощью **функции потерь** (задается параметрами **Метрика для ранней остановки (регрессия)** или **Метрика для ранней остановки (классификация)** в зависимости от решаемой задачи).
2. **Обратное распространение ошибки (Backward propagation)** направлено на минимизацию функции потерь путем корректировки параметров нейронной сети (весов и смещений нейронов). Скорректировать эти параметры позволяют методы оптимизации на основе градиентного спуска (задается параметром **Метод оптимизации**).

Процесс оптимизации будет выполняться до тех пор, пока не будет достигнут критерий остановки (задается параметрами **Ранняя остановка**, **Погрешность оптимизации** и **Количество итераций**).

### 3.2.5.8.7.3. Список параметров узла

**Список параметров узла** представлен в таблице ниже.

| Параметр                                  | Возможные значения и ограничения  | Описание   |
|---|---|--|
| <b>Название</b>                           | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                           | Ручной ввод<br>Ограничений на значение нет  | Описание узла  |
| <b>Количество нейронов в скрытом слое</b> | Для изменения количества нейронов и слоев необходимо: <ul style="list-style-type: none"> <li>Выбрать кнопку <b>Изменить</b></li> <li>Добавить необходимое количество слоев (Иконка + позволяет добавить скрытый слой)</li> <li>Указать необходимое число нейронов в каждом слое</li> <li>Выбрать кнопку <b>«Сохранить»</b><br/>По умолчанию — 1 слой, в котором 100 нейронов</li> </ul> | Данный параметр задает количество скрытых слоев и нейронов в скрытых слоях   |
| <b>Функция активации</b>                  | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>Identity</li> <li>logistic</li> <li>tanh</li> <li>relu (по умолчанию)</li> </ul>   | Данный параметр задает функцию активации для скрытых слоев.<br>Предусмотрены следующие функции: <ul style="list-style-type: none"> <li><b>Identity</b> – линейная функция, результат пропорционален переданному аргументу <math>f(x) = x</math></li> <li><b>logistic</b> – сигмоидная (логистическая) функция               <math display="block">\sigma(x) = \frac{1}{1 + e^{-x}}</math> </li> <li><b>tanh</b> – функция гиперболического тангенса               <math display="block">f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}}</math> </li> <li><b>relu</b> – функция активации ReLu (Rectified Linear Unit)               <math display="block">f(x) = \max(0, x)</math> </li> </ul> |
| <b>Метод оптимизации</b>                  | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>lbfgs</li> <li>sgd</li> </ul>  | Данный параметр задает метод оптимизации, который будет использоваться для обновления весов нейронов скрытых слоев нейронной сети. Предусмотрены следующие методы:   |

| Параметр                                     | Возможные значения и ограничения  | Описание   |
|--|---|--|
|  | <ul style="list-style-type: none"> <li>adam (по умолчанию)</li> </ul>   | <ul style="list-style-type: none"> <li><b>sgd</b> – стохастический градиентный спуск.</li> <li><b>lbfgs</b> – алгоритм оптимизации из семейства квазиньютоновских методов, который аппроксимирует алгоритм Бродена–Флетчера–Гольдфарба–Шанно (BFGS) с использованием ограниченного объема памяти.</li> <li><b>adam</b> – метод адаптивной оценки моментов</li> </ul> <div style="background-color: #e6f2ff; padding: 10px;"> <p><b>Примечание</b><br/>     Метод <b>SGD</b> стоит использовать на небольших сбалансированных наборах данных, в которых достаточно равномерно представлены элементы каждого класса. В случае несбалансированности исходной выборки, стохастический градиент не производит качественного распознавания редких значений признаков, также низкая скорость сходимости проявляется на большом объеме данных.<br/>     Метод <b>adam</b> довольно хорошо работает с большими наборами данных с точки зрения как времени обучения, так и валидационной оценки. Однако для небольших наборов данных метод <b>lbfgs</b> может сходиться быстрее и работать лучше.</p> </div> |
| <b>L2 регуляризация</b>                      | Ручной ввод<br>Неотрицательное число<br>По умолчанию — 0,0001   | Данный параметр задает значение L2-регуляризации   |
| <b>Правило определения скорости обучения</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>constant (по умолчанию)</li> </ul> | Данный параметр задает правило определения скорости обучения для обновления веса. Предусмотрены следующие варианты:  |

| Параметр  | Возможные значения и ограничения  | Описание   |
|---|---|--|
|   | <ul style="list-style-type: none"> <li>• invscalling</li> <li>• adaptive</li> </ul> | <ul style="list-style-type: none"> <li>• <b>constant</b> – постоянная скорость обучения, которая задается в <b>параметре Начальная скорость обучения</b></li> <li>• <b>invscalling</b> – постепенно снижает скорость обучения на каждом шаге, используя показатель обратный указанному в <b>параметре Показатель степени для изменения скорости обучения</b></li> <li>• <b>adaptive</b> – поддерживает скорость обучения постоянной на уровне значения, указанного в параметре <b>Начальная скорость обучения</b>, пока потери при обучении продолжают уменьшаться. Каждый раз, когда две последовательные эпохи не могут уменьшить потери при обучении по крайней мере на значение, указанное в параметре <b>Погрешность оптимизации</b>, или не могут увеличить оценку проверки по крайней мере на это же значение, если включена <b>Ранняя остановка</b>, текущая скорость обучения делится на 5.</li> </ul> <div style="background-color: #e0f2ff; padding: 10px; margin-top: 10px;"> <p><b>Примечание</b><br/>Данный параметр можно задать только в случае выбранного <b>метода оптимизации sgd</b>.</p> </div> |
| <b>Начальная скорость обучения</b>                        | Ручной ввод<br>Число больше 0<br>По умолчанию — 0,001                               | Данный параметр задает начальную скорость обучения, которая управляет размером шага при обновлении весов.<br><div style="background-color: #e0f2ff; padding: 10px; margin-top: 10px;"> <p><b>Примечание</b><br/>Данный параметр можно задать только в случае выбранного <b>метода оптимизации sgd или adam</b>.</p> </div>   |
| <b>Показатель степени для изменения скорости обучения</b> | Ручной ввод<br>Число больше 0<br>По умолчанию — 0,5                                 | Данный параметр задает показатель степени для изменения скорости обучения, который используется для обновления скорости обучения, когда для параметра <b>Правило определения скорости обучения</b> указано значение <b>invscalling</b> .   |

| Параметр  | Возможные значения и ограничения  | Описание  |
|---|---|---|
|   |   | <p><b>Примечание</b><br/>Данный параметр можно задать только в случае выбранного <b>метода оптимизации sgd</b>.</p>   |
| <b>Перемешивание наблюдений</b>                 | Чекбокс   | <p>Выбор данного чекбокса указывает на необходимость перемешивать наблюдения в каждой итерации.</p> <p><b>Примечание</b><br/>Данный параметр можно задать только в случае выбранного <b>метода оптимизации sgd или adam</b>.</p>  |
| <b>Seed</b>                                     | Ручной ввод числового значения<br>По умолчанию — 42                     | Начальное числовое значение для генератора случайных чисел.<br>Используется для воспроизведения результатов при повторном запуске узла  |
| <b>Максимальное количество итераций (эпох)</b>  | Ручной ввод<br>Число больше 0<br>По умолчанию — 200                     | Данный параметр задает максимальное количество итераций. Оптимизатор выполняет итерации до сходимости (определенной значением <b>параметра Погрешность оптимизации</b> ) или до указанного количества итераций.   |
| <b>Количество наблюдений для одной итерации</b> | Ручной ввод<br>Число больше 0<br>По умолчанию — 200                     | Данный параметр задает количество наблюдений для одной итерации   |
| <b>Ранняя остановка</b>                         | Чекбокс   | <p>Выбор данного чекбокса указывает на необходимость ранней остановки алгоритма, если валидационная оценка не улучшается на значение, указанное в параметре <b>Погрешность оптимизации</b>, в течение указанного в параметре <b>Количество итераций</b> без существенного улучшения значения.</p> <p><b>Примечание</b><br/>Ранняя остановка эффективна только в случае выбранного <b>метода оптимизации sgd или adam</b>.</p> |
| <b>Метрика для ранней остановки (регрессия)</b> | Раскрывающийся список со следующими значениями:<br>• MSE (по умолчанию) | Данный параметр задает метрику для ранней остановки (задача регрессии). Предусмотрены следующие метрики:<br>• <b>MSE</b> – среднеквадратичная ошибка  |

| Параметр   | Возможные значения и ограничения   | Описание   |
|--|--|--|
|  | <ul style="list-style-type: none"> <li>• MAE</li> <li>• MAPE</li> <li>• R2</li> </ul>  | <ul style="list-style-type: none"> <li>• <b>MAE</b> – средняя абсолютная ошибка</li> <li>• <b>MAPE</b> – средняя абсолютная ошибка в процентах</li> <li>• <b>R2</b> – коэффициент детерминации</li> </ul>  |
| <b>Метрика для ранней остановки (классификация)</b>          | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Logloss (по умолчанию)</li> <li>• Accuracy</li> <li>• MCC</li> <li>• AUC ROC</li> </ul> | Данный параметр задает метрику для ранней остановки (задача классификации). Предусмотрены следующие метрики: <ul style="list-style-type: none"> <li>• <b>Logloss</b> – логистическая функция потерь</li> <li>• <b>Accuracy</b> – доля правильных ответов алгоритма</li> <li>• <b>MCC</b> – коэффициент корреляции Мэттьюза</li> <li>• <b>AUC ROC</b> – площадь под ROC-кривой</li> </ul> |
| <b>Размер валидационной выборки для ранней остановки (%)</b> | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию — 0,1   | Данный параметр задает долю обучающих данных, которые нужно отложить в качестве валидационного набора для ранней остановки. Используется только в случае выбора чекбокса <b>Ранняя остановка</b> .   |
| <b>Количество итераций без существенного улучшения</b>       | Ручной ввод<br>Число больше 0<br>По умолчанию — 10   | Данный параметр задает количество итераций без существенного улучшения <div style="background-color: #e0f2ff; padding: 10px; margin-top: 10px;"> <b>Примечание</b><br/>           Данный параметр можно задать только в случае выбранного <b>метода оптимизации sgd или adam</b>.         </div>   |
| <b>Погрешность оптимизации</b>                               | Ручной ввод<br>Число больше 0<br>По умолчанию — 0,0001   | Данный параметр задает допустимую погрешность для оптимизации. В случае, если потери не улучшаются на указанное значение в течение указанного в параметре Количество итераций без существенного улучшения числа, то считается, что сходимость достигнута, и обучение останавливается.  |
| <b>Инерция</b>   | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию — 0,9   | Данный параметр задает инерцию для обновления градиентного спуска <div style="background-color: #e0f2ff; padding: 10px; margin-top: 10px;"> <b>Примечание</b><br/>           Данный параметр можно задать только в случае выбранного <b>метода оптимизации sgd</b> </div>  |
| <b>Использовать инерцию Нестерова</b>                        | Чекбокс  | Выбор данного чекбокса указывает на необходимость использовать инерцию Нестерова   |

| Параметр   | Возможные значения и ограничения  | Описание  |
|--|---|---|
|  |   | <p><b>Примечание</b><br/>Данный параметр можно задать только в случае выбранного <b>метода оптимизации <code>sgd</code></b> и значения <b>Инерции &gt; 0</b></p>  |
| <b>Beta 1</b>  | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию — 0,9  | Данный параметр задает скорость экспоненциального убывания для оценок вектора первого момента при <b>методе оптимизации <code>adam</code></b>   |
| <b>Beta 2</b>  | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию — 0,999  | Данный параметр задает скорость экспоненциального убывания для оценок вектора второго момента при <b>методе оптимизации <code>adam</code></b>   |
| <b>Epsilon</b>   | Ручной ввод<br>Число больше 0<br>По умолчанию — 1e-8  | Данный параметр задает значение числовой устойчивости при <b>методе оптимизации <code>adam</code></b>   |
| <b>Стандартизация</b>                                  | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• No (по умолчанию)</li> <li>• std</li> <li>• range</li> </ul> | <p>Данный параметр отвечает за выбор метода стандартизации числовых переменных.</p> <p><b>Стандартизация</b> – преобразование числовых наблюдений с целью приведения их к некоторой общей шкале. Необходимость стандартизации вызвана тем, что разные признаки из обучающего набора могут быть представлены в разных масштабах и изменяться в разных диапазонах, что влияет на выявление некорректных зависимостей моделью.</p> <p>Предусмотрены следующие методы:</p> <ul style="list-style-type: none"> <li>• <b>no</b> — стандартизация не нужна</li> <li>• <b>std</b> — стандартное отклонение - преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li> <li>• <b>range</b> — диапазон - линейно преобразует значения переменных в диапазон [0, 1].</li> </ul> |
| <b>Максимальное количество вызовов целевой функции</b> | Ручной ввод<br>целочисленного значения больше 0<br>По молчанию - 15000  | Данный параметр задает максимальное количество вызовов функции потерь. Т.е. алгоритм выполняет итерации до сходимости ( <b>параметр Погрешность оптимизации</b> ), если количество итераций достигает значения параметра Максимальное количество итераций или заданного значения параметра вызова функции потерь.   |

| Параметр | Возможные значения и ограничения | Описание   |
|----------|----------------------------------|--|
|          |                                  | <p><b>Примечание</b><br/>Данный параметр можно задать только в случае выбранного метода оптимизации <b>Ibfgs</b></p> |

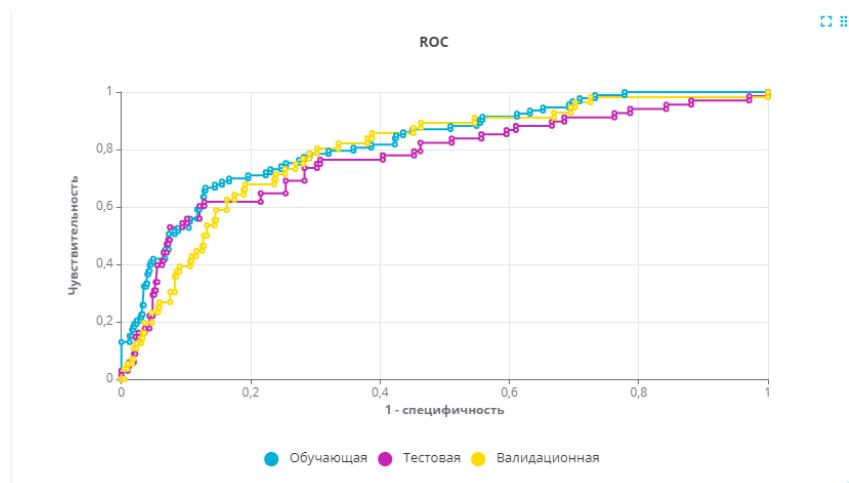
Таблица 35 Параметры узла «Нейронная сеть»

### 3.2.5.8.7.4. Результаты выполнения узла

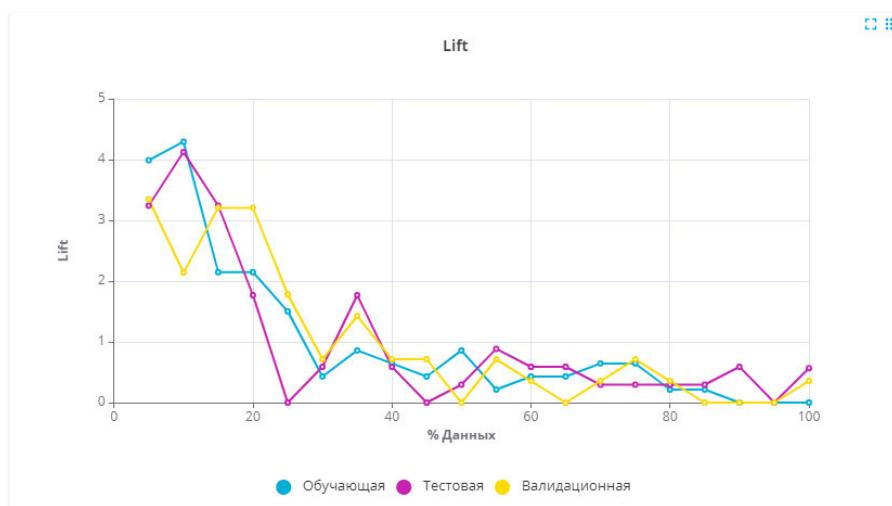
Узел «Нейронная сеть» имеет разные результаты в зависимости от решаемой задачи.

**Результаты бинарной классификации представлены следующими объектами:**

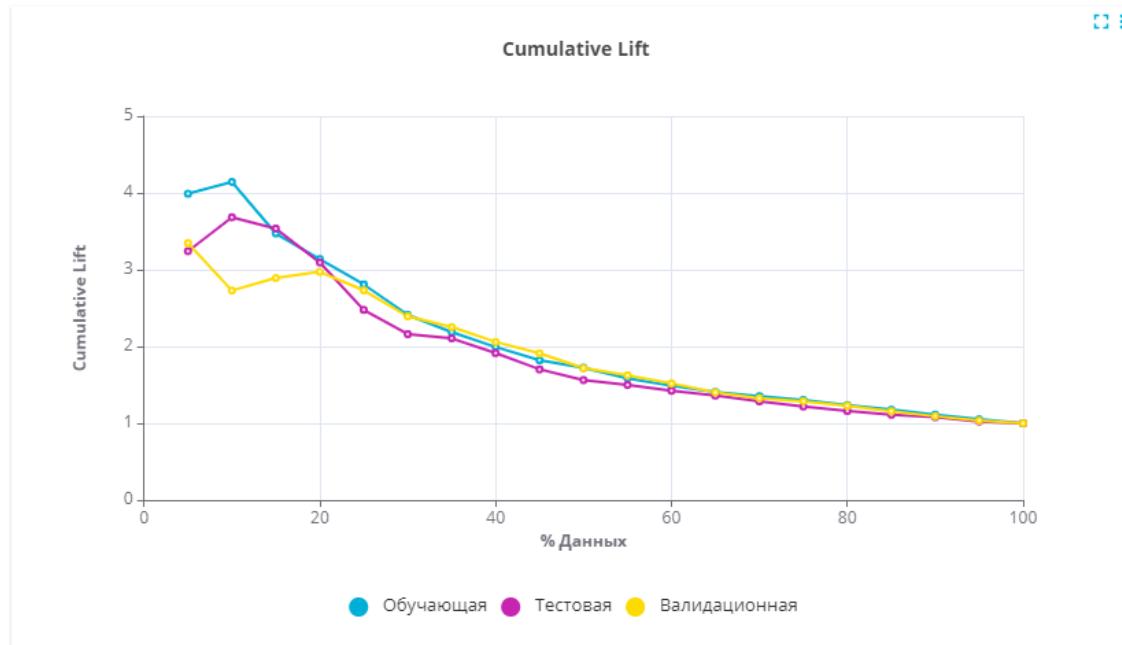
- График ROC.



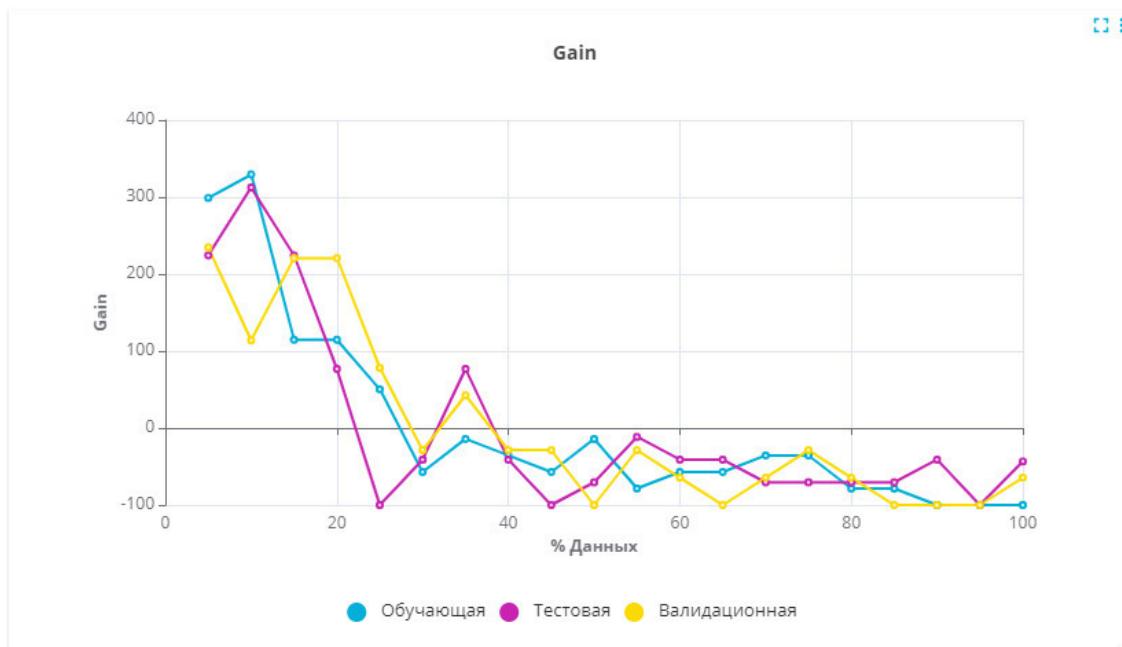
- График Lift.



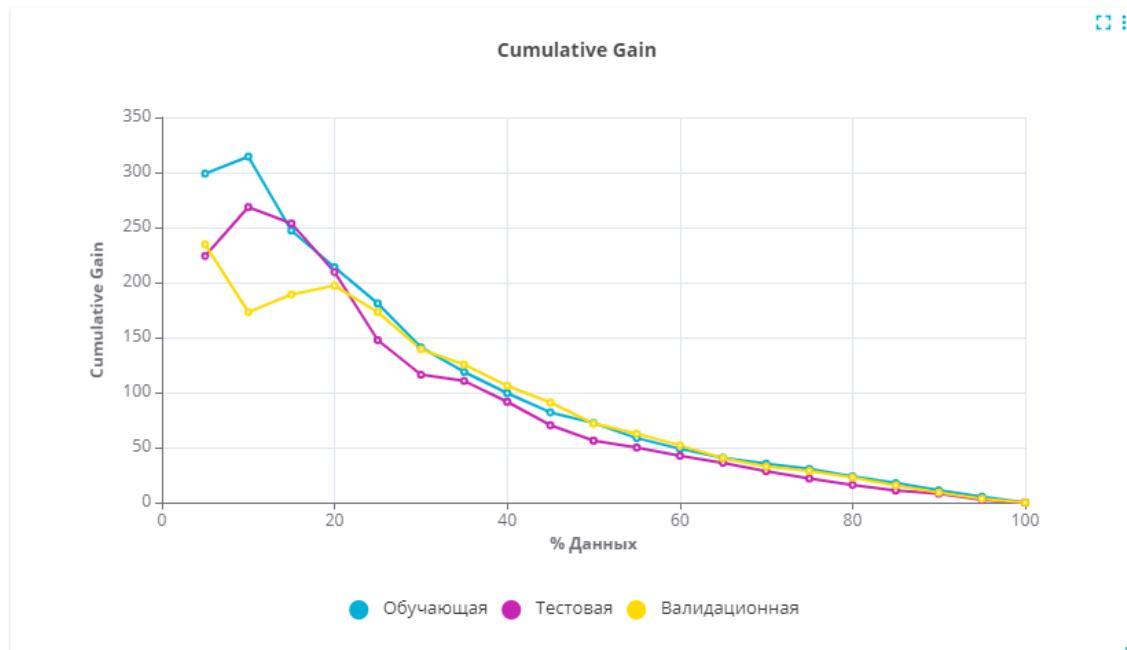
- График Cumulative Lift.



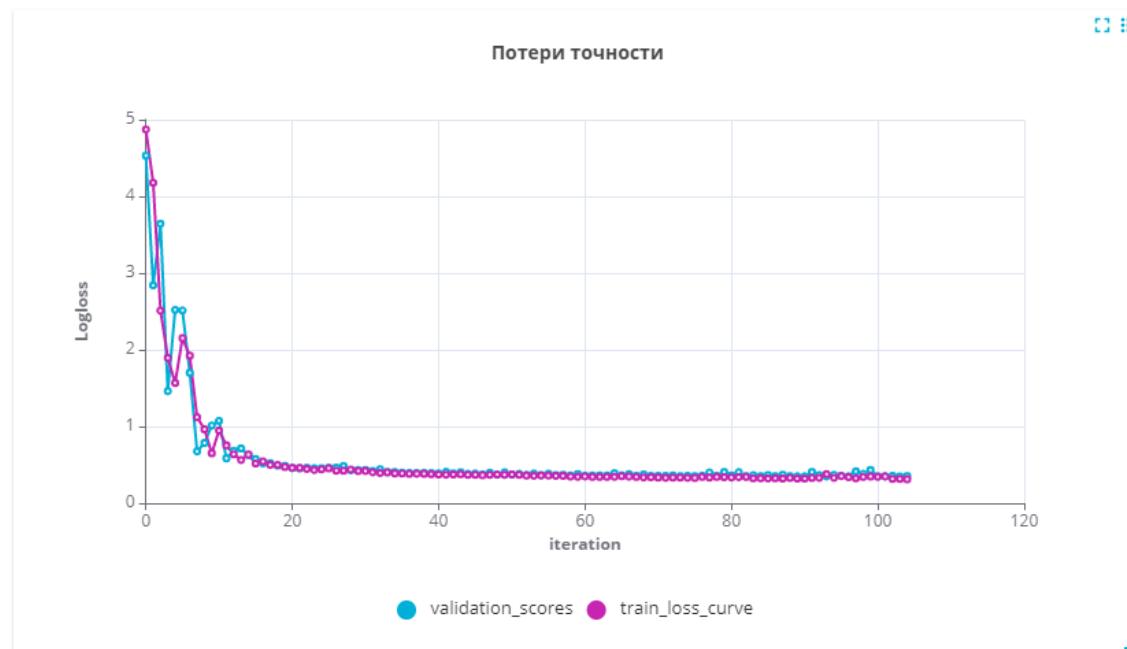
- График Gain.



- График Cumulative Gain.



- График потери точности.



- Таблица с метриками качества модели.

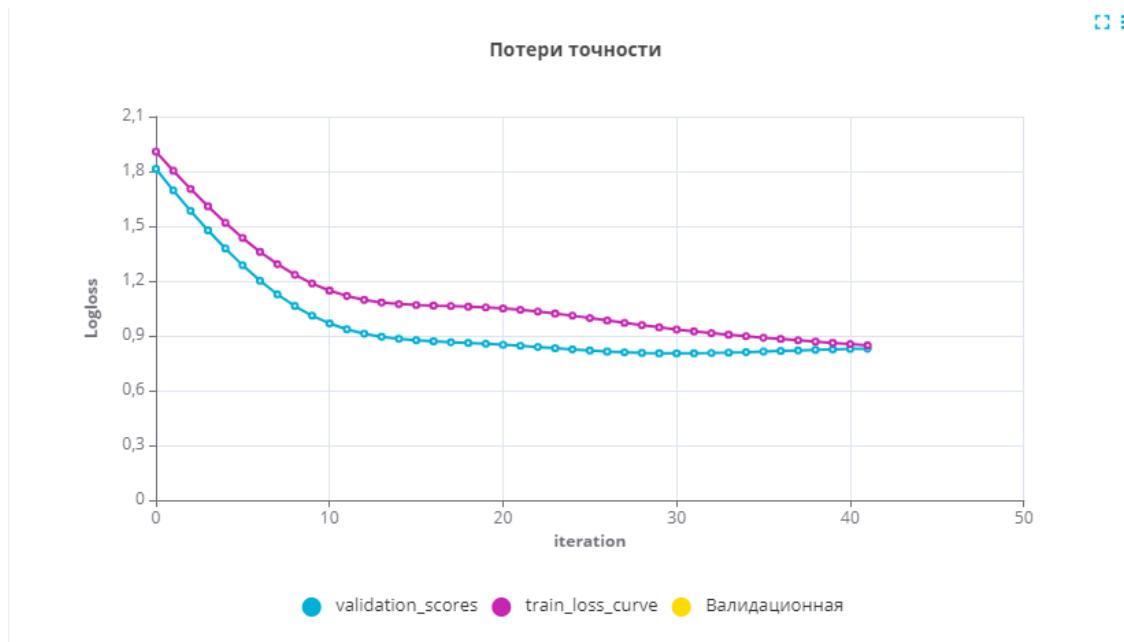
|               | AUC ROC   | gini      | log loss  | nobs      |
|---------------|-----------|-----------|-----------|-----------|
| Filter...     | Filter... | Filter... | Filter... | Filter... |
| Валидационная | 0.793     | 0.587     | 0.336     | 479       |
| Обучающая     | 0.822     | 0.644     | 0.318     | 639       |
| Тестовая      | 0.77      | 0.541     | 0.362     | 481       |

- Таблица с метриками качества модели задачи классификации.

|               | misclassification | mcc       | nobs      |
|---------------|-------------------|-----------|-----------|
| Filter...     | Filter...         | Filter... | Filter... |
| Валидационная | 0.121             | 0.168     | 479       |
| Обучающая     | 0.134             | 0.272     | 639       |
| Тестовая      | 0.149             | 0.088     | 481       |

**Результаты многоклассовой классификации представлены следующими объектами:**

- График потери точности.



- Таблица с метриками качества модели.

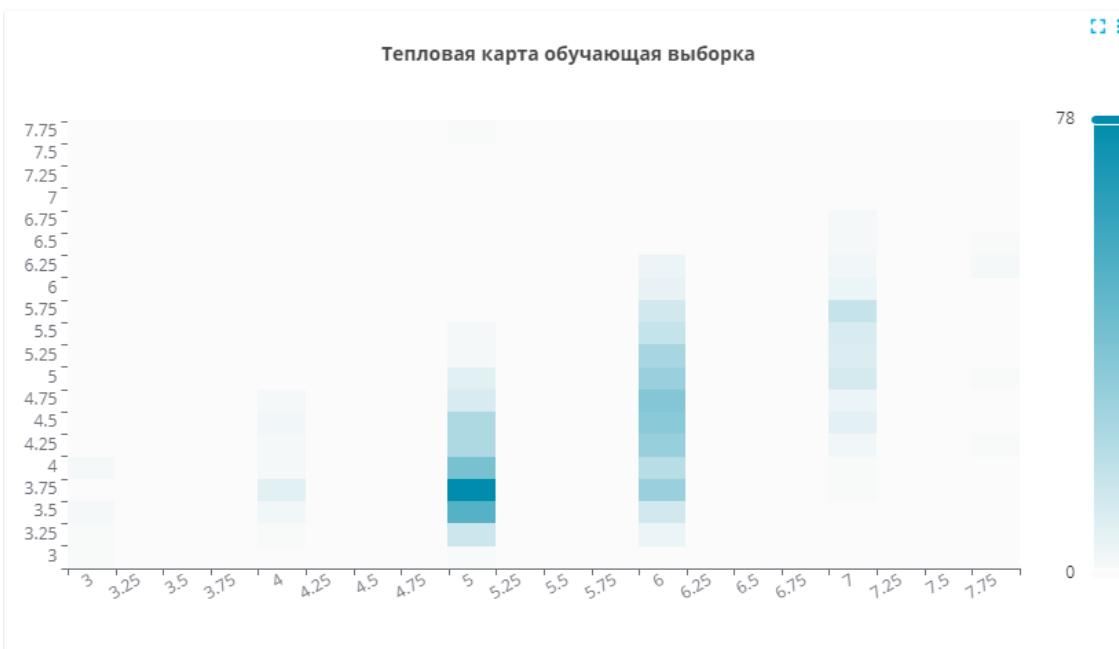
| Метрики модели |            |           |
|----------------|------------|-----------|
| ↑              | log loss ↑ | nobs ↑    |
| Filter...      | Filter...  | Filter... |
| Валидационная  | 1.052      | 35        |
| Обучающая      | 0.911      | 47        |
| Тестовая       | 1.003      | 37        |

- Таблица с метриками качества модели задачи классификации.

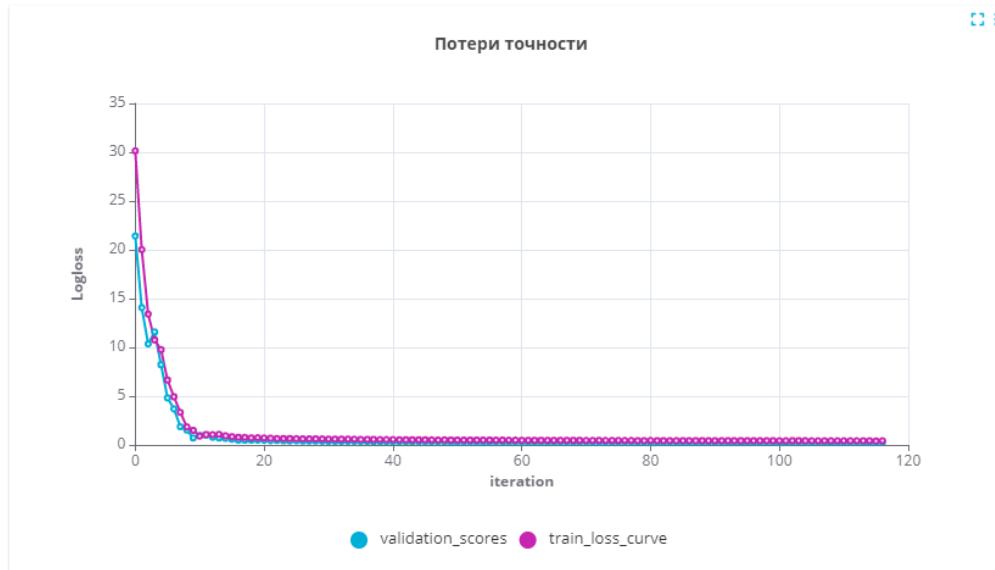
| <a href="#">↑</a>         | <a href="#">misclassification ↑</a> | <a href="#">mcc ↑</a>     | <a href="#">nobs ↑</a>    |
|---------------------------|-------------------------------------|---------------------------|---------------------------|
| <a href="#">Filter...</a> | <a href="#">Filter...</a>           | <a href="#">Filter...</a> | <a href="#">Filter...</a> |
| Валидационная             | 0.657                               | 0                         | 35                        |
| Обучающая                 | 0.489                               | 0                         | 47                        |
| Тестовая                  | 0.621                               | 0                         | 37                        |

### Результаты регрессии представлены следующими объектами:

- Тепловые карты для обучающей, валидационной и тестовой выборок.



- График потери точности.



- Таблица с метриками качества модели.

| Метрики модели     |           |           |           |           |           |           |
|--------------------|-----------|-----------|-----------|-----------|-----------|-----------|
| ↑                  | mse       | rmse      | mae       | mape      | r2        | nobs ↑    |
| Filter...          | Filter... | Filter... | Filter... | Filter... | Filter... | Filter... |
| CV (валидационная) | 0.562     | 0.745     | 0.566     | 0.104     | 0.084     | 127.8     |
| CV (обучающая)     | 0.509     | 0.705     | 0.533     | 0.098     | 0.218     | 511.2     |
| Валидационная      | 0.527     | 0.726     | 0.567     | 0.103     | 0.178     | 479       |
| Обучающая          | 0.412     | 0.642     | 0.487     | 0.09      | 0.361     | 639       |
| Тестовая           | 0.471     | 0.686     | 0.529     | 0.096     | 0.293     | 481       |

- Таблицы с метриками качества в разбиениях обучающей и валидационной выборках.

| Метрики в разбиениях обучающей выборки |           |           |           |           |           |           |           |
|--|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| Номер разбие...                        | fit_time  | mse       | rmse      | mae       | mape      | r2        | nobs      |
| Filter...                              | Filter... | Filter... | Filter... | Filter... | Filter... | Filter... | Filter... |
| 0                                      | 0.381     | 0.465     | 0.682     | 0.511     | 0.096     | 0.315     | 511       |
| 1                                      | 0.397     | 0.409     | 0.64      | 0.489     | 0.092     | 0.34      | 511       |
| 2                                      | 0.348     | 0.413     | 0.642     | 0.481     | 0.087     | 0.332     | 511       |
| 3                                      | 0.096     | 0.837     | 0.915     | 0.71      | 0.129     | -0.224    | 511       |
| 4                                      | 0.342     | 0.42      | 0.648     | 0.474     | 0.086     | 0.328     | 512       |

### 3.2.5.8.8. Узел «LDA»

**Узел «LDA»** представляет линейный дискриминантный анализ, который применяется для нахождения линейных комбинаций признаков, наилучшим образом разделяющих два или более класса объектов или событий.

**Список параметров узла** представлен в таблице ниже.

| Параметр                             | Возможные значения и ограничения  | Описание  |
|--------------------------------------|---|---|
| Название                             | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  |
| Описание                             | Ручной ввод<br>Ограничений на значение нет  | Описание узла   |
| Метод оптимизации                    | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• svd (по умолчанию)</li> <li>• lsqr</li> <li>• eigen</li> </ul> | Данный параметр задает метод, используемый для поиска матрицы объектов гиперплоскости LDA.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• svd – разложение по сингулярным значениям</li> <li>• lsqr – метод наименьших квадратов</li> <li>• eigen – разложение по собственным значениям</li> </ul> |
| Задать величину сжатия автоматически | Чекбокс   | Выбор данного чекбокса указывает на необходимость задать величину сжатия автоматически  |
| Сжатие                               | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию – 0  | Данный параметр задает значение сжатия  |
| Допустимая погрешность               | Ручной ввод<br>Неотрицательное число<br>По умолчанию – 0,0001   | Данный параметр задает допустимую погрешность для оптимизации<br>Данный параметр можно задать только в случае выбранного метода оптимизации svd   |

**Таблица 36 Параметры узла «LDA»**

#### Результаты выполнения узла:

Узел «LDA» имеет разные результаты в зависимости от решаемой задачи.

Результаты **многоклассовой классификации** представлены следующими объектами:

- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица с коэффициентами переменных.

Результаты **бинарной классификации** представлены следующими объектами:

- График ROC.
- График Lift.

- График Cumulative Lift.
- График Gain.
- График Cumulative Gain.
- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица с коэффициентами переменных.

### 3.2.5.8.9. Узел «Градиентный бустинг (XGBOOST)»

В основе **узла "XGBoost"** лежит алгоритм градиентного бустинга на деревьях поиска решений или линейных моделях. Используется для решения задач классификации и регрессии.

#### **Алгоритм работы градиентного бустинга:**

Градиентный бустинг – алгоритм машинного обучения, который строит модель предсказания в виде ансамбля слабых предсказывающих моделей (в основном Дерево решений, но также есть вариант с линейными моделями). На каждой итерации вычисляется отклонение предсказаний уже обученного ансамбля (всех предыдущих построенных моделей) на обучающей выборке. Следующая добавляемая в ансамбль модель будет сводить среднее отклонение предыдущей к минимуму.

Новые деревья добавляются в ансамбль до тех пор, пока уменьшается ошибка, либо пока не выполнится одно из правил «ранней остановки».

#### **Особенности реализации XGBoost (Extreme Gradient Boosting) в сравнении со стандартным алгоритмом градиентного бустинга:**

XGBoost поддерживает как алгоритм предварительной сортировки, так и алгоритм на основе гистограммы.

1. В данной реализации каждый лист дерева не проверяется на предмет того, что можно ли его разделить дальше. Здесь заранее задается максимальная глубина дерева, глубже которой узлы и листья не строятся. Далее после построения всех листьев до заданной глубины проверяется дает ли сплит (деление) в конкретном узле прирост информации (Information Gain). Для этого Information Gain сравнивается с заданным порогом, который настраивается экспериментальным путем. Если значение меньше порога, то листья этого узла отсекаются. Таким образом избегается переобучение
2. Используется регуляризация для избежания малого количества наблюдений в листе, что также позволяет бороться с переобучением.

**Список параметров узла** представлен в таблице ниже.

| <b>Параметр</b>                     | <b>Возможные значения и ограничения</b>   | <b>Описание</b>   | <b>Группа параметров</b> |
|-------------------------------------|---|---|--------------------------|
| <b>Название</b>                     | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  | Общий параметр           |
| <b>Описание</b>                     | Ручной ввод<br>Ограничений на значение нет  | Описание узла   | Общий параметр           |
| <b>Бустер</b>                       | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• gbtree (по умолчанию)</li> <li>• gblinear</li> <li>• DART</li> </ul>   | Данный параметр задает тип базового алгоритма для бустинга.<br>Предусмотрены следующие типы: <ul style="list-style-type: none"> <li>• gbtree – бустинг на основе деревьев</li> <li>• gblinear – бустинг на основе линейных моделей</li> <li>• DART – модификация gbtree (отбрасывает деревья, для предотвращения переобучения)</li> </ul> | Общий параметр           |
| <b>Количество оценочных функций</b> | Ручной ввод<br>Число больше 0<br>По умолчанию — 100   | Данный параметр задает число итераций градиентного бустинга   | Общий параметр           |
| <b>Скорость обучения</b>            | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию — 0,3  | Данный параметр задает скорость обучения модели и контролирует, с каким весом предсказания каждой следующей модели суммируются с предсказаниями ансамбля.<br>Значение по умолчанию (0,3) является слишком большим, обычно хорошо работают значения меньше 0.1   | Общий параметр           |
| <b>Цель обучения для регрессии</b>  | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Квадрат ошибки (по умолчанию)</li> <li>• Функция потерь Хьюбера</li> <li>• Пуассоновская регрессия</li> <li>• Регрессия Твида</li> </ul> | Данный параметр задает используемую при обучении функцию потерь. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• Квадрат ошибки – квадратичная функция потерь</li> <li>• Функция потерь Хьюбера – функция квадратична для малых значений остатка (разница между</li> </ul>  | Общий параметр           |

| Параметр | Возможные значения и ограничения                                  | Описание   | Группа параметров |
|----------|---|--|-------------------|
|          | <ul style="list-style-type: none"> <li>Гамма регрессия</li> </ul> | <p>наблюдаемым и предсказанным значением), и линейна для больших значений остатка</p> <ul style="list-style-type: none"> <li>Пуассоновская регрессия – предназначена для прогнозирования счетчиков (неотрицательных целых чисел) (например, количество дождевых явлений в год или количество событий прерывания производства в год). Поэтому использовать данную функцию следует в случае соответствия следующим условиям:</li> <li>переменная ответа имеет распределение Пуассона,</li> <li>метки не должны быть отрицательными</li> <li>бессмысленно использовать для дробных чисел</li> <li>Регрессия Твиди – предназначена для прогнозирования целевой переменной, имеющей распределение Твиди (например, общее количество осадков в год или общее время прерывания в год)</li> <li>Гамма регрессия – предназначена для прогнозирования целевой переменной,</li> </ul> |                   |

| <b>Параметр</b>                        | <b>Возможные значения и ограничения</b>   | <b>Описание</b>   | <b>Группа параметров</b>  |
|--|---|---|---|
|  |   | имеющей гамма-распределение (например, количество осадков на одно событие или продолжительность прерывания)   |   |
| <b>Цель обучения для классификации</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Бинарная логистическая регрессия (по умолчанию)</li> <li>• Бинарная с hinge loss</li> <li>• Мультиклассовая с softprob</li> <li>• Мультиклассовая с softmax</li> </ul> | Данный параметр задает используемую при обучении функцию потерь. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• Бинарная логистическая регрессия – возвращает прогнозируемую вероятность (не класс)</li> <li>• Бинарная с hinge loss – кусочно-линейная функция потерь для бинарной классификации. Возвращает принадлежность одному из двух классов – 0 или 1</li> <li>• Мультиклассовая с softprob – функция softprob для мультиклассовой классификации, возвращает матрицу со значениями вероятности каждого класса</li> <li>• Мультиклассовая с softmax – функция softmax для мультиклассовой классификации, возвращает класс с максимальной вероятностью принадлежности</li> </ul> | Общий параметр  |
| <b>Дисперсия распределения Твиди</b>   | Ручной ввод<br>Число больше 1 и меньше 2<br>По умолчанию — 1,5  | Данный параметр используется для управления дисперсией распределения Твиди. Значение ближе к 2 переходит в гамма-   | Параметр актуален при заданных Цели обучения для регрессии = Регрессия Твиди и Метриках для |

| <b>Параметр</b>  | <b>Возможные значения и ограничения</b>   | <b>Описание</b>  | <b>Группа параметров</b>                                |
|--|---|--|---|
|  |   | распределение, значение ближе к 1 – распределение Пуассона.  | валидации регрессии, связанными с распределением Твиди. |
| <b>Размер (%) валидационной выборки для ранней остановки</b> | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию – 0,1  | Данный параметр задает размер (%) валидационной выборки для ранней остановки   | Общий параметр  |
| <b>Количество итераций до ранней остановки</b>               | Ручной ввод<br>По умолчанию – 0   | Данный параметр задает количество итераций до ранней остановки   | Общий параметр  |
| <b>Метрика для валидации регрессии</b>                       | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• RMSE (по умолчанию)</li> <li>• MAE</li> <li>• MAPE</li> <li>• poisson nloglik</li> <li>• gamma nloglik</li> <li>• gamma deviance</li> <li>• tweedie nloglik</li> </ul> | Данный параметр задает метрику качества на валидационной выборке. Предусмотрены следующие метрики: <ul style="list-style-type: none"> <li>• RMSE – Среднеквадратическая ошибка</li> <li>• MAE – Средняя абсолютная ошибка</li> <li>• MAPE – Средняя абсолютная ошибка в процентах</li> <li>• poisson nloglik – отрицательная логарифмическая функция правдоподобия для регрессии Пуассона</li> <li>• gamma nloglik – отрицательная логарифмическая функция правдоподобия для гамма-регрессии</li> <li>• gamma deviance – остаточное отклонение для гамма-регрессии</li> <li>• tweedie nloglik – отрицательная логарифмическая функция правдоподобия для регрессии Твиди</li> </ul> | Общие параметры   |
| <b>Метрика для валидации</b>                                 | Раскрывающийся список со  | Данный параметр задает метрику качества на валидационной выборке.  | Общие параметры   |

| Параметр                        | Возможные значения и ограничения   | Описание  | Группа параметров  |
|---------------------------------|--|---|--|
| <b>классификации</b>            | следующими значениями: <ul style="list-style-type: none"> <li>• logloss (по умолчанию)</li> <li>• error</li> <li>• merror</li> <li>• mlogloss</li> <li>• auc</li> <li>• aucpr</li> </ul> | Предусмотрены следующие метрики: <ul style="list-style-type: none"> <li>• logloss – логистическая функция ошибки</li> <li>• error – частота ошибок бинарной классификации, рассчитывается как неправильно классифицированные объекты/все объекты. При прогнозировании положительными экземплярами будут считаться наблюдения со значением прогноза больше 0,5, остальные – как отрицательные</li> <li>• merror – частота ошибок мультиклассовой классификации, рассчитывается как неправильно классифицированные объекты/все объекты</li> <li>• mlogloss – мультиклассовая логистическая функция ошибки</li> <li>• auc – количественная интерпретация кривой ошибок, площадь, ограниченная ROC-кривой и осью доли ложных положительных классификаций</li> <li>• aucpr – площадь под PR-кривой (Precision-Recall curve)</li> </ul> |  |
| <b>Cutoff для метрики error</b> | Ручной ввод<br>По умолчанию – 0,5  | Данный параметр задает порог отсечения, чтобы относить новые примеры к одному из двух классов (задача бинарной классификации).  | Общие параметры<br>Актуален при выбранной метрике error для валидации классификации. |

| Параметр                                      | Возможные значения и ограничения  | Описание  | Группа параметров   |
|---|---|---|---|
| <b>Способ определения важности переменных</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• gain (по умолчанию)</li> <li>• weight</li> <li>• cover</li> <li>• total gain</li> <li>• total cover</li> </ul> | Данный параметр определяет метод оценки важности переменных. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• gain – средний выигрыш по всем разделениям (splits), в которых используется переменная</li> <li>• weight – количество раз, когда переменная используется для разделения данных по всем деревьям</li> <li>• cover – среднее количество наблюдений для каждой фичи по всем разделениям, в которых используется переменная</li> <li>• total gain – общий выигрыш по всем разделениям, в которых используется переменная</li> <li>• total cover – общее количество наблюдений для каждой фичи всех разделений, в которых используется переменная</li> </ul> | Параметры для Tree Booster (Бустер = gbtree). Для Бустер = gblinear значение параметра = 'weight' |
| <b>Количество параллельных потоков</b>        | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 0   | Данный параметр задает количество параллельных потоков, используемых для запуска xgboost<br>0 означает использование всех доступных потоков (CPU)   | Общие параметры   |
| <b>Seed</b>                                   | Ручной ввод числового значения<br>По умолчанию — 42   | Начальное числовое значение для генератора случайных чисел  | Общие параметры   |
| <b>Метод построения дерева</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• auto (по умолчанию)</li> </ul>   | Данный параметр задает метод построения дерева. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• auto – для небольшого набора</li> </ul>  | Параметры для Tree Booster и DART (Бустер = gbtree или dart)                                      |

| Параметр                               | Возможные значения и ограничения   | Описание  | Группа параметров  |
|--|--|---|--|
|  | <ul style="list-style-type: none"> <li>exact</li> <li>approx</li> <li>hist</li> </ul>  | <p>данных будет использован exact, для большого набора данных – approx</p> <ul style="list-style-type: none"> <li>exact – жадный алгоритм, который перебирает все наблюдения входного набора в ходе процедуры поиска разделения. Данный метод более точен среди других жадных методов, но медленнее в вычислительной производительности</li> <li>approx – приближенный жадный алгоритм, который использует quantile sketch (квантильные наброски) и gradient histogram (приближенные гистограммы статистики градиента)</li> <li>hist – более быстрый приближенный жадный алгоритм, оптимизированный для histogram.</li> </ul> |  |
| <b>Метод добавления узлов к дереву</b> | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>depthwise (по умолчанию)</li> <li>lossguide</li> </ul> | <p>Данный параметр задает способ добавления новых узлов к дереву. Предусмотрены следующие методы:</p> <ul style="list-style-type: none"> <li>depthwise – разделение в узлах, ближайших к корню</li> <li>lossguide – разделение на узлы с наибольшим изменением значения функции потерь</li> </ul>   | <p>Параметры для Tree Booster и DART (Бустер = gbtree или dart)<br/>Поддерживается только если метод построения дерева выбран как hist</p> |

| <b>Параметр</b>  | <b>Возможные значения и ограничения</b>                                | <b>Описание</b>  | <b>Группа параметров</b>   |
|--|--|--|--|
| <b>Количество параллельных деревьев</b>                          | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 1            | Данный параметр задает количество параллельных деревьев, построенных на каждой итерации  | Параметры для Tree Booster и DART (Бустер = gbtree или dart)   |
| <b>Максимальная глубина дерева</b>                               | Ручной ввод<br>Число больше 0<br>По умолчанию – 6                      | Данный параметр задает максимальную глубину дерева   | Параметры для Tree Booster и DART (Бустер = gbtree или dart)   |
| <b>Минимальное снижение потери для разбиения</b>                 | Ручной ввод<br>Число больше или равно 0, float<br>По умолчанию – 0     | Данный параметр задает значение минимального уменьшения функции потерь для разбиения   | Параметры для Tree Booster и DART (Бустер = gbtree или dart)   |
| <b>Максимальное количество листов</b>                            | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 0            | Данный параметр задает максимальное количество листов.   | Параметры для Tree Booster и DART (Бустер = gbtree или dart)<br>Актуален, если выбран Метод добавления узлов к дереву = lossguide  |
| <b>Максимальное количество бинов для интервальных переменных</b> | Ручной ввод<br>Число больше 0<br>По умолчанию – 256                    | Данный параметр задает максимальное количество бинов для интервальных переменных.<br>Увеличение этого числа повышает оптимальное разделение за счет увеличения времени вычислений.                                 | Параметры для Tree Booster и DART (Бустер = gbtree или dart)<br>Актуален, если выбран метод построения дерева = hist   |
| <b>Относительное количество бинов (sketch_eps)</b>               | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию – 0,03        | Данный параметр задает относительное количество бинов, примерно соответствующее $1 / \text{sketch\_eps}$ . По сравнению с прямым выбором количества бинов дает теоретическую гарантию точности набросков (sketch). | Параметры для Tree Booster и DART (Бустер = gbtree или dart)<br>Актуален, если выбран метод построения дерева = approx   |
| <b>Соотношение колонок для каждого дерева</b>                    | Ручной ввод<br>Число больше 0 и меньше или равно 1<br>По умолчанию – 1 | Данный параметр задает долю переменных, используемых на каждой итерации (при построении каждого дерева).<br>Подвыборка происходит один раз для каждого построенного дерева   | Параметры для Tree Booster и DART (Бустер = gbtree или dart)<br>Данные параметры работают кумулятивно.<br>Например, комбинация <code>{'colsample_bytree':0,5, 'colsample_bylevel':0,5, 'colsample_bynode':0,5}</code> с 64 функциями оставит 8 функций |

| <b>Параметр</b>                                  | <b>Возможные значения и ограничения</b>                                | <b>Описание</b>  | <b>Группа параметров</b>                                     |
|--|--|--|--|
|  |  |  | на выбор при каждом разбиении.                               |
| <b>Соотношение колонок для каждого уровня</b>    | Ручной ввод<br>Число больше 0 и меньше или равно 1<br>По умолчанию – 1 | Данный параметр задает долю подвыборки признаков, которые будут использованы для обучения каждого уровня. Подвыборка происходит один раз для каждого нового уровня глубины, достигнутого в дереве. Колонки выбираются из набора колонок, выбранных для текущего дерева.  |  |
| <b>Соотношение колонок для каждого разбиения</b> | Ручной ввод<br>Число больше 0 и меньше или равно 1<br>По умолчанию – 1 | Данный параметр задает долю подвыборки признаков, которые будут использованы для каждого разделения (узла). Подвыборка происходит каждый раз, когда оценивается новое разбиение. Колонки (столбцы) выбираются из набора колонок, выбранного для текущего уровня  |  |
| <b>Минимальный вес для потомка</b>               | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 0            | Данный параметр задает минимальный вес потомка, необходимый для разделения. Если шаг разделения дерева приводит к листу с суммой весов меньше, чем заданное данным параметром значение, то процесс построения откажется от дальнейшего разделения. В задаче линейной регрессии это соответствует минимальному количеству наблюдений, которые должны быть в каждом узле. Чем больше значение, тем более консервативен алгоритм. | Параметры для Tree Booster и DART (Бустер = gbtree или dart) |
| <b>Максимальный шаг в листе</b>                  | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 0            | Данный параметр задает максимальный шаг обновления для каждого выходного листа.  | Параметры для Tree Booster и DART (Бустер = gbtree или dart) |

| Параметр  | Возможные значения и ограничения  | Описание   | Группа параметров  |
|---|---|--|--|
|   |   | Если значение равно 0, это означает, что ограничения отсутствуют. Данный параметр может помочь при логистической регрессии, когда классы несбалансированы.   |  |
| <b>Соотношение случайной подвыборки в обучающей выборке</b> | Ручной ввод<br>Число больше 0 или меньше или равно 1<br>По умолчанию – 1  | Данный параметр задает долю объектов обучающей выборки, используемых на каждой итерации.<br>При значении равном 0,5 XGBoost будет случайным образом отбирать половину обучающих данных перед обучением (growing) дерева, что предотвращает переобучение.<br>Подвыборка будет происходить один раз в каждой итерации бустинга   | Бэггинг  |
| <b>Метод сэмплинга</b>                                      | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"><li>• uniform (по умолчанию)</li><li>• gradient based</li></ul> | Данный параметр задает метод сэмплинга.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"><li>• uniform – каждое наблюдение имеет равную вероятность быть выбранным. Для хороших результатов необходимо выбрать соотношение случайной подвыборки в обучающей выборке = 0,5</li><li>• gradient based – вероятность выбора наблюдения пропорциональна регуляризованному абсолютному значению градиентов</li></ul> | Бэггинг  |
| <b>L1 регуляризатор</b>                                     | Ручной ввод<br>По умолчанию – 0   | Регуляризация добавляет ограничения алгоритма относительно аспектов модели, которые не   | Регуляризация<br>Параметры для Linear Booster и Tree Booster (Бустер = |

| <b>Параметр</b>                   | <b>Возможные значения и ограничения</b>  | <b>Описание</b>  | <b>Группа параметров</b>   |
|-----------------------------------|--|--|--|
|                                   |  | зависят от данных для обучения. Регуляризация обычно используется, чтобы избежать переобучения. Регуляризация L1 применяется для получения максимально разреженной модели  | gblinear или Бустер = gbtree)  |
| <b>L2 регуляризатор</b>           | Ручной ввод<br>По умолчанию – 1  | Регуляризация L2 ограничивает чрезмерный рост какой-либо отдельной координаты весового вектора. Регуляризация L2 полезна в том случае, если целью является создание модели, имеющей в целом малые значения веса.   | Регуляризация<br>Параметры для Linear Booster и Tree Booster (Бустер = gblinear или Бустер = gbtree) |
| <b>Тип сэмплинга</b>              | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• uniform (по умолчанию)</li> <li>• weighted</li> </ul> | Данный параметр задает тип алгоритма сэмплинга. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• uniform –деревья имеют одинаковую вероятность быть отброшенными</li> <li>• weighted – отбрасываемые деревья выбираются пропорционально весу</li> </ul>   | Параметры для DART (Бустер = dart)   |
| <b>Тип алгоритма нормализации</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Дерево (по умолчанию)</li> <li>• Лес</li> </ul>       | Данный параметр задает тип алгоритма нормализации. Предусмотрены следующие варианты: <ul style="list-style-type: none"> <li>• Дерево – новые деревья имеют такой же вес, как и каждое из отбрасываемых деревьев. Вес новых деревьев равен <math>1/(k + \text{скорость обучения})</math>. Отбрасываемые деревья масштабируются с коэффициентом <math>k/(k + \text{скорость обучения})</math></li> </ul> | Параметры для DART (Бустер = dart)   |

| Параметр  | Возможные значения и ограничения  | Описание  | Группа параметров                                |
|---|---|---|--|
|   |   | <ul style="list-style-type: none"> <li>Лес – новые деревья имеют тот же вес, что и сумма отбрасываемых деревьев (леса). Вес новых деревьев равен <math>1/(k + \text{скорость обучения})</math>. Отбрасываемые деревья масштабируются с коэффициентом <math>k/(k + \text{скорость обучения})</math></li> </ul> |  |
| <b>Доля отбрасываемых деревьев</b>                | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию – 0  | Данный параметр задает долю отбрасываемых деревьев  | Параметры для DART (Бустер = dart)               |
| <b>Отбрасывать хотя бы одно дерево</b>            | Чекбокс   | Выбор данного чекбокса указывает на то, что по крайней мере одно дерево будет отбрасываться   | Параметры для DART (Бустер = dart)               |
| <b>Вероятность пропуска отбрасывания деревьев</b> | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию – 0  | Данный параметр задает вероятность пропуска процедуры отбрасывания деревьев.<br>Если процедура отбрасывания деревьев пропущена, новые деревья добавляются так же, как и gbtree.   | Параметры для DART (Бустер = dart)               |
| <b>Алгоритм</b>                                   | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• shotgun (по умолчанию)</li> <li>• coord descent</li> </ul> | Данный параметр задает алгоритм для обучения линейной модели.<br>Предусмотрены следующие алгоритмы: <ul style="list-style-type: none"> <li>• shotgun – алгоритм спуска с параллельными координатами</li> <li>• coord descent – алгоритм спуска по обычным координатам</li> </ul>                              | Параметры для Linear Booster (Бустер = gblinear) |
| <b>Метод выбора переменных</b>                    | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• cyclic (по умолчанию)</li> <li>• shuffle</li> </ul>        | Данный параметр задает метод выбора переменных.<br>Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• cyclic – детерминированны</li> </ul>  | Параметры для Linear Booster (Бустер = gblinear) |

| Параметр  | Возможные значения и ограничения  | Описание  | Группа параметров   |
|---|---|---|---|
|   | <ul style="list-style-type: none"> <li>random</li> <li>greedy</li> <li>thrifty</li> </ul> | <p>• выбор путем циклического перебора признаков по одному</p> <p>• shuffle – похоже на cyclic, но со случайным перемешиванием функций перед каждым обновлением</p> <p>• random – случайный (с заменой) селектор координат</p> <p>• greedy – выбирается координата с наибольшей величиной градиента.</p> <p>• thrifty – перед циклическими обновлениями функции переупорядочиваются по убыванию величины их одномерных изменений веса</p> |   |
| <b>Количество отбираемых переменных для жадных алгоритмов</b> | Числовое значение<br>По умолчанию – 0   | Данный параметр задает количество отбираемых переменных для жадных алгоритмов.<br>Значение, равное 0, означает использование всех функций.  | Параметры для Linear Booster (Бустер = gblinear)<br>Актуален, если выбраны метод выбора переменных = greedy или thrifty |

**Таблица 37 Параметры узла «Градиентный бустинг (XGBOOST)»**

#### **Результаты выполнения узла:**

Узел «Градиентный бустинг (XGBOOST)» имеет разные результаты в зависимости от решаемой задачи.

Результаты **регрессии** представлены следующими объектами:

- Термальные карты обучающей, валидационной и тестовой выборок.
- Таблица с метриками качества модели.
- Таблица со списком переменных, сорттированных по важности.

Результаты **многоклассовой классификации** представлены следующими объектами:

- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица со списком переменных, сортированных по важности.

Результаты **бинарной классификации** представлены следующими объектами:

- График ROC.
- График Lift.
- График Cumulative Lift.
- График Gain.
- График Cumulative Gain.
- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица со списком переменных, сортированных по важности.

### 3.2.5.8.10. Узел «Градиентный бустинг (LightGBM)»

В основе **узла «LightGBM»** лежит реализация алгоритма градиентного бустинга на деревьях поиска решений, который включает в себя две ключевые идеи: **GOSS** и **EFB**.

Используется для решения задач классификации и регрессии.

#### **Алгоритм работы градиентного бустинга:**

Градиентный бустинг – алгоритм машинного обучения, который строит модель предсказания в виде ансамбля слабых предсказывающих моделей (в основном Дерево решений). На каждой итерации вычисляется отклонение предсказаний уже обученного ансамбля на обучающей выборке. Следующая добавляемая в ансамбль модель будет сводить среднее отклонение предыдущей к минимуму. Новые деревья добавляются в ансамбль до тех пор, пока ошибка уменьшается, либо пока не выполняется одно из правил «ранней остановки».

#### **Особенности реализации LightGBM:**

- Используется **алгоритм роста дерева по листьям** (в каждом узле рассчитывается дает ли разделение листа прирост информации –  $\text{inf gain}$ ), поэтому также необходимо ограничивать построение дерева максимальной глубиной для предотвращения переобучения.
- LightGBM использует **алгоритмы на основе гистограмм**, которые разбивают значения непрерывных признаков (атрибутов) на дискретные ячейки.
- **Градиентная односторонняя выборка (Gradient-based One-Side Sampling - GOSS)**

Для расчета берутся не все наблюдения, а лишь с большой ошибкой и небольшая часть наблюдений с маленькой ошибкой. Данный подход позволяет уменьшить расчетную часть алгоритма, за счет чего растет скорость расчета.

При расчете прироста информации GOSS вводит постоянный множитель для экземпляров данных с небольшими градиентами. Т. е. GOSS сначала сортирует наблюдения в соответствии с абсолютными значениями их градиентов и выбирает верхние экземпляры  $a * 100\%$  ( $a$  - **коэффициент сохранения большого градиента**). Затем он случайным образом отбирает наблюдения  $b * 100\%$  ( $b$  - **коэффициент сохранения большого градиента**) из остальных данных. После этого GOSS усиливает выборочные данные с небольшими градиентами на константу  $(1-a)/b$  при расчете прироста информации. Это позволяет учитывать исходное распределение данных.

- Встроенная кодировка категориальных переменных:**

Считается сумма градиентов по каждой категории и категории объединяются в группы.

- Объединение взаимоисключающих признаков (Exclusive Feature Bundling - EFB)**

Данный подход подразумевает объединение разреженных (в основном нулевых взаимоисключающих признаков, таких как категориальные переменные, закодированные унитарным кодированием.

Переменные, которые не принимают одновременно ненулевые значения, можно «связать». Уменьшаем количество переменных, объединяя взаимоисключающие в 1 переменную, увеличивается скорость расчета.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения   | Описание   | Группа параметров |
|-----------------|--|--|-------------------|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет   | Название узла, которое будет отображаться в интерфейсе   | Общий параметр    |
| <b>Описание</b> | Ручной ввод<br>Ограничений на значение нет   | Описание узла  | Общий параметр    |
| <b>Бустер</b>   | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• gbtree (по умолчанию)</li> <li>• DART</li> <li>• GOSS</li> <li>• random forest</li> </ul> | Данный параметр задает тип базового алгоритма для бустинга.<br>Предусмотрены следующие типы: <ul style="list-style-type: none"> <li>• gbtree – бустинг на основе деревьев</li> <li>• DART – модификация gbtree (отбрасывает деревья, для предотвращения переобучения)</li> </ul> | Общий параметр    |

| Параметр                            | Возможные значения и ограничения  | Описание   | Группа параметров |
|-------------------------------------|---|--|-------------------|
|                                     |   | <ul style="list-style-type: none"> <li>• GOSS – используются только те экземпляры, которые приводят к большому градиенту ошибки, для обновления модели и удаления остальных экземпляров.</li> <li>• random forest – Случайный лес</li> </ul>   |                   |
| <b>Линейные модели в листах</b>     | Чекбокс   | Выбор данного чекбокса указывает на необходимость обучения линейных моделей в листах каждого дерева  | Общий параметр    |
| <b>Количество оценочных функций</b> | Ручной ввод<br>Целочисленное значение больше или равно 0<br>По умолчанию - 100  | Данный параметр задает число итераций градиентного бустинга  | Общий параметр    |
| <b>Скорость обучения</b>            | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию — 0,3  | Данный параметр задает скорость обучения модели и контролирует, с каким весом предсказания каждой следующей модели суммируются с предсказаниями ансамбля.  | Общий параметр    |
| <b>Цель обучения для регрессии</b>  | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Среднеквадратическая ошибка (по умолчанию)</li> <li>• Средняя абсолютная ошибка</li> <li>• Функция потерь Хьюбера</li> <li>• Регрессия Твиди</li> <li>• Гамма регрессия</li> <li>• Пуассоновская регрессия</li> <li>• Квантильная регрессия</li> <li>• MAPE</li> </ul> | Данный параметр задает используемую при обучении функцию потерь. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• Среднеквадратическая ошибка</li> <li>• Средняя абсолютная ошибка</li> <li>• Функция потерь Хьюбера – функция квадратична для малых значений остатка (разница между наблюдаемым и предсказанным значением), и линейна для</li> </ul> | Общий параметр    |

| Параметр | Возможные значения и ограничения | Описание  | Группа параметров |
|----------|----------------------------------|---|-------------------|
|          |                                  | <p>больших значений остатка</p> <ul style="list-style-type: none"> <li>• Регрессия Твида – предназначена для прогнозирования целевой переменной, имеющей распределение Твида (например, общее количество осадков в год или общее время прерывания в год).</li> <li>• Гамма регрессия – предназначена для прогнозирования целевой переменной, имеющей гамма-распределение (например, количество осадков на одно событие или продолжительность прерывания)</li> <li>• Пуассоновская регрессия – предназначена для прогнозирования счетчиков (неотрицательных целых чисел) (например, количество дождевых явлений в год или количество событий прерывания производства в год).</li> <li>• Квантильная регрессия – предназначена для моделирования взаимосвязи между набором предикторов и</li> </ul> |                   |

| Параметр                               | Возможные значения и ограничения   | Описание  | Группа параметров |
|--|--|---|-------------------|
|  |  | <p>определенными процентилями (квантилями) целевой переменной</p> <ul style="list-style-type: none"> <li>• MAPE - средняя абсолютная ошибка в процентах</li> </ul>  |                   |
| <b>Цель обучения для классификации</b> | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>• Бинарная логистическая регрессия (по умолчанию)</li> <li>• Мультиклассовая с softmax</li> </ul>  | <p>Данный параметр задает используемую при обучении функцию потерь. Предусмотрены следующие:</p> <ul style="list-style-type: none"> <li>• Бинарная логистическая регрессия – возвращает прогнозируемую вероятность (не класс)</li> <li>• Мультиклассовая с softmax – функция softmax для мультиклассовой классификации, возвращает класс с максимальной вероятностью принадлежности</li> </ul>  | Общий параметр    |
| <b>Метрика для валидации регрессии</b> | <p>Раскрывающийся список со следующими значениями:</p> <p>Среднеквадратическая ошибка (по умолчанию)<br/>Средняя абсолютная ошибка<br/>RMSE<br/>MAPE<br/>Пуассоновская регрессия<br/>Функция потерь Хьюбера<br/>Квантильная регрессия<br/>Гамма регрессия<br/>Гамма-отклонение<br/>Регрессия Твида</p> | <p>Данный параметр задает метрику качества на валидационной выборке. Предусмотрены следующие метрики:</p> <ul style="list-style-type: none"> <li>• Среднеквадратическая ошибка (по умолчанию)</li> <li>• Средняя абсолютная ошибка</li> <li>• RMSE</li> <li>• MAPE</li> <li>• Пуассоновская регрессия</li> <li>• Функция потерь Хьюбера</li> <li>• Квантильная регрессия</li> <li>• Гамма регрессия</li> <li>• Гамма-отклонение</li> <li>• Регрессия Твида</li> </ul> | Общий параметр    |

| Параметр                                   | Возможные значения и ограничения  | Описание   | Группа параметров |
|--|---|--|-------------------|
| <b>Метрика для валидации классификации</b> | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>• Бинарная логистическая регрессия (по умолчанию)</li> <li>• Бинарная ошибка</li> <li>• Мультиклассовая логистическая регрессия</li> <li>• Мультиклассовая ошибка</li> <li>• AUC</li> </ul> | <p>Данный параметр задает метрику качества на валидационной выборке. Предусмотрены следующие метрики:</p> <ul style="list-style-type: none"> <li>• Бинарная логистическая регрессия – логистическая функция ошибки</li> <li>• Бинарная ошибка – частота ошибок бинарной классификации, рассчитывается как неправильно классифицированные объекты/все объекты. При прогнозировании положительными экземплярами будут считаться наблюдения со значением прогноза больше 0,5, остальные – как отрицательные</li> <li>• Мультиклассовая логистическая регрессия – мультиклассовая логистическая функция ошибки</li> <li>• Мультиклассовая ошибка – частота ошибок мультиклассовой классификации, рассчитывается как неправильно классифицированные объекты/все объекты</li> <li>• AUC – количественная интерпретация кривой ошибок, площадь, ограниченная ROC-кривой и осью доли ложных положительных классификаций</li> </ul> | Общий параметр    |

| <b>Параметр</b>  | <b>Возможные значения и ограничения</b>                                 | <b>Описание</b>  | <b>Группа параметров</b> |
|--|---|--|--------------------------|
| <b>Количество итераций до ранней остановки</b>               | Ручной ввод<br>По умолчанию - 0   | Данный параметр задает количество итераций до ранней остановки   | Общий параметр           |
| <b>Размер (%) валидационной выборки для ранней остановки</b> | Ручной ввод Число больше 0 и меньше 1<br>По умолчанию – 0,1             | Данный параметр задает размер (%) валидационной выборки для ранней остановки   | Общий параметр           |
| <b>Seed</b>  | Ручной ввод числового значения<br>По умолчанию – 42                     | Начальное числовое значение для генератора случайных чисел   | Общий параметр           |
| <b>Количество параллельных потоков</b>                       | Ручной ввод<br>По умолчанию – -1  | Данный параметр задает количество параллельных потоков, используемых для запуска   | Общий параметр           |
| <b>Максимальная глубина дерева</b>                           | Ручной ввод<br>По умолчанию – -1  | Данный параметр задает максимальную глубину дерева.<br>Значение меньше или равно 0 означает отсутствие ограничений       | Построение дерева        |
| <b>Максимальное количество листов</b>                        | Ручной ввод<br>Целочисленное значение больше 0<br>По умолчанию - 31     | Данный параметр задает максимальное количество листов в дереве   | Построение дерева        |
| <b>Минимальное количество наблюдений для создания бина</b>   | Ручной ввод<br>Целочисленное значение больше 0<br>По умолчанию – 200000 | Данный параметр задает количество наблюдений для создания бина   | Построение дерева        |
| <b>Максимальное количество бинов</b>                         | Ручной ввод<br>Целочисленное значение больше 1<br>По умолчанию – 255    | Данный параметр задает максимальное количество бинов, в которые будут группироваться значения признаков                  | Построение дерева        |
| <b>Минимальное количество наблюдений в бине</b>              | По умолчанию – 3  | Данный параметр задает минимальное число наблюдений в бине   | Построение дерева        |
| <b>Использовать EFB</b>                                      | Чекбокс   | Выбор данного чекбокса указывает на необходимость использовать EFB   | Построение дерева        |
| <b>Минимальный gain для разбиения</b>                        | Ручной ввод<br>По умолчанию – 0   | Данный параметр задает минимальное снижение потерь, необходимое для создания дальнейшего раздела на листовом узле дерева | Построение дерева        |
| <b>Минимальный вес для разбиения</b>                         | По умолчанию – 0,0001   | Данный параметр задает минимальную сумму весов экземпляров   | Построение дерева        |

| Параметр   | Возможные значения и ограничения   | Описание   | Группа параметров                  |
|--|--|--|------------------------------------|
|  |  | (Гессиан), необходимую для дочернего элемента (листа)  |                                    |
| <b>Минимальное количество наблюдений в листе</b>           | По умолчанию – 20  | Данный параметр задает минимальное количество наблюдений в листе   | Построение дерева                  |
| <b>Доля выборки признаков в каждой итерации</b>            | По умолчанию – 1   | Данный параметр задает долю переменных, используемых на каждой итерации (при построении каждого дерева).   | Построение дерева                  |
| <b>Доля выборки признаков в каждом узле дерева</b>         | По умолчанию – 1   | Данный параметр задает долю подвыборки признаков, которые будут использованы для каждого разделения (узла).  | Построение дерева                  |
| <b>Использовать экстремально рандомизированные деревья</b> | Чекбокс  | <p>Выбор данного чекбокса указывает на необходимость использования экстремально рандомизированных деревьев.</p> <p>В экстремально рандомизированных деревьях используется случайное подмножество объектов-кандидатов, но вместо поиска наиболее отличительных пороговых значений пороги выбираются случайным образом для каждого объекта-кандидата, и лучший из этих случайно сгенерированных пороговых значений выбирается в качестве правила разделения. Обычно это позволяет еще немного уменьшить дисперсию модели за счет чуть большего увеличения смещения</p> | Построение дерева                  |
| <b>Доля отбрасываемых деревьев</b>                         | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию – 0,1 | Данный параметр задает долю отбрасываемых деревьев   | Параметры для DART (Бустер = dart) |

| Параметр  | Возможные значения и ограничения   | Описание   | Группа параметров                  |
|---|--|--|------------------------------------|
| <b>Максимальное количество отбрасываемых деревьев</b>       | По умолчанию – 50  | Данный параметр задает максимальное число отбрасываемых деревьев   | Параметры для DART (Бустер = dart) |
| <b>Вероятность отсутствия отбрасывания</b>                  | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию – 0,5 | Данный параметр задает вероятность пропуска процедуры отбрасывания деревьев.<br>Если процедура отбрасывания деревьев пропущена, новые деревья добавляются также, как и gbtree  | Параметры для DART (Бустер = dart) |
| <b>Использовать XGBoost DART</b>                            | Чекбокс  | Выбор данного чекбокса указывает на необходимость использования XGBoost DART   | Параметры для DART (Бустер = dart) |
| <b>Использовать uniform drop</b>                            | Чекбокс  | Выбор данного чекбокса указывает на то, что отбрасываемые деревья выбираются равномерно.   | Параметры для DART (Бустер = dart) |
| <b>Коэффициент сохранения большого градиента</b>            | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию – 0,2 | Данный параметр задает коэффициент сохранения наблюдений с большим градиентом  | Параметры для GOSS (Бустер = GOSS) |
| <b>Коэффициент сохранения маленького градиента</b>          | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию – 0,1 | Данный параметр задает коэффициент сохранения наблюдений с маленьким градиентом  | Параметры для GOSS (Бустер = GOSS) |
| <b>Соотношение случайной подвыборки в обучающей выборке</b> | Ручной ввод<br>Число больше или равно 0 и меньше или равно 1<br>По умолчанию – 1   | Данный параметр задает соотношение случайной подвыборки в обучающей выборке.<br><br>Для того, чтобы процесс бэггинга был запущен необходимо указать значение меньше 1.   | Параметры бэггинга                 |
| <b>Частота случайной подвыборки</b>                         | Ручной ввод целочисленного значения<br>По умолчанию – 0                            | Данный параметр задает выполнение бэггинга на каждой <b>k</b> (заданной значением) итерации. Каждая <b>k</b> -я итерация LightGBM будет случайным образом отбирать ' <b>Соотношение случайной подвыборки в обучающей выборке</b> '*100% данных для использования в следующих <b>k</b> итерациях.<br><br>Значение параметра | Параметры бэггинга                 |

| Параметр   | Возможные значения и ограничения                             | Описание  | Группа параметров                                |
|--|--|---|--|
|  |  | равное 0 означает отключение бэггинга   |  |
| <b>L1 регуляризатор</b>  | По умолчанию – 0   | Регуляризация добавляет ограничения алгоритма относительно аспектов модели, которые не зависят от данных для обучения. Регуляризация обычно используется, чтобы избежать переобучения.<br>Регуляризация L1 применяется для получения максимально разреженной модели | Параметры регуляризации                          |
| <b>L2 регуляризатор</b>  | По умолчанию – 1   | Регуляризация L2 ограничивает чрезмерный рост какой-либо отдельной координаты весового вектора. Регуляризация L2 полезна в том случае, если целью является создание модели, имеющей в целом малые значения веса   | Параметры регуляризации                          |
| <b>L2 регуляризатор для линейной регрессии</b>                       | По умолчанию – 0   | Данный параметр задает линейную регуляризацию дерева  | Параметры регуляризации                          |
| <b>L2 регуляризатор в категориальном разбиении</b>                   | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 10 | Данный параметр используется для задания регуляризации категориального разделения (сплита)  | Параметры регуляризации                          |
| <b>Минимальное количество наблюдений в категориальной переменной</b> | Ручной ввод<br>Число больше 0<br>По умолчанию – 100          | Данный параметр задает минимальное количество наблюдений в категориальной переменной  | Параметры для обработки категориальных признаков |
| <b>Лимит количества разбиений для категориальной переменной</b>      | По умолчанию – 32  | Данный параметр задает ограничение на количество разбиений (сплитов) для категориальной переменной  | Параметры для обработки категориальных признаков |
| <b>Псевдосчет в сглаживании Лапласа</b>                              | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 10 | Данный параметр позволяет уменьшить влияние шумов на категориальные переменные, особенно для категорий с небольшим количеством данных.  | Параметры для обработки категориальных признаков |

| Параметр  | Возможные значения и ограничения                  | Описание   | Группа параметров                                |
|---|---|--|--|
| <b>Максимальное количество категорий для One-Hot Encoding</b> | Ручной ввод<br>Число больше 0<br>По умолчанию – 4 | Когда количество категорий одной переменной меньше или равно данному значению, будет использоваться алгоритм разделения one-vs-other («один против другого») | Параметры для обработки категориальных признаков |

**Таблица 38 Параметры узла «LightGBM»**

#### Результаты выполнения узла:

Узел «LightGBM» имеет разные результаты в зависимости от решаемой задачи.

Результаты **регрессии** представлены следующими объектами:

- Тепловые карты обучающей, валидационной и тестовой выборок.
- Таблица с метриками качества модели.
- Таблица со списком переменных, сортированных по важности.

Результаты **многоклассовой классификации** представлены следующими объектами:

- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица со списком переменных, сортированных по важности.

Результаты **бинарной классификации** представлены следующими объектами:

- График ROC.
- График Lift.
- График Cumulative Lift.
- График Gain.
- График Cumulative Gain.
- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица со списком переменных, сортированных по важности.

#### 3.2.5.8.11. Узел «Градиентный бустинг (CatBoost)»

В основе **узла «CatBoost»** лежит реализация алгоритма градиентного бустинга, которая оптимизирована под работу с категориальными признаками и хорошо работает с параметрами по умолчанию.

Используется для решения задач классификации и регрессии.

## **Алгоритм работы градиентного бустинга:**

Градиентный бустинг – алгоритм машинного обучения, который строит модель предсказания в виде ансамбля слабых предсказывающих моделей (в основном Дерево решений). На каждой итерации вычисляется отклонение предсказаний уже обученного ансамбля на обучающей выборке. Следующая добавляемая в ансамбль модель будет сводить среднее отклонение предыдущей к минимуму. Новые деревья добавляются в ансамбль до тех пор, пока ошибка уменьшается, либо пока не выполняется одно из правил «ранней остановки».

### **Особенности реализации CatBoost:**

- CatBoost строит симметричные (сбалансированные) деревья, т.е. на каждом шаге листья предыдущего дерева разделяются по одному и тому же условию.
- Алгоритм сам выбирает лучший способ обработки категориальных признаков в зависимости от решаемой задачи
- Хорошо работает с гиперпараметрами по умолчанию

**Список параметров узла** представлен в таблице ниже.

| Параметр                           | Возможные значения и ограничения  | Описание  | Группа параметров |
|------------------------------------|---|---|-------------------|
| <b>Название</b>                    | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  | Общий параметр    |
| <b>Описание</b>                    | Ручной ввод<br>Ограничений на значение нет  | Описание узла   | Общий параметр    |
| <b>Количество деревьев</b>         | Ручной ввод<br>По умолчанию - 100   | Данный параметр задает максимальное количество деревьев   | Общий параметр    |
| <b>Скорость обучения</b>           | Ручной ввод<br>По умолчанию – 0,03  | Данный параметр задает скорость обучения, которая определяет насколько быстро или медленно модель будет учиться   | Общий параметр    |
| <b>Цель обучения для регрессии</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Среднеквадратическая ошибка</li> <li>• Средняя абсолютная ошибка</li> <li>• Функция потерь Хьюбера</li> <li>• Регрессия Твиди</li> <li>• Пуассоновская регрессия</li> <li>• Квантильная регрессия</li> <li>• MAPE</li> </ul> | Данный параметр задает регрессионный показатель, используемый для обучения.<br>Предусмотрены следующие показатели: <ul style="list-style-type: none"> <li>• Среднеквадратическая ошибка</li> <li>• Средняя абсолютная ошибка</li> <li>• Функция потерь Хьюбера</li> <li>• Регрессия Твиди</li> <li>• Пуассоновская регрессия</li> </ul> | Общий параметр    |

| Параметр                                   | Возможные значения и ограничения  | Описание  | Группа параметров         |
|--|---|---|---------------------------|
|  |   | <ul style="list-style-type: none"> <li>• Квантильная регрессия</li> <li>• MAPE</li> </ul>   |                           |
| <b>Цель обучения для классификации</b>     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Logloss (по умолчанию)</li> <li>• CrossEntropy</li> <li>• MultiClass</li> </ul>  | Данный параметр задает классификационный показатель, используемый для обучения.<br>Предусмотрены следующие показатели: <ul style="list-style-type: none"> <li>• Logloss</li> <li>• CrossEntropy</li> <li>• MultiClass</li> </ul>  | Общий параметр            |
| <b>Метрика для валидации регрессии</b>     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Среднеквадратическая ошибка (по умолчанию)</li> <li>• Средняя абсолютная ошибка</li> <li>• MAPE</li> <li>• SMAPE</li> <li>• R2</li> <li>• Пуассоновская регрессия</li> <li>• Функция потерь Хьюбера</li> <li>• Квантильная регрессия</li> <li>• Регрессия Твиди</li> </ul> | Данный параметр задает метрику, используемую для обнаружения переобучения.<br>Предусмотрены следующие: <ul style="list-style-type: none"> <li>• Среднеквадратическая ошибка</li> <li>• Средняя абсолютная ошибка</li> <li>• MAPE</li> <li>• SMAPE</li> <li>• R2</li> <li>• Пуассоновская регрессия</li> <li>• Функция потерь Хьюбера</li> <li>• Квантильная регрессия</li> <li>• Регрессия Твиди</li> </ul> | Общий параметр            |
| <b>Метрика для валидации классификации</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Logloss (по умолчанию)</li> <li>• CrossEntropy</li> <li>• MultiClass</li> <li>• Precision</li> <li>• F1</li> <li>• Accuracy</li> <li>• AUC</li> </ul>  | Данный параметр задает метрику, используемую для обнаружения переобучения   | Общий параметр            |
| <b>L2 регуляризатор</b>                    | Ручной ввод<br>По умолчанию – 3   | Данный параметр задает коэффициент при члене регуляризации L2 функции потерь  | Общий параметр            |
| <b>Детектор переобучения</b>               | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• IncToDec (по умолчанию)</li> <li>• Iter</li> </ul>   | Данный параметр задает тип детектора переобучения.<br>Предусмотрены: <ul style="list-style-type: none"> <li>• IncToDec – обучение</li> </ul>  | Параметр ранней остановки |

| <b>Параметр</b>                                  | <b>Возможные значения и ограничения</b>  | <b>Описание</b>   | <b>Группа параметров</b>  |
|--|--|---|---|
|  |  | <ul style="list-style-type: none"> <li>останавливается при достижении порога</li> <li>Iter – обучение останавливается после указанного количества итераций, начиная с итерации с оптимальным значением метрики</li> </ul>   |   |
| <b>Порог для IncToDec детектора переобучения</b> | Ручной ввод<br>По умолчанию - 0  | Данный параметр задает порог для IncToDec детектора переобучения.   | Параметр ранней остановки<br>Активен при заданном Детекторе переобучения = IncToDec |
| <b>Количество итераций до ранней остановки</b>   | Ручной ввод<br>По умолчанию - 0  | Данный параметр задает количество итераций до ранней остановки  | Параметр ранней остановки   |
| <b>Количество потоков для обучения</b>           | Ручной ввод<br>По умолчанию - 0  | Данный параметр задает количество потоков для обучения  | Параметр производительности   |
| <b>Коэффициент регуляризации размера модели</b>  | Ручной ввод<br>Целое число больше или равно 0<br>По умолчанию – 0,5  | Данный параметр задает коэффициент регуляризации модели. Чем больше значение, тем меньше размер модели  | Параметр производительности   |
| <b>Seed</b>                                      | Ручной ввод<br>Целочисленное значение<br>По умолчанию - 42   | Начальное числовое значение для генератора случайных чисел  | Общий параметр  |
| <b>Максимальная глубина дерева</b>               | Ручной ввод<br>Целочисленное значение<br>По умолчанию - 6  | Данный параметр задает максимальную глубину дерева  | Параметры построения дерева   |
| <b>Стратегия построения дерева</b>               | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>SymmetricTree (по умолчанию)</li> <li>Depthwise</li> <li>Lossguide</li> </ul> | Данный параметр определяет, как будет применяться жадный алгоритм поиска.<br>Предусмотрены: <ul style="list-style-type: none"> <li>SymmetricTree – дерево строится уровень за уровнем, пока не достигнет необходимой глубины. На каждом шаге листья с предыдущего дерева</li> </ul> | Параметры построения дерева   |

| <b>Параметр</b>                                  | <b>Возможные значения и ограничения</b>   | <b>Описание</b>   | <b>Группа параметров</b>    |
|--|---|---|-----------------------------|
|  |   | <ul style="list-style-type: none"> <li>разделяются с тем же условием.</li> <li>• Depthwise – дерево строится шаг за шагом, пока не достигнет необходимой глубины. Листья разделяются с использованием условия, которое приводит к лучшему уменьшению потерь.</li> <li>• Lossguide – дерево строится по листьям до тех пор, пока не будет достигнуто заданное количество листьев. На каждом шаге разделяется нетерминальный лист с лучшим уменьшением потерь.</li> </ul> |                             |
| <b>Минимальное количество наблюдений в листе</b> | Ручной ввод<br>По умолчанию – 1   | Данный параметр задает минимальное количество обучающих наблюдений в листе  | Параметры построения дерева |
| <b>Скоринговая функция</b>                       | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Cosine (по умолчанию)</li> <li>• L2</li> </ul> | Данный параметр задает тип оценки, используемый для выбора следующего разбиения при построении дерева. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• Cosine</li> <li>• L2</li> </ul>  | Параметры построения дерева |
| <b>Доля выборки признаков в каждом разбиении</b> | Ручной ввод<br>По умолчанию – 1   | Данный параметр задает долю переменных, используемых в каждом разбиении   | Параметры построения дерева |
| <b>Максимальное количество бинов</b>             | Ручной ввод<br>По умолчанию – 254   | Данный параметр задает максимальное количество бинов, в которые будут группироваться значения признаков   | Параметры построения дерева |

| Параметр   | Возможные значения и ограничения   | Описание   | Группа параметров           |
|--|--|--|-----------------------------|
| <b>Тип квантилизации интервальных переменных</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• GreedyLogSum (по умолчанию)</li> <li>• Median</li> <li>• Uniform</li> <li>• UniformAndQuantiles</li> <li>• MaxLogSum</li> <li>• MinEntropy</li> </ul> | Данный параметр задает тип квантования интервальных переменных. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• GreedyLogSum</li> <li>• Median – включение примерно одинакового количества объектов в bucket</li> <li>• Uniform – разбиение происходит разделением сегмента (минимальное значение переменной – максимальное значение переменной) на подсегменты одинаковой длины. В этом случае используются абсолютные значения признака</li> <li>• UniformAndQuantiles – комбинируются сплиты, полученные в типах Median и Uniform, предварительно уменьшив размер квантования вдвое</li> <li>• MaxLogSum</li> <li>• MinEntropy</li> </ul> | Параметры построения дерева |
| <b>Максимальное количество листов</b>            | Ручной ввод<br>По умолчанию – 31   | Данный параметр задает максимальное количество листьев в результирующем дереве. Актуален, если Стратегия построения дерева = Lossguide   | Параметры построения дерева |
| <b>Метод обработки пропусков</b>                 | Раскрывающийся список со следующими значениями:  | Данный параметр задает метод работы с пропущенными значениями.   | Параметры построения дерева |

| <b>Параметр</b>                           | <b>Возможные значения и ограничения</b>   | <b>Описание</b>   | <b>Группа параметров</b>    |
|---|---|---|-----------------------------|
|   | <ul style="list-style-type: none"> <li>• Forbidden (по умолчанию)</li> <li>• Min</li> <li>• Max</li> </ul>  | <p>Предусмотрены следующие:</p> <ul style="list-style-type: none"> <li>• Forbidden – наличие пропущенных значений вызовет ошибку</li> <li>• Min – пропущенные значения будут приняты за максимальные значения для данного признака</li> <li>• Max – пропущенные значения будут приняты как минимальные значения для данного признака.</li> </ul>  |                             |
| <b>Метод вычисления значений в листах</b> | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>• Auto (по умолчанию)</li> <li>• Newton</li> <li>• Gradient</li> <li>• Exact</li> </ul> | <p>Данный параметр задает метод вычисления значений в листах.</p> <p>Данный параметр зависит от режима и выбранной функции потерь:</p> <ul style="list-style-type: none"> <li>• Регрессия с функциями потерь Quantile или MAE — одна точная итерация (Exact).</li> <li>• Регрессия с любой функцией потерь, кроме Quantile или MAE — одна итерация градиента (Gradient).</li> <li>• Режим классификации – Десять ньютоновских итераций (Newton).</li> <li>• Режим мультиклассификации – одна итерация по Ньютону (Newton).</li> </ul> | Параметры построения дерева |

| <b>Параметр</b>   | <b>Возможные значения и ограничения</b>   | <b>Описание</b>  | <b>Группа параметров</b>    |
|---|---|--|-----------------------------|
| <b>Количество итераций в листах</b>                         | Ручной ввод<br>По умолчанию – 0   | Данный параметр регулирует количество шагов, выполняемых в каждом дереве при вычислении значений листьев   | Параметры построения дерева |
| <b>Тип отступа при вычислении значений в листах</b>         | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• AnyImprovement (по умолчанию)</li> <li>• No</li> <li>• Armijo</li> </ul>             | Данный параметр задает тип бэктрекинга, использующийся при градиентном спуске. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• AnyImprovement – уменьшает шаг спуска до того, как значение функции потерь будет меньшим, чем оно было на последней итерации.</li> <li>• No</li> <li>• Armijo – уменьшает шаг спуска до тех пор, пока не будет выполнено условие Вольфе.</li> </ul> | Параметры построения дерева |
| <b>Тип Bootstrap</b>  | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Bernoulli (по умолчанию)</li> <li>• Bayesian</li> <li>• MVS</li> <li>• No</li> </ul> | Данный параметр определяет метод семплинга весов объектов  | Бэггинг                     |
| <b>Соотношение случайной подвыборки в обучающей выборке</b> | Ручной ввод<br>По умолчанию – 0,8   | Данный параметр задает соотношение случайной подвыборки в обучающей выборке  | Бэггинг                     |
| <b>Частота сэмплирования весов при построении деревьев</b>  | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• PerTreeLevel (по умолчанию)</li> <li>• PerTree</li> </ul>                            | Данный параметр задает частоту выборки весов и объектов при построении деревьев. Предусмотрены: <ul style="list-style-type: none"> <li>• PerTreeLevel – перед построением каждого нового дерева</li> <li>• PerTree – перед выбором каждого нового</li> </ul>   | Бэггинг                     |

| Параметр   | Возможные значения и ограничения  | Описание  | Группа параметров  |
|--|---|---|--------------------|
|  |   | разделения дерева   |                    |
| <b>Распределение весов в Байесовском Bootstrap</b> | Ручной ввод<br>Целое число больше или равно 0<br>По умолчанию – 1   | Данный параметр определяет распределение из которого выбираются веса для байесовского типа Bootstrap. Веса выбираются из экспоненциального распределения, если значение этого параметра установлено на 1. Все веса равны 1, если значение этого параметра установлено на 0.         | Бэггинг            |
| <b>Вес знаменателя</b>                             | Ручной ввод<br>Число больше или равно 0<br>По умолчанию – 1   | Данный параметр влияет на вес знаменателя и может использоваться для балансировки между сэмплинга на основе важности и Бернулли сэмплинг. Значение приближенное к 0 подразумевает использование сэмплинга на основе важности.   | Бэггинг            |
| <b>Тип бустинга</b>                                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Ordered (по умолчанию)</li> <li>• Plain</li> </ul>       | Данный параметр задает схему бустинга. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• Ordered – упорядоченная схема, обеспечивает лучшее качество на небольших наборах данных</li> <li>• Plain – простая для классической схемы градиентного бустинга</li> </ul> | Параметры бустинга |
| <b>Тип сжатия модели</b>                           | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Constant (по умолчанию)</li> <li>• Decreasing</li> </ul> | Данный параметр определяет как фактический коэффициент сжатия модели рассчитывается на каждой итерации.   | Параметры бустинга |
| <b>Коэффициент сжатия модели</b>                   | Ручной ввод<br>Число в диапазоне (0, 1]<br>По умолчанию – 1   | Данный параметр задает константу, используемую для расчета коэффициента   | Параметры бустинга |

| Параметр  | Возможные значения и ограничения  | Описание   | Группа параметров                                 |
|---|---|--|---|
|   |   | умножения модели на каждой итерации.   |   |
| <b>Максимальное количество категорий для One-hot encoding</b>       | Ручной ввод<br>По умолчанию – 2   | Данный параметр задает использование one-hot encoding для всех категориальных признаков с количеством уникальных значений меньшим или равным заданному значению. | Параметры для обработки категориальных переменных |
| <b>Максимальное количество объединяемых категорий</b>               | Ручной ввод<br>Целое число больше или равно 1 и меньше или равно 16<br>По умолчанию – 4 | Данный параметр задает максимальное количество объединяемых категорий  | Параметры для обработки категориальных переменных |
| <b>Максимальное количество листов с категориальными переменными</b> | Ручной ввод<br>Целое число больше или равно 0<br>По умолчанию – 0                       | Данный параметр задает максимальное количество листьев с категориальными признаками. Если количество превышает указанное значение, часть листьев отбраковывается | Параметры для обработки категориальных переменных |

Таблица 39 Параметры узла «LightGBM»

#### Результаты выполнения узла:

Узел «CatBoost» имеет разные результаты в зависимости от решаемой задачи.

Результаты **регрессии** представлены следующими объектами:

- Тепловые карты обучающей, валидационной и тестовой выборок.
- Таблица с метриками качества модели.
- Таблица со списком переменных, сортированных по важности.

Результаты **многоклассовой классификации** представлены следующими объектами:

- Таблица с метриками качества модели.
- Таблица с метриками качества модели задачи классификации.
- Таблица со списком переменных, сортированных по важности.

Результаты **бинарной классификации** представлены следующими объектами:

- График ROC.
- График Lift.
- График Cumulative Lift.
- График Gain.

- График Cumulative Gain.
- Таблица с метриками качества модели.

### 3.2.5.8.12. Узел «GLM»

Обобщенная линейная модель (GLM), которая лежит в основе данного узла, обобщает линейную регрессию и допускает наличие у зависимой переменной распределения, отличающегося от нормального. GLM связывает зависимую переменную с факторами посредством задаваемой функции связи.

**Список параметров узла** представлен в таблице ниже.

| Параметр                                | Возможные значения и ограничения  | Описание   |
|---|---|--|
| <b>Название</b>                         | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                         | Ручной ввод<br>Ограничений на значение нет  | Описание узла  |
| <b>Тип распределения</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Распределение Пуассона</li> <li>• Гамма распределение</li> <li>• Распределение Твиди</li> </ul>  | Данный параметр задает тип распределения обобщенной линейной модели. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• Распределение Пуассона</li> <li>• Гамма распределение</li> <li>• Распределение Твиди</li> </ul>   |
| <b>Функция связи</b>                    | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• auto</li> <li>• identity – для нормального распределения</li> <li>• log – для типа распределения Пуассона и гамма распределения</li> </ul> | Данный параметр задает функцию связи. Предусмотрены следующие: <ul style="list-style-type: none"> <li>• auto</li> <li>• identity – для нормального распределения</li> <li>• log – для типа распределения Пуассона и гамма распределения</li> </ul>   |
| <b>Распределение целевой переменной</b> | Ручной ввод<br>Число больше или равное 0 и меньше или равное 3<br>По умолчанию - 0  | Данный параметр задает распределение целевой переменной, где: <ul style="list-style-type: none"> <li>• 0 задает нормальное распределение</li> <li>• 1 задает распределение Пуассона</li> <li>• (1,2) задает составное распределение Пуассона - гамма</li> <li>• 2 задает гамма распределение</li> <li>• 3 задает обратное гауссово распределение</li> </ul> Данный параметр доступен при выборе Типа распределения = Распределение Твиди |
| <b>Добавить константу в модель</b>      | Чекбокс   | Выбор данного чекбокса добавит константу в модель  |

| Параметр                      | Возможные значения и ограничения   | Описание  |
|-------------------------------|--|---|
| <b>Стандартизация</b>         | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Нет</li> <li>• Стандартное отклонение (по умолчанию)</li> <li>• Диапазон</li> </ul> | Данный параметр отвечает за выбор метода стандартизации данных. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>• Нет.</li> <li>• Стандартное отклонение – преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li> <li>• Диапазон – линейно преобразует значения переменных в диапазон [0, 1].</li> </ul> |
| <b>L2</b>                     | Ручной ввод<br>Число больше или равное 0<br>По умолчанию – 0   | Данный параметр задает константу, которая определяет силу регуляризации.<br>$L2 = 0$ эквивалентна GLM без штрафов   |
| <b>Количество итераций</b>    | Ручной ввод<br>Целое число больше или равно 1<br>По умолчанию – 100  | Данный параметр задает максимальное количество итераций после достижения которого алгоритм останавливается  |
| <b>Допустимая погрешность</b> | Ручной ввод<br>Число больше 0<br>По умолчанию – 0,0001   | Данный параметр задает допустимую погрешность, после достижения которой алгоритм останавливается  |

**Таблица 40 Параметры узла «LightGBM»**

#### **Результаты выполнения узла:**

- Тепловые карты обучающей, валидационной и тестовой выборок.
- Таблица с метриками качества модели.
- Таблица с коэффициентами переменных

#### 3.2.5.8.13. Узел «Нейронная сеть (PyTorch)»

**Узел "Нейронная сеть (PyTorch)"** позволяет строить нейросеть типа **MLP (multilayer perceptron, многослойный перцептрон)**, в которой входной сигнал преобразуется в выходной, проходя последовательно через скрытые слои.

Данный узел отличается от узла "Нейронная сеть" (подробнее [Узел «Нейронная сеть»](#)) возможностью более глубокой настройки конфигурации слоев нейросети и широким выбором оптимизаторов.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения           | Описание   | Группа параметров |
|-----------------|--|--|-------------------|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе | Общие параметры   |

| Параметр                     | Возможные значения и ограничения           | Описание  | Группа параметров |
|------------------------------|--|---|-------------------|
| Описание                     | Ручной ввод<br>Ограничений на значение нет | Описание узла   | Общие параметры   |
| Конфигурация слоев нейросети | Кнопка                                     | <p>При выборе данной кнопки откроется окно <b>Конфигурация слоев нейросети</b>, где можно задавать слои, количество нейронов в слоях и функции активации для каждого узла. Для этого необходимо выбрать кнопку <b>Добавить</b> и в появившемся списке выбрать необходимую функцию активации и настроить параметры.</p> <p>Предусмотрены:</p> <ul style="list-style-type: none"> <li>• <b>Функция активации ReLU (Rectified Linear Unit)</b> –<br/> <math>f(x) = \max(0, x)</math></li> <li>• <b>Функция активации CELU</b> –<br/> <math>f(x) = \max(0, x) + \min(0, \alpha * (\exp(x/\alpha) - 1))</math><br/>         Дополнительно необходимо задать <math>\alpha</math> <ul style="list-style-type: none"> <li>• <b>Функция активации ELU</b> –<br/> <math>f(x) = \begin{cases} \alpha * (\exp(x) - 1) &amp; \text{for } x \leq 0 \\ x &amp; \text{for } x &gt; 0 \end{cases}</math><br/>         Дополнительно необходимо задать <math>\alpha</math> <ul style="list-style-type: none"> <li>• <b>Функция активации Sigmoid</b> –<br/> <math>\sigma(x) = \frac{1}{1 + e^{-x}}</math></li> <li>• <b>Функция активации Softmax</b> –<br/> <math>\sigma(x_i) = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}} \text{ for } i = 1, 2, \dots, K</math></li> <li>• <b>Линейный слой</b> –<br/>         линейная функция, результат пропорционален переданному аргументу<br/> <math>f(x) = x</math><br/>         Дополнительно необходимо задать <b>Количество выходных переменных</b> и при необходимости выбрать чекбокс <b>Добавить константу</b> <ul style="list-style-type: none"> <li>• <b>Функция активации Logsigmoid</b> –<br/> <math>\text{LogSigmoid}(x) = \log(\frac{1}{1 + e^{-x}})</math></li> <li>• <b>Исключение (Dropout)</b> –<br/>         Во время обучения случайным образом обнуляет некоторые</li> </ul> </li> </ul> </li> </ul> </li> </ul> | Общие параметры   |

| Параметр                    | Возможные значения и ограничения  | Описание  | Группа параметров |
|-----------------------------|---|---|-------------------|
|                             |   | <p>элементы входного тензора с вероятностью <math>p</math>, используя выборки из распределения Бернулли. Дополнительно необходимо задать <b>Вероятность исключения</b></p> <ul style="list-style-type: none"> <li>• <b>Tanh</b> – функция гиперболического тангенса</li> </ul> $f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}}$ <p>После настройки необходимой конфигурации слоев нейросети выбрать кнопку <b>Сохранить</b>.</p>   |                   |
| <b>Функция потерь</b>       | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>• MAE</li> <li>• MSE</li> <li>• Poisson loss</li> <li>• Negative log likelihood</li> <li>• Cross entropy</li> </ul>   | <p>Данный параметр задает функцию потерь. Предусмотрены:</p> <ul style="list-style-type: none"> <li>• <b>MAE</b> (средняя абсолютная ошибка)</li> <li>• <b>MSE</b> (среднеквадратическая ошибка)</li> <li>• <b>Poisson loss</b></li> <li>• <b>Negative log likelihood</b> (отрицательное логарифмическое правдоподобие)</li> <li>• <b>Cross entropy</b> (перекрестная энтропия)</li> </ul>  | Общие параметры   |
| <b>Алгоритм оптимизации</b> | <p>Раскрывающийся список со следующими значениями:</p> <ul style="list-style-type: none"> <li>• SGD</li> <li>• Adam (по умолчанию)</li> <li>• Adadelta</li> <li>• Adamax</li> <li>• LBFGS</li> <li>• AdamW</li> <li>• ASGD</li> <li>• NAdam</li> <li>• RAdam</li> <li>• Rprop</li> <li>• RMSprop</li> </ul> | <p>Данный параметр задает метод оптимизации, который будет использоваться для обновления весов нейронов скрытых слоев нейронной сети. Предусмотрены следующие методы:</p> <ul style="list-style-type: none"> <li>• <b>SGD (stochastic gradient descent)</b> – стохастический градиентный спуск. Данный метод делает шаг постоянной величины в направлении, указанном градиентом в текущей точке</li> <li>• <b>Adam (adaptive moment estimation)</b> – адаптивная оценка момента. Данный метод сочетает в себе идею метода Momentum о накоплении градиента и идею методов Adadelta и RMSProp об экспоненциальном сглаживании информации о предыдущих значениях квадратов градиентов.</li> <li>• <b>Adadelta (adaptive learning rate)</b> – метод адаптивной скорости обучения</li> </ul> | Общие параметры   |

| Параметр | Возможные значения и ограничения | Описание  | Группа параметров |
|----------|----------------------------------|---|-------------------|
|          |                                  | <ul style="list-style-type: none"> <li>• <b>Adamax</b> – модификация метода Adam, основанная на бесконечной норме (max)</li> <li>• <b>LBFGS (limited-memory BFGS)</b> – BFGS с ограниченной памятью.</li> <li>• <b>AdamW</b>. Данный метод основан на адаптивной оценке моментов первого и второго порядка с добавленным методом уменьшения весов</li> <li>• <b>ASGD (average stochastic gradient descent)</b> – усредненный стохастический градиентный спуск. Данный метод усредняет веса, вычисляемые на каждой итерации.</li> <li>• <b>Nadam (Nesterov-accelerated adaptive momentum)</b>. Данный метод представляет собой модификацию оптимизатора Adam с добавлением момента Нестерова при вычислении градиентов.</li> <li>• <b>RAdam (Rectified Adam)</b>. Данный метод является модификацией Adam, более устойчивой к изменению значений скорости обучения.</li> <li>• <b>Rprop (resilient backpropagation)</b> – устойчивый алгоритм обратного распространения. Данный метод использует только знаки частных производных для подстройки весовых коэффициентов. Также Rprop поддерживает отдельные дельты для каждого веса и смещения и адаптирует эти дельты во время обучения.</li> <li>• <b>RMSPROP (root mean square propagation)</b> – среднеквадратичное распространение корня. Данный метод использует усредненный по истории квадрат градиента.</li> </ul> |                   |

| <b>Параметр</b>  | <b>Возможные значения и ограничения</b>  | <b>Описание</b>  | <b>Группа параметров</b> |
|--|--|--|--------------------------|
| <b>Скорость обучения</b>                               | Ручной ввод<br>По умолчанию - 0,001  | Данный параметр задает скорость обучения, которая управляет размером шага при обновлении весов.  | Общие параметры          |
| <b>Количество эпох</b>                                 | Ручной ввод<br>По умолчанию - 10   | Данный параметр задает сколько раз алгоритм обучения будет обрабатывать весь набор обучающих данных.   | Общие параметры          |
| <b>Размер пакета</b>                                   | Ручной ввод<br>По умолчанию - 128  | Данный параметр определяет количество выборок, которые необходимо обработать перед обновлением параметров модели.  | Общие параметры          |
| <b>Seed</b>  | Ручной ввод<br>По умолчанию - 42   | Начальное числовое значение для генератора случайных чисел. Используется для воспроизведения результатов при повторном запуске узла.   | Общие параметры          |
| <b>L2 регуляризация</b>                                | Ручной ввод<br>По умолчанию - 0,1  | Данный параметр задает значение L2-регуляризации   | Общие параметры          |
| <b>Доля валидационной выборки</b>                      | Ручной ввод<br>По умолчанию - 0,1  | Данный параметр задает долю валидационной выборки, которая будет отобрана из исходной тестовой   | Общие параметры          |
| <b>Количество итераций без существенного улучшения</b> | Ручной ввод<br>По умолчанию - 5  | Данный параметр задает количество эпох без улучшения, после которых скорость обучения будет снижена.<br>Пример: если значение параметра = 2, то первые 2 эпохи без улучшений loss будут проигнорированы, и только после 3-й эпохи скорость обучения уменьшится.    | Общие параметры          |
| <b>Режим динамического расчета порога</b>              | Раскрывающийся список со следующими значениями:<br><ul style="list-style-type: none"> <li>• rel</li> <li>• abs</li> </ul>                                | Данный параметр задает режим расчета порога<br>Предусмотрены: <ul style="list-style-type: none"> <li>• <b>rel</b> – <code>dynamic_threshold = best * ( 1 - threshold )</code></li> <li>• <b>abs</b> – <code>dynamic_threshold = best - threshold</code></li> </ul> | Общие параметры          |
| <b>Порог</b>   | Ручной ввод<br>По умолчанию - 0,0001   | Данный параметр задает порог расчета нового оптимума, чтобы сосредоточиться только на значительных изменениях  | Общие параметры          |
| <b>Стандартизация</b>                                  | Раскрывающийся список со следующими значениями:<br><ul style="list-style-type: none"> <li>• no</li> <li>• std (по умолчанию)</li> <li>• range</li> </ul> | Данный параметр отвечает за выбор метода стандартизации числовых переменных.<br><b>Стандартизация</b> – преобразование числовых наблюдений с целью приведения их к некоторой общей шкале.<br>Необходимость стандартизации вызвана тем, что разные признаки         | Общие параметры          |

| Параметр                             | Возможные значения и ограничения                             | Описание   | Группа параметров  |
|--------------------------------------|--|--|--|
|                                      |  | <p>из обучающего набора могут быть представлены в разных масштабах и изменяться в разных диапазонах, что влияет на выявление некорректных зависимостей моделью.</p> <p>Предусмотрены следующие методы:</p> <ul style="list-style-type: none"> <li>• <b>no</b> — стандартизация не нужна</li> <li>• <b>std</b> — стандартное отклонение - преобразует наблюдения таким образом, чтобы их среднее значение равнялось нулю, а стандартное отклонение равнялось 1.</li> <li>• <b>range</b> — диапазон - линейно преобразует значения переменных в диапазон <math>[0, 1]</math>.</li> </ul> |  |
| <b>Beta1</b>                         | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,9   | Данный параметр задает коэффициент, используемый для управления скоростью затухания скользящих средних значений градиента (первого момента)  | Параметры алгоритма оптимизации <b>Adam</b> , <b>Adamax</b> , <b>AdamW</b> , <b>RAdam</b> , <b>Nadam</b> |
| <b>Beta2</b>                         | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,999 | Данный параметр задает коэффициент, используемый для управления скоростью затухания средних значений вторых моментов градиентов (некентрированной дисперсии)   | Параметры алгоритма оптимизации <b>Adam</b> , <b>Adamax</b> , <b>AdamW</b> , <b>RAdam</b> , <b>Nadam</b> |
| <b>Epsilon</b>                       | Ручной ввод числа с плавающей точкой<br>По умолчанию - 1e-8  | Данный параметр задает значение, добавляемое к знаменателю для улучшения числовой стабильности Minimal decay applied to lr. If the difference between new and old lr is smaller than eps, the update is ignored. Default: 1e-8.  | Параметры алгоритма оптимизации <b>Adam</b> , <b>Adamax</b> , <b>AdamW</b> , <b>RAdam</b> , <b>Nadam</b> |
| <b>Использовать алгоритм AMSGrad</b> | Чекбокс  | Выбор данного чекбокса указывает, что необходимо использовать вариант AMSGrad этого алгоритма.<br>Разница между AMSgrad и Adam заключается в рассчитанном векторе второго момента, который используется для обновления параметров.   | Параметры алгоритма оптимизации <b>Adam</b> , <b>AdamW</b>   |

| <b>Параметр</b>                               | <b>Возможные значения и ограничения</b>                        | <b>Описание</b>   | <b>Группа параметров</b>                                 |
|---|--|---|--|
| <b>Импульс (Momentum)</b>                     | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0       | Данный параметр задает коэффициент импульса (запоминает скорость на предыдущем шаге и добавляет в указанное число раз меньшую величину на следующем шаге)   | Параметры алгоритма оптимизации <b>SGD, RMSProp</b>      |
| <b>Dampening</b>                              | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0       | Данный параметр задает демпфирование импульса. Демпфирование гарантирует, что оптимизатор не сделает слишком больших шагов, что может произойти, если использовать только импульс. Чем выше градиент, тем больше демпфирование уменьшает размер шага. | Параметры алгоритма оптимизации <b>SGD</b>               |
| <b>Момент Нестерова</b>                       | Чекбокс  | Выбор данного чекбокса включает импульс Нестерова (использует производную не в текущей точке, а в следующей, если бы мы продолжали двигаться в этом же направлении без изменений)   | Параметры алгоритма оптимизации <b>SGD</b>               |
| <b>Rho</b>                                    | Ручной ввод<br>По умолчанию - 0,9                              | Данный параметр задает коэффициент, используемый для вычисления скользящего среднего квадратов градиентов   | Параметры алгоритма оптимизации <b>Adadelta</b>          |
| <b>Epsilon</b>                                | По умолчанию - 1e-8  | Данный сглаживающий параметр задает значение, предотвращающее деление на 0.   | Параметры алгоритма оптимизации <b>Adadelta, RMSProp</b> |
| <b>Максимум итераций за шаг оптимизации</b>   | Ручной ввод целочисленного значения<br>По умолчанию - 20       | Данный параметр задает максимальное число итераций за шаг оптимизации   | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Максимум вычислений за шаг оптимизации</b> | Ручной ввод целочисленного значения<br>По умолчанию - 1        | Данный параметр задает максимальное число вычислений функции за шаг оптимизации   | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Tolerance grad</b>                         | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,00001 | Данный параметр задает допуск завершения при оптимальности первого порядка  | Параметры алгоритма оптимизации <b>LBFGS</b>             |
| <b>Tolerance change</b>                       | Ручной ввод числа с плавающей точкой<br>По умолчанию - 1e-9    | Данный параметр задает допуск завершения при изменении значения/параметра функции   | Параметры алгоритма оптимизации <b>LBFGS</b>             |

| <b>Параметр</b>                                  | <b>Возможные значения и ограничения</b>  | <b>Описание</b>   | <b>Группа параметров</b>                     |
|--|--|---|--|
| <b>Количество запоминаемых шагов оптимизации</b> | Ручной ввод целочисленного значения<br>По умолчанию - 100  | Данный параметр задает количество запоминаемых шагов оптимизации  | Параметры алгоритма оптимизации <b>LBFGS</b> |
| <b>Line search</b>                               | Список: <ul style="list-style-type: none"><li>• no (по умолчанию)</li><li>• strong Wolfe</li></ul> | Данный параметр задает метод линейного поиска   | Параметры алгоритма оптимизации <b>LBFGS</b> |
| <b>Lambda</b>                                    | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,0001                                      | Данный параметр задает затухание  | Параметры алгоритма оптимизации <b>ASGD</b>  |
| <b>Alpha</b>                                     | Ручной ввод числа с плавающей точкой<br>По умолчанию - 0,75  | Данный параметр задает мощность для обновления скорости обучения  | Параметры алгоритма оптимизации <b>ASGD</b>  |
| <b>t0</b>  | Ручной ввод<br>По умолчанию - 100000   | Данный параметр задает точку, с которой начинается усреднение. Если требуемое количество итераций меньше данного значения, то усреднение не произойдет.   | Параметры алгоритма оптимизации <b>ASGD</b>  |
| <b>Сокращение импульса</b>                       | Ручной ввод<br>По умолчанию - 0,004  | Данный параметр задает значение сокращения импульса   | Параметры алгоритма оптимизации <b>Nadam</b> |
| <b>Коэффициент уменьшения (eta minus)</b>        | Ручной ввод<br>Число больше 0 и меньше 1<br>По умолчанию - 0,5                                     | Данный параметр задает мультипликативный коэффициент уменьшения.<br>Если на текущем шаге частная производная по соответствующему весу поменяла свой знак, значит последнее изменение было большим, и алгоритм проскочил локальный минимум.<br>Следовательно, величину коррекции необходимо уменьшить на значение данного параметра и вернуть предыдущее значение весового коэффициента. | Параметры алгоритма оптимизации <b>Rprop</b> |
| <b>Коэффициент увеличения (eta plus)</b>         | Ручной ввод<br>Число больше 1<br>По умолчанию - 1,2  | Данный параметр задает мультипликативный коэффициент увеличения.<br>Если на текущем шаге частная производная по соответствующему весу не поменяла свой знак, значит нужно увеличить величину коррекции на значение данного  | Параметры алгоритма оптимизации <b>Rprop</b> |

| Параметр                        | Возможные значения и ограничения       | Описание  | Группа параметров                              |
|---------------------------------|--|---|--|
|                                 |  | параметра для достижения более быстрой сходимости.  |  |
| <b>Минимальный размер шага</b>  | Ручной ввод<br>По умолчанию - 0,000001 | Данный параметр задает минимальный размер шага. Он необходим, чтобы не допустить слишком маленьких значений весов, ограничивает величину коррекции снизу. | Параметры алгоритма оптимизации <b>Rprop</b>   |
| <b>Максимальный размер шага</b> | Ручной ввод<br>По умолчанию - 50       | Данный параметр задает максимальный размер шага. Он необходим, чтобы не допустить слишком больших значений весов, ограничивает величину коррекции сверху. | Параметры алгоритма оптимизации <b>Rprop</b>   |
| <b>Alpha</b>                    | Ручной ввод<br>По умолчанию - 0,99     | Данный параметр задает константу сглаживания  | Параметры алгоритма оптимизации <b>RMSProp</b> |
| <b>Центрировать</b>             | Чекбокс                                | Выбор данного чекбокса указывает, что необходимо вычислить центрированный RMSProp, градиент которого нормализуется по оценке его дисперсии                | Параметры алгоритма оптимизации <b>RMSProp</b> |

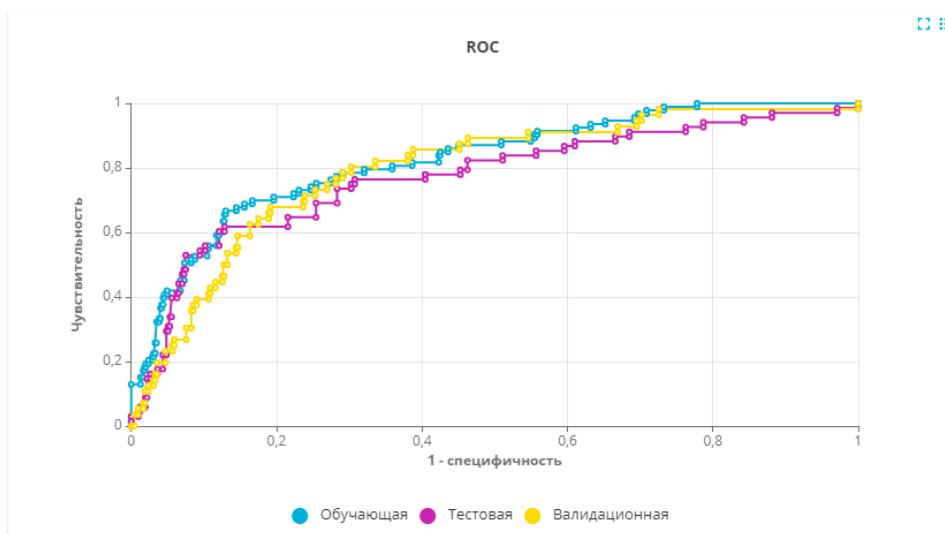
Таблица 41 Параметры узла «Нейронная сеть (PyTorch)»

### 3.2.5.8.13.1. Результаты выполнения узла

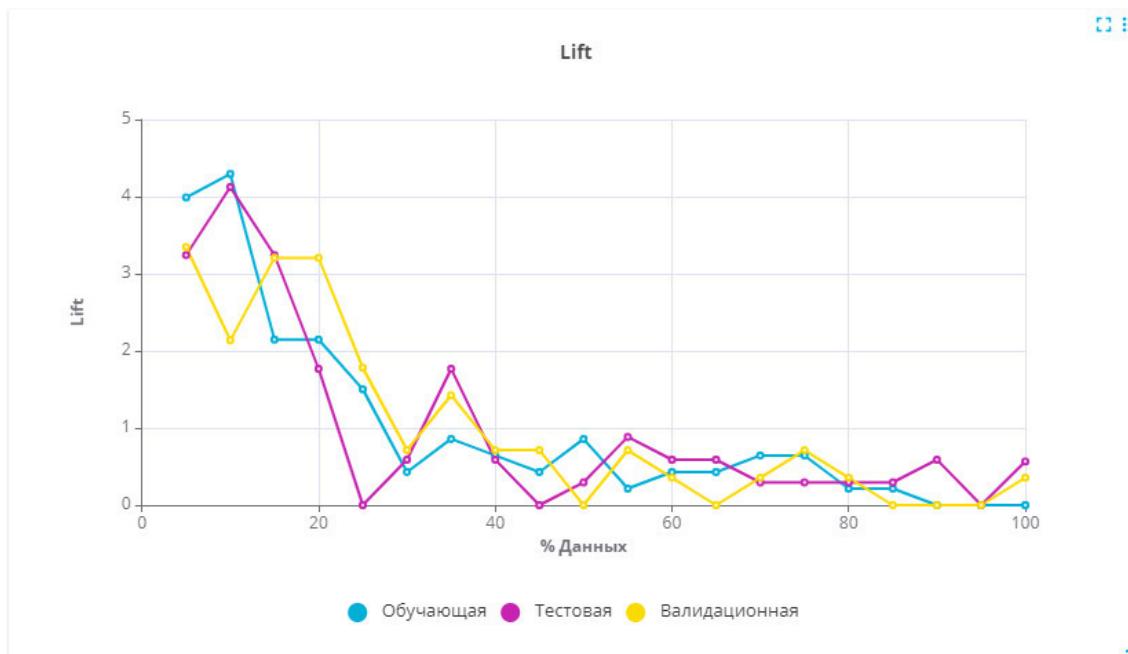
Узел «Нейронная сеть» имеет разные результаты в зависимости от решаемой задачи.

**Результаты бинарной классификации представлены следующими объектами:**

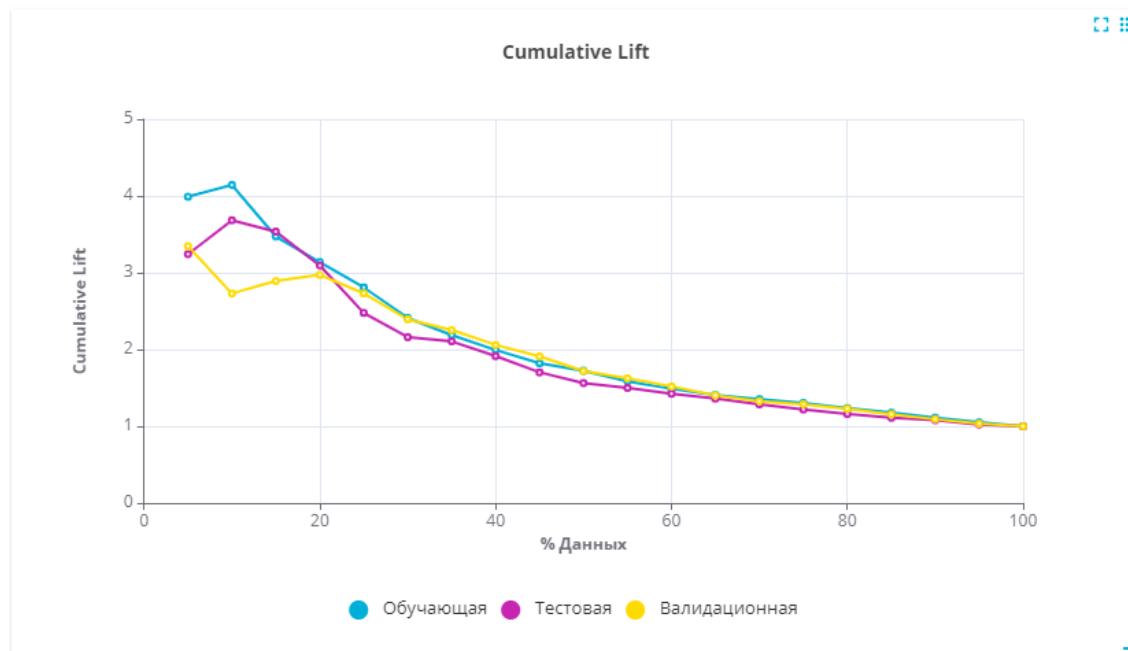
- График ROC.



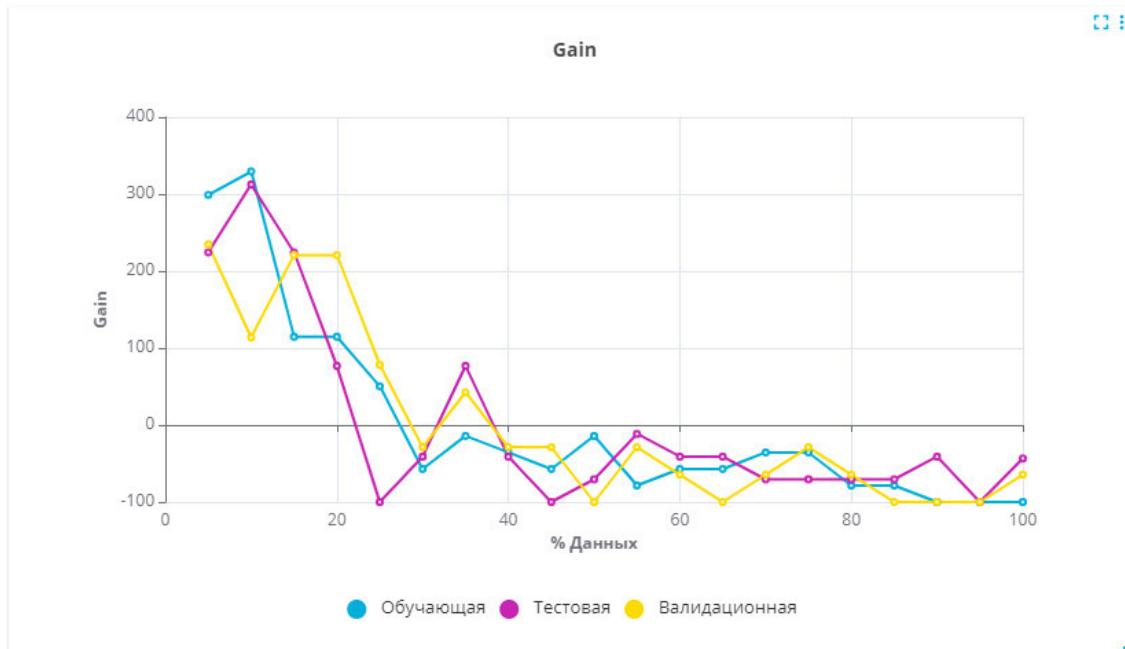
- График Lift.



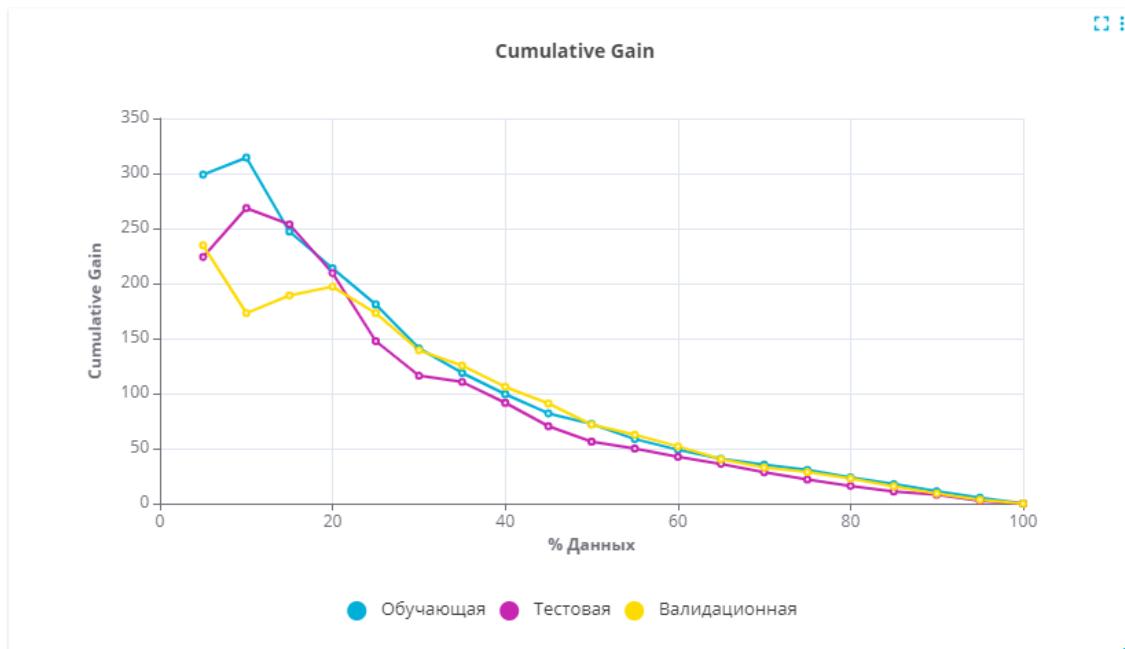
- График Cumulative Lift.



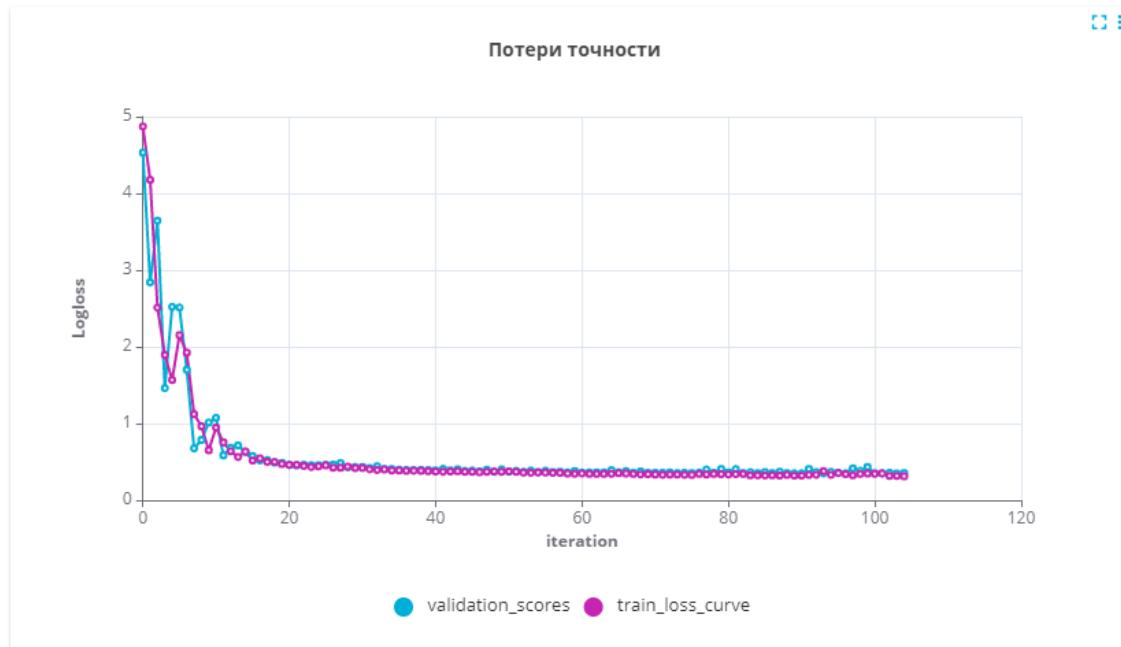
- График Gain.



- График Cumulative Gain.



- График потери точности.



- Таблица с метриками качества модели.

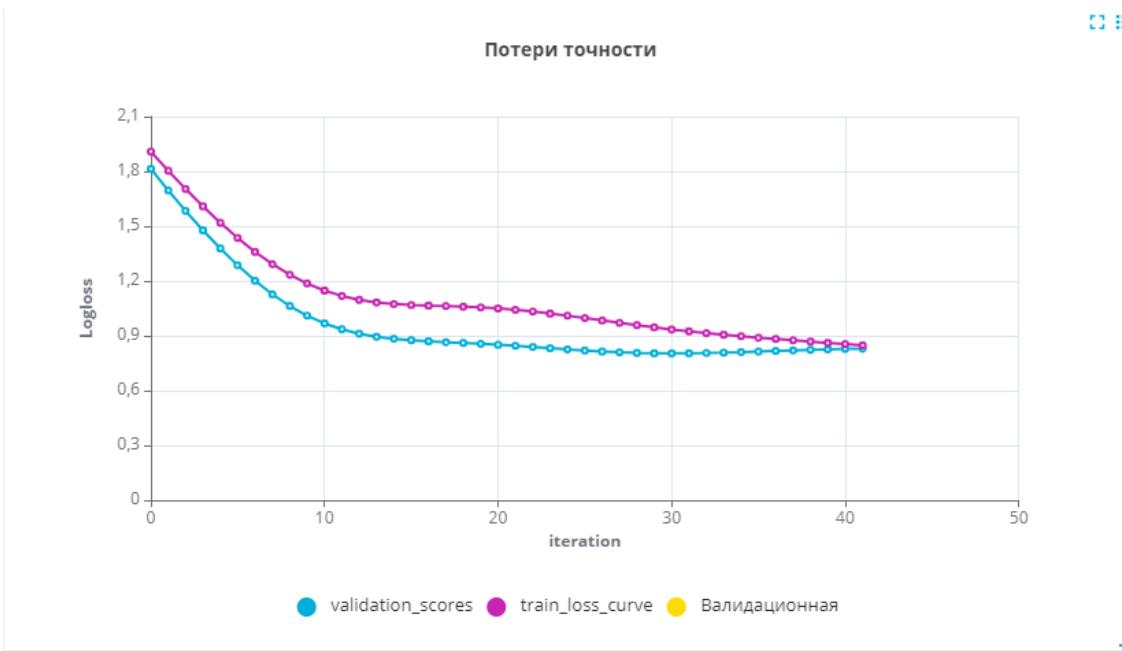
| ↑             | AUC ROC   | gini      | log loss ↑ | nobs ↑    |
|---------------|-----------|-----------|------------|-----------|
| Filter...     | Filter... | Filter... | Filter...  | Filter... |
| Валидационная | 0.793     | 0.587     | 0.336      | 479       |
| Обучающая     | 0.822     | 0.644     | 0.318      | 639       |
| Тестовая      | 0.77      | 0.541     | 0.362      | 481       |

- Таблица с метриками качества модели задачи классификации.

| <b>↑</b>      | <b>misclassification ↑</b> | <b>mcc ↑</b> | <b>nobs ↑</b> |
|---------------|----------------------------|--------------|---------------|
| Filter...     | Filter...                  | Filter...    | Filter...     |
| Валидационная | 0.121                      | 0.168        | 479           |
| Обучающая     | 0.134                      | 0.272        | 639           |
| Тестовая      | 0.149                      | 0.088        | 481           |

**Результаты многоклассовой классификации представлены следующими объектами:**

- График потери точности.



- Таблица с метриками качества модели.

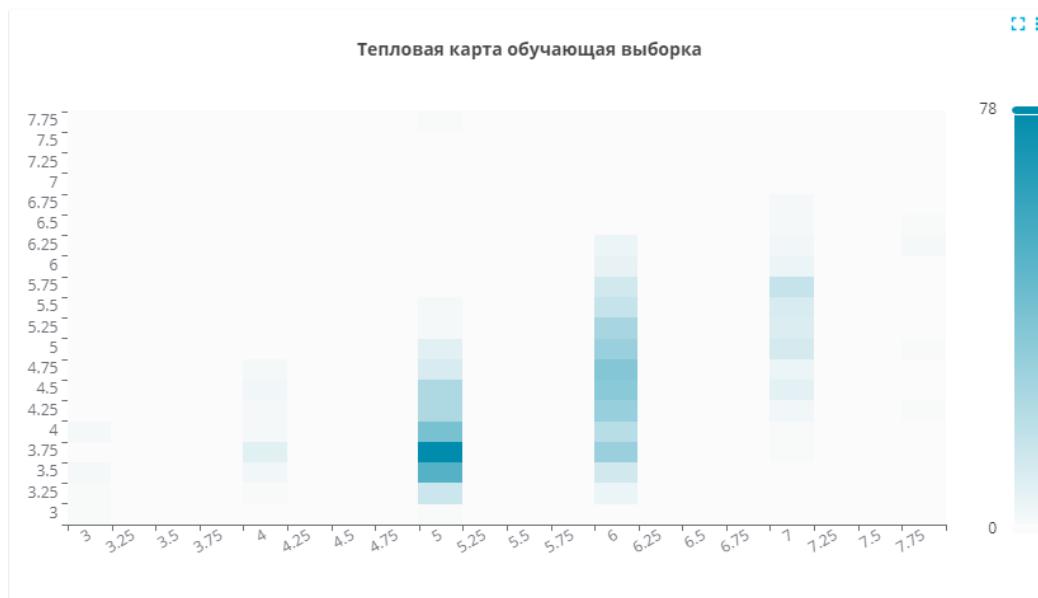
|               | log loss  | nobs      |
|---------------|-----------|-----------|
| Filter...     | Filter... | Filter... |
| Валидационная | 1.052     | 35        |
| Обучающая     | 0.911     | 47        |
| Тестовая      | 1.003     | 37        |

- Таблица с метриками качества модели задачи классификации.

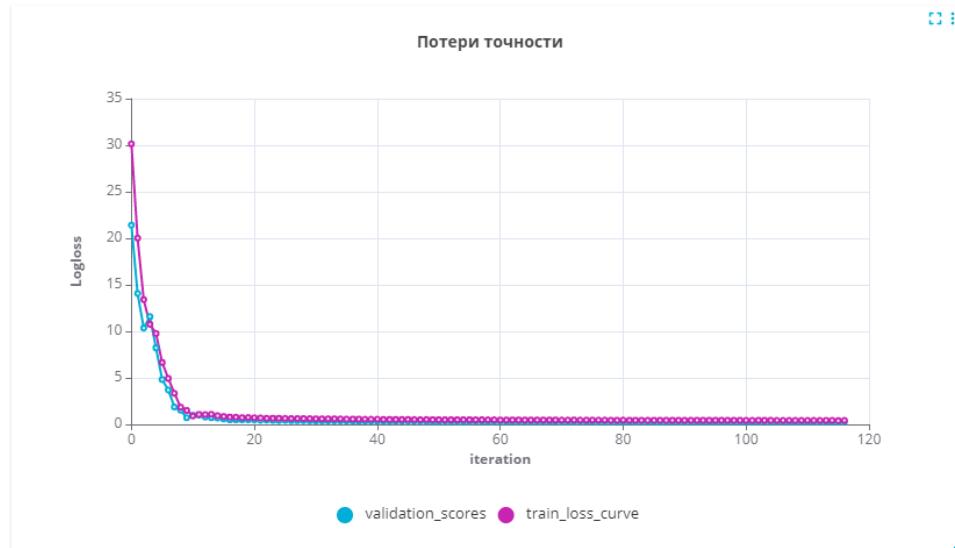
|               | misclassification | mcc       | nobs      |
|---------------|-------------------|-----------|-----------|
| Filter...     | Filter...         | Filter... | Filter... |
| Валидационная | 0.657             | 0         | 35        |
| Обучающая     | 0.489             | 0         | 47        |
| Тестовая      | 0.621             | 0         | 37        |

### Результаты регрессии представлены следующими объектами:

- Тепловые карты для обучающей, валидационной и тестовой выборок.



- График потери точности.



- Таблица с метриками качества модели.

| Метрики модели     |       |       |       |       |       |       |
|--------------------|-------|-------|-------|-------|-------|-------|
|                    | mse   | rmse  | mae   | mape  | r2    | nobs  |
| CV (валидационная) | 0.562 | 0.745 | 0.566 | 0.104 | 0.084 | 127.8 |
| CV (обучающая)     | 0.509 | 0.705 | 0.533 | 0.098 | 0.218 | 511.2 |
| Валидационная      | 0.527 | 0.726 | 0.567 | 0.103 | 0.178 | 479   |
| Обучающая          | 0.412 | 0.642 | 0.487 | 0.09  | 0.361 | 639   |
| Тестовая           | 0.471 | 0.686 | 0.529 | 0.096 | 0.293 | 481   |

- Таблицы с метриками качества в разбиениях обучающей и валидационной выборках.

| Метрики в разбиениях обучающей выборки |          |       |       |       |       |        |      |
|--|----------|-------|-------|-------|-------|--------|------|
| Номер разбие...                        | fit_time | mse   | rmse  | mae   | mape  | r2     | nobs |
| 0                                      | 0.381    | 0.465 | 0.682 | 0.511 | 0.096 | 0.315  | 511  |
| 1                                      | 0.397    | 0.409 | 0.64  | 0.489 | 0.092 | 0.34   | 511  |
| 2                                      | 0.348    | 0.413 | 0.642 | 0.481 | 0.087 | 0.332  | 511  |
| 3                                      | 0.096    | 0.837 | 0.915 | 0.71  | 0.129 | -0.224 | 511  |
| 4                                      | 0.342    | 0.42  | 0.648 | 0.474 | 0.086 | 0.328  | 512  |

### 3.2.5.8.14. Узел «AutoML»

**Узел "AutoML"** позволяет в автоматизированном режиме построить алгоритмы машинного обучения с различными параметрами и выбрать из них лучшую, исходя из заданной метрики качества.

Автоматическое машинное обучение (AutoML) позволяет автоматизировать процесс дизайна ML-пайплайнов. Это позволяет сократить требуемые ресурсы на процесс отбора оптимальных гиперпараметров модели и на тестирование.

AutoML состоит из трех шагов:

- предобработка данных;
- предобработка признаков;
- обучение и тестирование модели



**Рисунок 157** Принцип работы узла "AutoML"

В результате запуска узла AutoML Система автоматически перебирает различные методы балансировки, предобработки и моделирования и объединяет модели в ансамбли.

При необходимости пользователь может посмотреть, какие модели выбрала система и какие модели тестировались.

На рисунке ниже представлен пример результата AutoML.

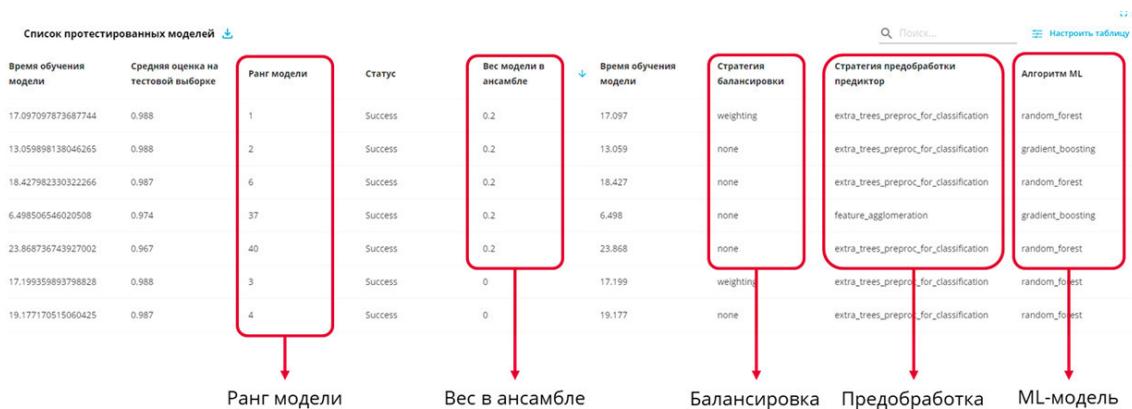


Рисунок 158 Пример результата узла AutoML

В представленном примере итоговая модель состоит из 5 моделей с различными предобработками и балансировками. Результаты каждой модели имеют вес 0,2.

|         | Балансировка | Предобработка           | ML-модель           |
|---------|--------------|-------------------------|---------------------|
| $f_1 =$ | Веса классов | Деревья решений         | Случайный лес       |
| $f_2 =$ | —            | Деревья решений         | Градиентный бустинг |
| $f_3 =$ | —            | Деревья решений         | Случайный лес       |
| $f_4 =$ | —            | Кластеризация признаков | Градиентный бустинг |
| $f_5 =$ | —            | Деревья решений         | Случайный лес       |

$$Model = 0,2 \cdot f_1 + 0,2 \cdot f_2 + 0,2 \cdot f_3 + 0,2 \cdot f_4 + 0,2 \cdot f_5$$

Рисунок 159 Пример итоговой модели

Список параметров узла представлен в таблице ниже.

| Параметр   | Возможные значения и ограничения           | Описание   |
|--|--|--|
| <b>Название</b>  | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>  | Ручной ввод<br>Ограничений на значение нет | Описание узла  |
| <b>Время на построение всех моделей, с</b>             | По умолчанию - 300                         | Количество времени в секундах, которым ограничено выполнение данного узла                      |
| <b>Время на построение одной модели, с</b>             | По умолчанию - 30                          | Количество времени в секундах, которым ограничено время построения одной модели                |
| <b>Количество конфигураций подбора гиперпараметров</b> | По умолчанию - 25                          | Количество конфигураций подбора гиперпараметров для ускорения подбора с помощью алгоритма SMAC |

| <b>Параметр</b>                                 | <b>Возможные значения и ограничения</b>  | <b>Описание</b>  |
|---|--|--|
| <b>Размер ансамблевых моделей</b>               | По умолчанию - 5   | Максимальное количество моделей в ансамбле   |
| <b>Количество лучших моделей в ансамбле</b>     | По умолчанию - 50  | Количество лучших моделей, участвующих в отборе в ансамбль   |
| <b>Максимальное количество моделей на диске</b> | По умолчанию - 50  | Максимальное количество моделей, которое будет построено в рамках временных ограничений  |
| <b>Seed</b>                                     | По умолчанию - 42  | Начальное числовое значение для генератора случайных чисел   |
| <b>Лимит памяти модели</b>                      | По умолчанию - 3072  | Лимит памяти для обучения модели   |
| <b>Типы предобработки переменных</b>            | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• densifier</li> <li>• extra_trees_preproc</li> <li>• fast_ica</li> <li>• feature_agglomeration</li> <li>• kernel_pca</li> <li>• kitchen_sinks</li> <li>• no.preprocessing</li> <li>• nystroem_sampler</li> <li>• pca</li> <li>• polynomial</li> <li>• random_trees_embedding</li> <li>• select_percentile</li> <li>• select_rates</li> <li>• truncatedSVD</li> </ul> | Набор методов предобработки переменных, который будет использован при переборе вариантов моделей.<br>Предусмотрены: <ul style="list-style-type: none"> <li>• densifier (Перевод данных из разреженного представления в плотное)</li> <li>• extra_trees_preproc (Отбор признаков с помощью ансамбля сильно рандомизированных деревьев)</li> <li>• fast_ica (Быстрый метод независимых компонент)</li> <li>• feature_agglomeration (Кластеризация признаков)</li> <li>• kernel_pca (Ядерный метод главных компонент)</li> <li>• kitchen_sinks (Аппроксимация ядерной функции методом случайных признаков Фурье (метод кухонных раковин))</li> <li>• no.preprocessing (Без предобработки)</li> <li>• nystroem_sampler (Аппроксимация ядерной функции методом Nyström'a)</li> <li>• pca (Метод главных компонент)</li> <li>• polynomial (Полиноминальные предикторы)</li> <li>• random_trees_embedding (Бинарное кодирование сильно рандомизированными деревьями)</li> </ul> |

| Параметр                         | Возможные значения и ограничения   | Описание  |
|----------------------------------|--|---|
|                                  |  | <ul style="list-style-type: none"> <li>select_percentile (Отбор признаков по перцентилю)</li> <li>select_rates (Отбор признаков по частотам ошибок)</li> <li>truncatedSVD (Усечённое сингулярное разложение)</li> </ul>   |
| <b>Типы используемых моделей</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>adaboost</li> <li>ard_regression</li> <li>bernoulli_nb</li> <li>decision_tree</li> <li>extra_trees</li> <li>gaussian_nb</li> <li>gaussian_process</li> <li>gradient_boosting</li> <li>k_nearest_neigbors</li> <li>Ida</li> <li>liblinear_svm</li> <li>libsvm_svc</li> <li>mlp</li> <li>multinomial_nb</li> <li>passive_aggressive</li> <li>qda</li> <li>random_forest</li> <li>sgd</li> </ul> | Набор алгоритмов, который будет использован при построении вариантов моделей.<br>Предусмотрены: <ul style="list-style-type: none"> <li>adaboost (Адаптивный бустинг (AdaBoost))</li> <li>ard_regression (Байесовская линейная регрессия с автоматическим определением актуальности (ARD Regression))</li> <li>bernoulli_nb (Наивный байесовский классификатор с распределением Бернулли)</li> <li>decision_tree (Дерево решений)</li> <li>extra_trees (Ансамбль сильно рандомизированных деревьев решений)</li> <li>gaussian_nb (Наивный байесовский классификатор с нормальным распределением)</li> <li>gaussian_process (Гауссовский процесс)</li> <li>gradient_boosting (Градиентный бустинг)</li> <li>k_nearest_neigbors (Метод К-ближайших соседей)</li> <li>Ida (Линейный дискриминантный анализ)</li> <li>liblinear_svm (Линейная машина опорных векторов)</li> <li>libsvm_svc (Машина опорных векторов)</li> <li>mlp (Многослойный перцептрон)</li> <li>multinomial_nb (Наивный байесовский классификатор с полиномиальным распределением)</li> </ul> |

| <b>Параметр</b>                        | <b>Возможные значения и ограничения</b>  | <b>Описание</b>  |
|--|--|--|
|  |  | <ul style="list-style-type: none"> <li>• <code>passive_aggressive</code> (Пассивно-агрессивный алгоритм)</li> <li>• <code>qda</code> (Квадратичный дискриминантный анализ)</li> <li>• <code>random_forest</code> (Случайный лес)</li> <li>• <code>sgd</code> (Линейная модель, обучаемая методом стохастического градиентного спуска)</li> </ul>   |
| <b>Тип ресемплинга</b>                 | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• <code>holdout</code></li> <li>• <code>cv</code></li> </ul>  | Тип ресемплинга выборки. Предусмотрены: <ul style="list-style-type: none"> <li>• <code>holdout</code> (Отложенная выборка)</li> <li>• <code>cv</code> (Кросс-валидация)</li> </ul>   |
| <b>% обучающей выборки для AutoML</b>  | По умолчанию - 0,67  | % Обучающей выборки, который будет использован для AutoML  |
| <b>Перемешать наблюдения</b>           | Чекбокс  | При выборе этой опции наблюдения будут перемешаны заново   |
| <b>Количество параллельных потоков</b> | По умолчанию - 0   | Параллельное исполнение  |
| <b>Метрика</b>                         | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• AUC ROC</li> <li>• Logloss</li> <li>• Accuracy</li> <li>• Balanced accuracy</li> <li>• F1</li> <li>• Logloss</li> <li>• MSE</li> <li>• MAE</li> <li>• R2</li> </ul> | Метрика, которую оптимизирует AutoML. Предусмотрены:<br>Для задачи бинарной классификации: <ul style="list-style-type: none"> <li>• AUC ROC</li> <li>• Logloss</li> <li>• Accuracy</li> <li>• Balanced accuracy</li> <li>• F1</li> </ul> Для задачи многоклассовой классификации: <ul style="list-style-type: none"> <li>• Logloss</li> </ul> Для задачи регрессии: <ul style="list-style-type: none"> <li>• MSE</li> <li>• MAE</li> <li>• R2</li> </ul> |
| <b>Коэффициент сжатия датасета</b>     | По умолчанию - 1   | Сжатие датасета для оптимизации памяти. Сжатие выполняется за счёт сокращения точности чисел с плавающей точкой.   |

**Таблица 42 Параметры узла «AutoML»**

### 3.2.5.9. Группа узлов «Работа с моделями»

#### 3.2.5.9.1. Узел «Сравнение моделей»

**Узел «Сравнение моделей»** оценивает построенные модели по метрикам качества и выбирает лучшую.

**Список параметров узла** представлен в таблице ниже.

| Параметр                                    | Возможные значения и ограничения   | Описание   |
|---|--|--|
| <b>Название</b>                             | Ручной ввод<br>Ограничений на значение нет   | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>                             | Ручной ввод<br>Ограничений на значение нет   | Описание узла  |
| <b>Метрика для сравнения. Регрессия</b>     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• MSE</li> <li>• RMSE</li> <li>• MAE</li> <li>• MAPE</li> <li>• R2</li> </ul> | Данный параметр задает метрику, по которой будет производиться сравнение регрессионных моделей |
| <b>Метрика для сравнения. Классификация</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• log loss</li> <li>• Gini</li> <li>• AUC ROC</li> </ul>                      | Данный параметр задает метрику, по которой будет производиться сравнение моделей классификации |
| <b>Выборка для сравнения</b>                | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Обучающая</li> <li>• Валидационная</li> <li>• Тестовая</li> </ul>           | Данный параметр задает выборку для сравнения   |

**Таблица 43 Параметры узла «Сравнение моделей»**

#### Результаты выполнения узла:

- Таблица отранжированных по выбранной метрике моделей с указанием модели-победителя.

| Результаты сравнения |                               |
|----------------------|-------------------------------|
| Узел                 | Результат (MISCLASSIFICATION) |
| Filter...            | Filter...                     |
| BASE_RF              | 0.137                         |
| BaseLDA              | 0.149                         |
| DT-crossval5         | 0.229                         |
| RF_CUTOFF_06         | 0.149                         |
| Победитель: BASE_RF  | 0.137                         |

**Рисунок 160 Пример таблицы с результатами сравнения моделей**

### 3.2.5.9.2. Узел «Регистрация модели»

**Узел «Регистрация модели»** позволяет сохранить построенную модель в выбранном проекте репозитория Model Manager.

**ВАЖНО!** Перед тем как зарегистрировать модель из MD в Репозиторий (при помощи узла «Регистрация модели») необходимо создать **Проект ММ**.

Является конечным узлом сценария моделирования.

**Список параметров узла** представлен в таблице ниже.

| Параметр          | Возможные значения и ограничения             | Описание  |
|-------------------|--|---|
| <b>Название</b>   | Ручной ввод<br>Ограничений на значение нет   | Название узла, которое будет отображаться в интерфейсе        |
| <b>Описание</b>   | Ручной ввод<br>Ограничений на значение нет   | Описание узла   |
| <b>Имя модели</b> | Ручной ввод<br>Ограничений на значение нет   | Наименование модели, которое будет отображаться в репозитории |
| <b>Проект ММ</b>  | Раскрывающийся список с доступными проектами | Выбор проекта ММ, в котором будет сохранена модель            |

Таблица 44 Параметры узла «Регистрация модели»

#### Результаты выполнения узла:

- Графические и табличные представления не предусмотрены.

### 3.2.5.9.3. Узел «Интерпретация»

Методы **узла «Интерпретация»** позволяют объяснить, как отдельные признаки и элементы модели влияют на целевую переменную. Рекомендуется устанавливать узел **«Фильтр»** перед узлом **«Интерпретация»** для отбора наблюдений, которые необходимо интерпретировать. Большое количество наблюдений значительно увеличивает время расчета.

**Список параметров узла** представлен в таблице ниже.

| Параметр                 | Возможные значения и ограничения  | Описание  |
|--------------------------|---|---|
| <b>Название</b>          | Ручной ввод<br>Ограничений на значение нет  | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>          | Ручной ввод<br>Ограничений на значение нет  | Описание узла   |
| <b>Тип интерпретации</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>LIME</li> <li>PD</li> <li>SHAP</li> <li>ICE</li> </ul> | Данный параметр задает тип интерпретации. Предусмотрены следующие методы: <ul style="list-style-type: none"> <li>LIME – Локальные интерпретируемые объяснения модели</li> <li>PD – Частичная зависимость</li> <li>SHAP – Аддитивные объяснения Шэпли</li> <li>ICE – Индивидуальные условные ожидания</li> </ul> |

Таблица 45 Параметры узла «Интерпретация»

Каждый из указанных методов интерпретации имеет свои параметры и результаты.

- **Метод LIME** (локально интерпретируемое объяснение модели).

Данный метод строит модель линейной регрессии, чтобы аппроксимировать предсказания исходной неинтерпретируемой модели локально, а не глобально.

Для этого создается новый набор данных из наблюдений, которые находятся вокруг выбранного для интерпретации наблюдения. Затем этот новый набор данных используется для обучения интерпретируемой линейной модели.

Коэффициенты этой линейной регрессии позволяют оценить важность и направление влияния каждого из предикторов при построении прогноза для выбранного наблюдения.

**Список параметров** метода LIME представлен в таблице ниже.

| Параметр                              | Возможные значения и ограничения                           | Описание   |
|---------------------------------------|--|--|
| <b>Количество признаков</b>           | Ручной ввод целочисленного значения<br>По умолчанию — 5    | Данный параметр задает количество признаков                |
| <b>Размер выборки</b>                 | Ручной ввод целочисленного значения<br>По умолчанию — 5000 | Данный параметр задает размер выборки                      |
| <b>Количество объясняемых классов</b> | Ручной ввод целочисленного значения<br>По умолчанию — 0    | Данный параметр задает количество объясняемых классов      |
| <b>Seed</b>                           | Ручной ввод числового значения<br>По умолчанию — 42        | Начальное числовое значение для генератора случайных чисел |

Таблица 46 Параметры метода LIME

#### Результаты выполнения метода:

В окне с результатами во вкладке **Выбор наблюдений** необходимо указать наблюдение для интерпретации. Во вкладках **Графические результаты** и **Табличные результаты** будут отображены список объяснений, отражающих вклад каждой функции в прогноз выборки данных в виде графиков и таблицы соответственно (рисунок ниже). Это обеспечивает локальную интерпретируемость, а также позволяет определить, какие изменения характеристик окажут наибольшее влияние на прогноз.

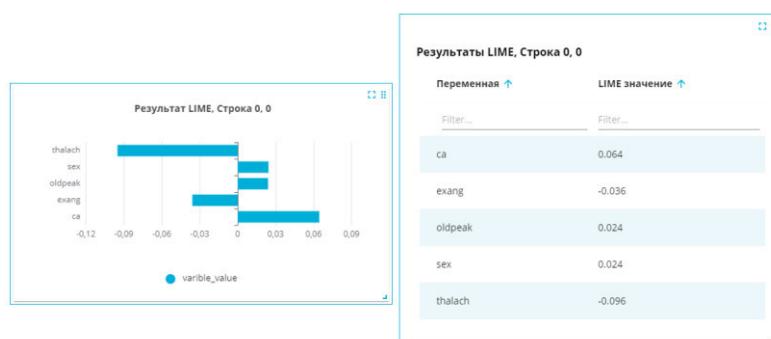


Рисунок 161 Пример графических и табличных результатов метода LIME

- **Метод PD.**

График частичной зависимости (Partial Dependence Plots – PDP) показывает, как меняется средний прогноз при изменении одного из предикторов.

**Алгоритм работы метода PD:** выбирается переменная и непрерывно изменяется ее значение. На график наносятся эти значения переменной и соответствующее среднее значение прогноза по выборке. Таким образом получается график зависимости прогнозируемых результатов от значений переменной.

Данный график частичной зависимости может показать, является ли отношение между целью и признаком линейным, монотонным или более сложным. Например, при применении к модели линейной регрессии графики частичной зависимости всегда показывают линейную зависимость.

**Список параметров** метода PD представлен в таблице ниже.

| Параметр                | Возможные значения и ограничения                           | Описание   |
|-------------------------|--|--|
| <b>Размер выборки</b>   | Ручной ввод целочисленного значения<br>По умолчанию — 5000 | Данный параметр задает размер выборки  |
| <b>Количество бинов</b> | Ручной ввод целочисленного значения<br>По умолчанию — 5    | Данный параметр задает количество бинов  |
| <b>Квантили</b>         | Таблица со значениями левого и правого квантилей           | Для редактирования значений квантилей необходимо выбрать ссылку <b>Редактировать</b> . В открывшемся окне задать новые значения квантилей. |
| <b>Seed</b>             | Ручной ввод числового значения<br>По умолчанию — 42        | Начальное числовое значение для генератора случайных чисел   |
| <b>Полный биннинг</b>   | Чекбокс  | Выбор данного чекбокса указывает на необходимость полного биннинга   |

Таблица 47 Параметры метода PD

**Результаты выполнения метода:**

В окне с результатами во вкладке **Графические результаты** выбрать предиктор и класс, чтобы график отразил связь между предиктором и прогнозом.

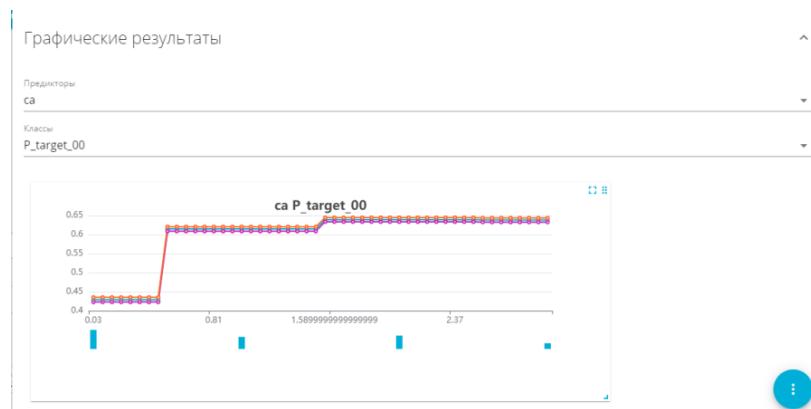


Рисунок 162 Пример графических результатов метода PD

- **Метод SHAP (Shapley Additive Explanation).**

**SHAP value** показывает средний вклад каждого предиктора в формирование прогноза для наблюдения, выбранного для интерпретации. Усреднение делается по всевозможным комбинациям всех остальных предикторов.

**Алгоритм расчета SHAP values:** Для каждого возможного упорядочивания признаков берутся все признаки, стоящие перед i-м признаком, и считается величина (прирост эффективности от добавления признака i в комбинацию признаков), равная разности между фактическим прогнозом и прогнозом с учетом текущего набора значений признака. После чего полученные значения усредняются по всем упорядочиваниям. Это означает, что SHAP values описывают ожидаемый прирост выходного значения модели при добавлении i-го признака в текущем примере.

**Список параметров** метода SHAP представлен в таблице ниже.

| Параметр                | Возможные значения и ограничения                          | Описание                                       |
|-------------------------|---|--|
| Размер выборки          | Ручной ввод целочисленного значения<br>По умолчанию — 100 | Данный параметр задает размер выборки          |
| Количество перестановок | Ручной ввод целочисленного значения<br>По умолчанию — 10  | Данный параметр задает количество перестановок |

**Таблица 48 Параметры метода SHAP**

#### Результаты выполнения метода:

В окне с результатами во вкладке **Выбор наблюдений** необходимо указать наблюдение для интерпретации. Во вкладках **Графические результаты** и **Табличные результаты** будут отображены список объяснений, отражающих вклад каждой функции в прогноз выборки данных в виде графиков и таблицы соответственно.



**Рисунок 163 Пример графических и табличных результатов метода SHAP**

- Метод ICE — Индивидуальные условные ожидания.**

**График индивидуального условного ожидания (ICE)** показывает, как меняется прогноз модели на одном или нескольких наблюдениях при изменении одного из предикторов. Однако, в отличие от PDP, который показывает средний эффект входной характеристики, график ICE визуализирует зависимость прогноза от характеристики для каждой выборки отдельно с одной строкой на выборку.

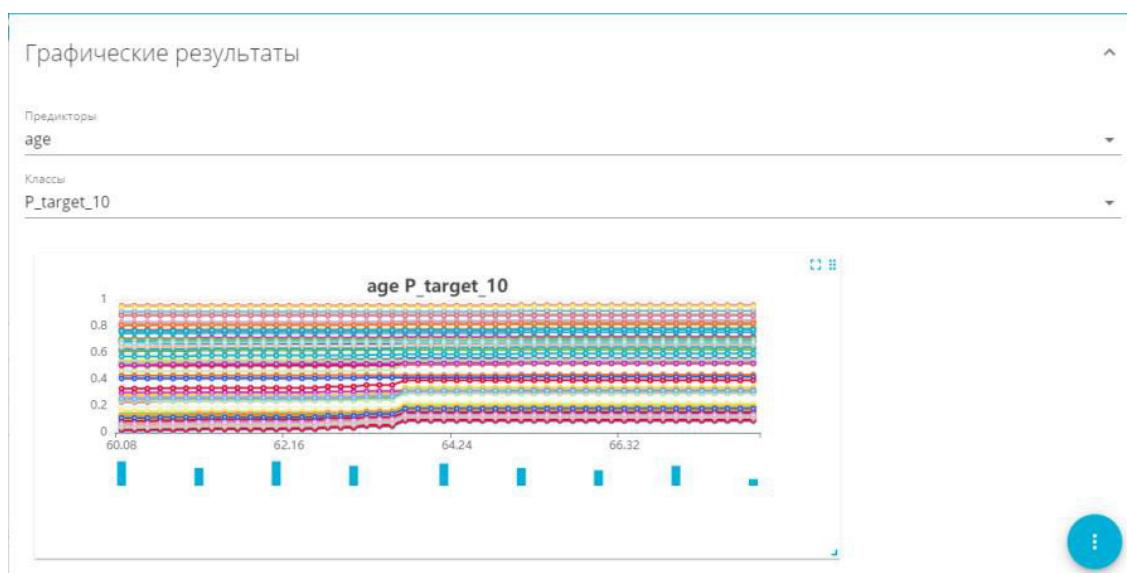
**Список параметров** метода ICE представлен в таблице ниже.

| Параметр         | Возможные значения и ограничения                           | Описание   |
|------------------|--|--|
| Размер выборки   | Ручной ввод целочисленного значения<br>По умолчанию — 5000 | Данный параметр задает размер выборки  |
| Количество бинов | Ручной ввод целочисленного значения<br>По умолчанию — 5    | Данный параметр задает количество бинов  |
| Квантили         | Таблица со значениями левого и правого квантилей           | Для редактирования значений квантилей необходимо выбрать ссылку <b>Редактировать</b> . В открывшемся окне задать новые значения квантилей. |
| Seed             | Ручной ввод числового значения<br>По умолчанию — 42        | Начальное числовое значение для генератора случайных чисел   |
| Полный биннинг   | Чекбокс  | Выбор данного чекбокса указывает на необходимость полного биннинга   |

**Таблица 49 Параметры метода ICE**

**Результаты выполнения метода:**

В окне с результатами во вкладке **Графические результаты** выбрать предиктор и класс, чтобы график отразил связь между предиктором и значением.



**Рисунок 164 Пример графических результатов метода ICE**

### 3.2.5.9.4. Узел «Подбор отсечки (Cut off)»

**Узел «Подбор отсечки»** позволяет определить оптимальный порог отсечения. Порог отсечения нужен для того, чтобы относить новые примеры к одному из двух классов (задача бинарной классификации).

**Список параметров узла** представлен в таблице ниже.

| Параметр                    | Возможные значения и ограничения   | Описание  |
|-----------------------------|--|---|
| <b>Название</b>             | Ручной ввод<br>Ограничений на значение нет   | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>             | Ручной ввод<br>Ограничений на значение нет   | Описание узла   |
| <b>Количество разбиений</b> | Ручной ввод целочисленного значения<br>По умолчанию — 100  | Данный параметр задает количество разбиений             |
| <b>Размер бина для Lift</b> | Ручной ввод целочисленного значения<br>По умолчанию — 20   | Данный параметр задает размер бина для расчета Lift     |
| <b>Критерий отсечки</b>     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Число (по умолчанию)</li> <li>• Колмогорова-Смирнова</li> </ul> | Данный параметр задает критерий выбора порога отсечения |
| <b>Отсечка</b>              | Ручной ввод<br>Число больше 0 меньше 1<br>По умолчанию — 0,5   | Данный параметр задает порог отсечки                    |

Таблица 50 Параметры узла «Подбор отсечки (Cut off)»

#### Результаты выполнения узла:

- Таблица со статистикой разбиения, в которой отражены следующие показатели:
  - Cutoff value – значение порога отсечения (считается в соответствии с заданным параметром Количество разбиений).
  - True Positives – количество верно классифицированных положительных примеров.
  - False Positives – количество неверно классифицированных положительных примеров.
  - True Negatives – количество верно классифицированных отрицательных примеров.
  - False Negatives – количество неверно классифицированных отрицательных примеров.
  - Predicted Positive – сумма True Positives и False Positives.
  - Predicted Negative – сумма True Negatives и False Negatives.
  - False Positives and Negatives – сумма False Positives и False Negatives
  - True Positives and Negatives – сумма True Positives и True Negatives

- Accuracy – результат деления суммы True Positive и True Negative на сумму всех значений (доля правильных ответов).
- True Positive Rate (Recall, полнота) – доля найденных объектов класса к общему числу объектов класса. Показывает, насколько хорошо классификатор находит объекты из класса.
- True Negative Rate – доля верно классифицированных отрицательных примеров от общего количества отрицательных примеров.
- False Positive Rate – доля неверных срабатываний классификатора к общему числу объектов за пределами класса. Показывает, насколько часто классификатор ошибается при отнесении того или иного объекта к классу.
- Precision (точность) – показывает долю объектов класса среди объектов, выделенных классификатором.
- F-score – гармоническое среднее между точностью и полнотой.
- KS (Kolmogorov-Smirnov statistic).
- partition\_id (если был в выборке).

| Статистика разбиений |                |                 |                |                 |                    |                    |                               |                              |            |                    |              |               |
|----------------------|----------------|-----------------|----------------|-----------------|--------------------|--------------------|-------------------------------|------------------------------|------------|--------------------|--------------|---------------|
| Cutoff ↑             | True Positives | False Positives | True Negatives | False Negatives | Predicted Positive | Predicted Negative | False Positives and Negatives | True Positives and Negatives | Accuracy ↑ | True Positive Rate | True Negativ | True Neat Rat |
| Filter... 0.01       | 47             | 40              | 0              | 0               | 87                 | 0                  | 40                            | 47                           | 0.54       | 1                  | 0            | 0             |
| Filter... 0.02       | 47             | 40              | 0              | 0               | 87                 | 0                  | 40                            | 47                           | 0.54       | 1                  | 0            | 0             |
| Filter... 0.03       | 47             | 38              | 2              | 0               | 85                 | 2                  | 38                            | 49                           | 0.563      | 1                  | 0            | 0.01          |
| Filter... 0.04       | 47             | 38              | 2              | 0               | 85                 | 2                  | 38                            | 49                           | 0.563      | 1                  | 0            | 0.01          |

**Рисунок 165 Пример таблицы со статистикой разбиения**

- Таблица с метриками.

| Метрики классификации |                     |       |        |
|-----------------------|---------------------|-------|--------|
| ↑                     | misclassification ↑ | mcc ↑ | nobs ↑ |
| Filter... 0           | 0.12                | 0.757 | 116    |
| Filter... 1           | 0.16                | 0.707 | 87     |
| Filter... 2           | 0.149               | 0.704 | 87     |

**Рисунок 166 Пример таблицы с метриками**

### 3.2.6. Кросс-валидация

Для повышения надежности моделей предусмотрена кросс-валидация.

**Кросс-валидация (перекрестная проверка)** — метод обучения и оценки модели, при котором исходное множество данных разделяется на несколько блоков и модель обучается на этих блоках.

Для использования кросс-валидации необходимо найти в боковой панели раздел Кросс-валидация и выбрать чекбокс **Использовать кросс-валидацию**. В результате должны появиться следующие параметры в соответствии с таблицей ниже.

| Параметр                            | Возможные значения и ограничения           | Описание   |
|-------------------------------------|--|--|
| <b>Количество разбиений</b>         | Ручной ввод<br>Неотрицательное целое число | Задает количество блоков, на которые будет разделено исходное множество данных |
| <b>Использовать стратификацию</b>   | Чекбокс                                    | Выбор данного чекбокса указывает на необходимость использовать стратификацию   |
| <b>Переменные для стратификации</b> | Список переменных набора данных            | Данный параметр позволяет выбрать переменные для стратификации.                |

Таблица 51 Параметры кросс-валидации

### 3.2.7. Автоподбор гиперпараметров

**Подбор параметров** — одна из важных задач для построения модели машинного обучения. Изменение параметров модели может принципиально повлиять на ее качество. Перебор этих параметров вручную может занять колоссальное количество времени. Для этого предусмотрен **автоподбор параметров** модели.

Для использования автоподбора параметров необходимо найти в боковой панели раздел **Автоподбор гиперпараметров (общее)** и выбрать чекбокс **Автоподбор параметров**. В результате должны появиться следующие параметры в соответствии с таблицей ниже.

| Параметр                                    | Возможные значения и ограничения  | Описание   |
|---|---|--|
| <b>Количество разбиений кросс-валидации</b> | Ручной ввод целочисленного значения<br>По умолчанию — 5   | Данный параметр задает количество разбиений кросс-валидации                  |
| <b>Использовать стратификацию</b>           | Чекбокс   | Выбор данного чекбокса указывает на необходимость использовать стратификацию |
| <b>Переменные для стратификации</b>         | Открывается в случае выбора чекбокса <b>Использовать стратификацию</b><br>Список переменных набора данных | Данный параметр позволяет выбрать переменные для стратификации               |
| <b>Количество итераций подбора</b>          | Ручной ввод целочисленного значения<br>По умолчанию — 10  | Данный параметр задает количество итераций подбора                           |

| Параметр                                     | Возможные значения и ограничения  | Описание   |
|--|---|--|
| <b>Метрика для оптимизации классификации</b> | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• LoglossNominal</li> <li>• LoglossBinary</li> </ul>   | Данный параметр задает метрику для оптимизации классификации |
| <b>Метрика для оптимизации регрессии</b>     | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• MSE</li> <li>• MAPE</li> </ul>   | Данный параметр задает метрику для оптимизации регрессии     |
| <b>Оценка точности автоподбора</b>           | Раскрывающийся список со следующими значениями: <ul style="list-style-type: none"> <li>• Заданная валидационная выборка</li> <li>• Новая валдационная выборка</li> <li>• Кросс-валидация</li> </ul> | Данный параметр задает оценку точности автоподбора           |
| <b>Seed</b>                                  | Ручной ввод числового значения<br>По умолчанию — 42   | Начальное числовое значение для генератора случайных чисел   |

**Таблица 52 Параметры автоподбора параметров модели**

Параметры последующей секции **Автоподбор гиперпараметров** зависят от типа модели. Для них задается диапазон параметров, в котором и производиться автоподбор.

С подробным описанием параметров моделей можно ознакомиться в описании конкретного узла.

### 3.2.8. Результаты моделирования

Результаты предусмотрены для следующих узлов:

- Разделение выборки.
- Sample.
- Фильтр.
- One-hot encoding.
- Заполнение пропусков.
- Трансформация.
- Биннинг/Энкодинг.
- Дисперсионный анализ.
- Стандартизация.
- Веса классов.
- Кластерный анализ (k-means).
- Иерархическая кластеризация.
- Дерево решений.
- Случайный лес.

- Байесовская регрессия.
- Линейная регрессия.
- Логистическая регрессия.
- Линейные модели.
- Нейронная сеть.
- LDA.
- Градиентный бустинг (XGBoost).
- Градиентный бустинг (LightGBM).
- Градиентный бустинг (CatBoost).
- GLM
- Сравнение моделей.
- Интерпретация.
- Подбор отсечки (Cut off).

Для просмотра результатов выполнения узла необходимо:

- Правой кнопкой мыши нажать по интересующему узлу.
- В открывшемся списке выбрать «**Посмотреть результат**».
- Откроется окно **Результаты выполнения узла** с табличным и/или графическим представлением результатов (наполнение окна зависит от узла, подробнее о результатах можно узнать в разделе с описанием каждого узла).

### 3.3. Пример базового сценария

В таблице ниже представлен пример базового сценария с разбивкой по шагам.

| № | Шаг                              | Узлы   |
|---|----------------------------------|--|
| 1 | <b>Подключение набора данных</b> | В любом из сценариев первоначальным узлом всегда должен быть « <b>Набор данных</b> ». Он задает набор данных, с которым далее будет производиться моделирование.   |
| 2 | <b>Задание метаданных</b>        | Узел « <b>Метаданные</b> » позволяет изменить метаданные (например, задать предикторы и целевую переменную). Целевая переменная определяет решаемую задачу: <ul style="list-style-type: none"> <li>• Классификация (Binary – бинарная классификация, Nominal – многоклассовая классификация)</li> <li>• Регрессия (Interval)</li> <li>• Кластеризация - нет необходимости указывать целевую переменную</li> <li>• Ассоциативные правила (необходимо задать переменные с ролями <b>Идентификатор</b> и <b>Предмет</b>)</li> </ul> Также узел « <b>Метаданные</b> » может использоваться на последующих этапах моделирования (например, при подготовке данных), если необходимо править метаданные (например, убрать атрибут из сценария или изменить тип атрибута). |

| <b>№</b> | <b>Шаг</b>                   | <b>Узлы</b>  |
|----------|------------------------------|--|
| 3        | <b>Подготовка данных</b>     | <p>Данный этап включает в себя процессы подготовки данных:</p> <ul style="list-style-type: none"> <li>• Заполнение пропущенных значений (узел «<b>Заполнение пропусков</b>»)</li> <li>• Фильтрация данных (узел «<b>Фильтр</b>»)</li> <li>• Создание новых расчетных атрибутов (узел «<b>Трансформация</b>»)</li> <li>• Преобразование категориальных переменных в числовые (узел «<b>One-hot encoding</b>»)</li> <li>• Бинаризация интервальных переменных и кодирование категориальных (узел «<b>Биннинг/энкодинг</b>»)</li> <li>• Дисперсионный анализ</li> <li>• Корректировка неравномерного распределения классов в исходном наборе данных (узел «<b>Sample</b>»)</li> <li>• Разделение набора данных на обучающую, валидационную и тестовую выборки (узел «<b>Разделение выборки</b>»)</li> <li>• Исследование данных для выяснения статистических характеристик переменных (узел «<b>Профилирование</b>»)</li> <li>• Преобразование числовых наблюдений с целью приведения их к общей шкале (узел «<b>Стандартизация</b>»)</li> <li>• Корректировка дисбаланса классов при помощи задания весов (узел «<b>Веса классов</b>»)</li> <li>• Уменьшение размерности - преобразование большого набора переменных в меньший (узлы «<b>PCA</b>» и «<b>Автоэнкодер (PyTorch)</b>»)</li> </ul>   |
| 4        | <b>Построение ML моделей</b> | <p>Узлы моделирования подразделяются в зависимости от решаемой задачи:</p> <ul style="list-style-type: none"> <li>• Классификация <ul style="list-style-type: none"> <li>Узел «<b>Дерево решений</b>»</li> <li>Узел «<b>Случайный лес</b>»</li> <li>Узел «<b>Логистическая регрессия</b>»</li> <li>Узел «<b>Линейные модели</b>»</li> <li>Узел «<b>Нейронная сеть</b>»</li> <li>Узел «<b>LDA</b>»</li> <li>Узел «<b>Градиентный бустинг (XGBoost)</b>»</li> <li>Узел «<b>Градиентный бустинг (LightGBM)</b>»</li> <li>Узел «<b>Градиентный бустинг (CatBoost)</b>»</li> <li>Узел «<b>Нейронная сеть (PyTorch)</b>»</li> <li>Узел «<b>AutoML</b>»</li> </ul> </li> <p>В результате выполнения данных узлов рассчитываются новые переменные, одна из которых – класс наблюдения, остальные – вероятность принадлежности к одному из классов (для каждого класса рассчитывается своя вероятность).</p> <ul style="list-style-type: none"> <li>• Регрессия <ul style="list-style-type: none"> <li>Узел «<b>Дерево решений</b>»</li> <li>Узел «<b>Случайный лес</b>»</li> <li>Узел «<b>Байесовская регрессия</b>»</li> <li>Узел «<b>Линейная регрессия</b>»</li> <li>Узел «<b>Линейные модели</b>»</li> <li>Узел «<b>Нейронная сеть</b>»</li> <li>Узел «<b>Градиентный бустинг (XGBoost)</b>»</li> <li>Узел «<b>Градиентный бустинг (LightGBM)</b>»</li> <li>Узел «<b>Градиентный бустинг (CatBoost)</b>»</li> <li>Узел «<b>GLM</b>»</li> <li>Узел «<b>Нейронная сеть (PyTorch)</b>»</li> <li>Узел «<b>AutoML</b>»</li> </ul> </li> </ul> </ul> |

| <b>№</b> | <b>Шаг</b>                | <b>Узлы</b>  |
|----------|---------------------------|--|
|          |                           | <p>В результате решения задачи регрессии в данных узлах рассчитывается переменная с результирующим значением.</p> <ul style="list-style-type: none"> <li>• Кластеризация           <ul style="list-style-type: none"> <li>Узел «<b>Кластерный анализ (k-means)</b>»</li> <li>Узел «<b>Иерархическая кластеризация</b>»</li> </ul> </li> </ul> <p>В результате решения задачи кластеризации рассчитывается переменная, в которой указывается кластер, к которому относится данное наблюдение.</p> <ul style="list-style-type: none"> <li>• Ассоциативные правила           <ul style="list-style-type: none"> <li>Узел «<b>Ассоциативные правила</b>»</li> </ul> </li> <li>• Обнаружение аномалий           <ul style="list-style-type: none"> <li>Узел «<b>Детекция аномалий</b>»</li> </ul> </li> </ul> <p>Каждый узел из группы «Обучение с учителем» имеет также параметры для кросс-валидации (метод оценки модели) и автоподбору гиперпараметров.</p> |
| 5        | <b>Работа с моделями</b>  | <p>После процесса построения моделей и перебора гиперпараметров, идет этап интерпретации и сравнения построенных моделей. Для этого предусмотрены следующие узлы:</p> <ul style="list-style-type: none"> <li>• Узел «<b>Сравнение моделей</b>» – сравнение полученных моделей и выбор лучшей</li> <li>• Узел «<b>Интерпретация</b>» включает в себя методы для интерпретации предсказаний модели – PD, LIME, ICE, SHAP</li> <li>• Узел «<b>Подбор отсечки (Cut off)</b>» позволяет подобрать отсечку для разделения на классы при бинарной классификации</li> </ul>  |
| 6        | <b>Регистрация модели</b> | <p>На конечном этапе модель-победитель можно зарегистрировать в Репозитории ММ (узел «<b>Регистрация модели</b>»).</p>   |

**Таблица 53 Пример базового сценария**

## 4. Компонент Разработка решений (Decision Manager, DM)

**Разработка решений** включает в себя обширный пласт работ. Первоначально необходимо создать переменные решения, указать доступные для цепочки Подключения, а затем настроить последовательность шагов решения, используя различные типы узлов. Перед регистрацией решения необходимо выполнить тестирования и убедиться, что результат соответствует ожиданиям.

Компонент включает в себя:

- Главный экран со списком доступных проектов решений (проектов DM).
- Конструктор цепочек решений.
- Вспомогательные окна настройки вида, создания проекта и т.д.

### 4.1. Главный экран DM

#### 4.1.1. Интерфейс главного экрана DM

Главный экран компонента **Разработка решений** открывается при выборе иконки в левом верхнем меню и представляет собой список доступных пользователю **Проектов** в табличном виде.

| Проекты (25) <a href="#">Добавить</a> |  |  |                   | <a href="#">Поиск...</a> | <a href="#">Настроить таблицу</a> |
|---------------------------------------|--|--|-------------------|--------------------------|-----------------------------------|
| <input type="checkbox"/>              | Название   | Описание   | Дата создания     | Дата обновления          |                                   |
| <input type="checkbox"/>              | тест   |  | 07.09.2022, 16:33 | 07.09.2022, 16:33        |                                   |
| <input type="checkbox"/>              | Цепочка решений для мониторинга характеристик модели | Цепочка решений для мониторинга характеристик модели                                   | 09.08.2022, 12:08 | 09.08.2022, 12:08        |                                   |
| <input type="checkbox"/>              | Цепочка решений для мониторинга SQL и Python         | Цепочка решений для мониторинга SQL и Python   | 17.08.2022, 14:07 | 17.08.2022, 14:07        |                                   |
| <input type="checkbox"/>              | Тестовый проект                                      | Тестовый проект  | 02.08.2022, 14:33 | 02.08.2022, 14:33        |                                   |
| <input type="checkbox"/>              | Тестирование формул                                  | Тестирование формул  | 21.09.2022, 13:20 | 21.09.2022, 13:20        |                                   |
| <input type="checkbox"/>              | Стратегия принятия решений по кредитному скринингу   | Принятие решений по заявкам на кредит физических лиц – клиентов финансовой организации | 26.08.2022, 14:54 | 26.08.2022, 14:54        |                                   |
| <input type="checkbox"/>              | Проект 2508  | Проект 2508  | 25.08.2022, 16:00 | 25.08.2022, 16:00        |                                   |
| <input type="checkbox"/>              | Пример ветвление 3                                   | Пример ветвление 3   | 29.08.2022, 11:15 | 29.08.2022, 11:15        |                                   |
| <input type="checkbox"/>              | Сумма качества данных: для пропущенных значений      | Регулярная проверка доли пропущенных значений в таблицах из всех каналов               | 14.09.2022, 11:33 | 14.09.2022, 11:33        |                                   |

Рисунок 167 Главный экран компонента Разработка решений (Decision Manager)

**Проект DM** – это набор цепочек решений, настраиваемых пользователем для решения конкретной задачи.

Таблица с доступными проектами имеет гибкие настройки отображения. Так, пользователь может:

- Изменить ширину любого столбца (для этого необходимо перетащить границу его заголовка до нужной ширины).

- Сортировать таблицу (для этого необходимо выбрать иконку рядом с заголовком сортируемого столбца).
- Скрывать/отображать столбцы и изменять их порядок в окне **Вид таблицы** (для открытия окна необходимо выбрать **Настроить Таблицу** в правом верхнем углу таблицы; при выборе иконки столбец скрывается, при наведении на иконку активируется возможность перемещения столбца).
- Сбросить внесенные изменения также в окне **Вид таблицы** (для этого выбрать кнопку «**Сбросить**»).

Для быстрого поиска объекта в таблице предусмотрено поле  в правой верхней части таблицы.

Объекты таблицы можно выгрузить в формате Excel. Для этого нужно выбрать иконку **Экспорта в excel** .

#### **4.1.2. Создание проекта DM**

Для создания нового проекта необходимо выполнить следующие шаги:

- Выбрать кнопку «**Добавить**» в верхней части таблицы с Проектами.
- В открывшемся окне **Создание проекта** задать название и описание.
- Сохранить изменения.

#### **4.1.3. Удаление проекта DM**

Для удаления проекта необходимо выполнить следующие шаги:

- Выбрать чекбокс рядом с удаляемым проектом.
- Выбрать кнопку «**Удалить**».

#### **4.1.4. Копирование проекта DM**

Для копирования проекта необходимо выполнить следующие шаги:

- С правой стороны от проекта выбрать иконку в виде трех вертикальных точек и в выпадающем меню выбрать **Копировать**.
- В открывшемся окне **Копирование проекта** задать название и описание нового проекта.
- Сохранить изменения.

#### **4.1.5. Редактирование проекта DM**

Для переименования и изменения описания проекта необходимо выполнить следующие шаги:

- С правой стороны от проекта выбрать иконку в виде трех вертикальных точек и в выпадающем меню выбрать **Редактировать**.

- В открывшемся окне **Редактирование проекта** внести необходимые изменения в названии и описании проекта.
- Сохранить изменения.

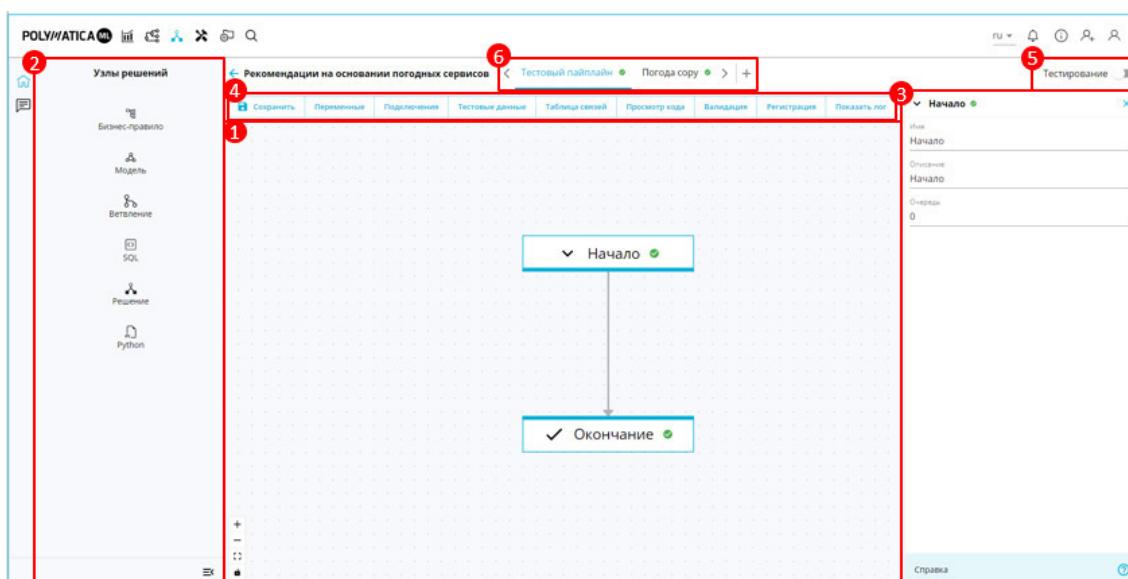
#### 4.1.6. Работа с проектом DM

Для начала работы с проектом необходимо выбрать его из списка доступных.

### 4.2. Конструктор цепочек решений

#### 4.2.1. Интерфейс Конструктора цепочек решений

Данный экран открывается при выборе любого из доступных проектов и представляет собой **Конструктор цепочек решений**.



**Рисунок 168 Элементы Конструктора цепочек решений**

Интерфейс **Конструктора цепочек решений** состоит из следующих элементов:

1. Рабочей области, в которую помещаются узлы решений между исходными узлами **Начало** и **Окончание**.
2. Левой боковой панели с узлами решений.
3. Правой боковой панели с настройками узла (открывается при выборе узла, размещенного в рабочем поле).
4. Верхней панели с кнопками для настройки цепочки решения.
5. Кнопка-переключатель режима Тестирования.
6. Панель создания нескольких сценариев в рамках одного проекта DM.

Для удобства пользования рабочей областью в левом нижнем углу предусмотрены кнопки масштабирования (кнопки

+

и

-

для приближения и отдаления, соответственно, кнопка

::

для масштабирования на объектах, расположенных в рабочем поле) и кнопка, блокирующая перемещение объектов в рабочем пространстве (

■

).

Левую боковую панель можно скрыть, выбрав кнопку

≡

в нижней части панели. Раскрыть скрытую панель можно выбрав

≡

также в нижней части панели.

Правая боковая панель открывается при выборе узла из рабочего поля. Скрыть эту панель можно выбрав иконку

×

в правом верхнем углу.

Панель с кнопками включает в себя:

- Кнопку **«Сохранить»**, для сохранения изменений в построенной цепочке решения.
- Кнопку **«Переменные»**, для добавления переменных, которые будут фигурировать в цепочке решений.
- Кнопку **«Подключения»**, для задания и отображения используемых в цепочке решений подключений к БД.
- Кнопку **«Тестовые данные»**, для задания тестовых данных на основе которых будет производиться тестирование цепочки (сценария).
- Кнопку **«Таблица связей»**, для корректирования связей между узлами решений.
- Кнопку **«Просмотр кода»**, для просмотра кода, лежащего за цепочкой решения.
- Кнопку **«Регистрация»**, для регистрации решения в Репозиторий ММ.

## 4.2.2. Узлы решений

### 4.2.2.1. Описание узлов решений

#### 4.2.2.1.1. Узел "Начало"

---

**Узел «Начало»** является базовым узлом цепочки решения.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения                 | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет       | Название узла, которое будет отображаться в интерфейсе |
| <b>Описание</b> | Ручной ввод<br>Ограничений на значение нет       | Описание узла  |
| <b>Очередь</b>  | Ручной ввод<br>Числовое неотрицательное значение | Порядок отображения при тестировании решения           |

**Таблица 54 Параметры узла "Начало"**

#### 4.2.2.1.2. Узел "Бизнес-правило"

---

**Узел "Бизнес-правила"** используется для реализации бизнес-логики в процессе принятия решения.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения                 | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет       | Название узла, которое будет отображаться в интерфейсе                 |
| <b>Описание</b> | Ручной ввод<br>Ограничений на значение нет       | Описание узла  |
| <b>Очередь</b>  | Ручной ввод<br>Числовое неотрицательное значение | Порядок отображения при тестировании решения                           |
| <b>Правила</b>  | Кнопка   | При выборе кнопки открывается окно <b>Правила</b> для задания условий. |

**Таблица 55 Параметры узла "Бизнес-правило"**

#### 4.2.2.1.3. Узел "Модель"

---

**Узел "Модель"** позволяет интегрировать публикацию версии модели из компонента **Управление моделями и решениями (ММ)**.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения           | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе |

| Параметр                | Возможные значения и ограничения                 | Описание   |
|-------------------------|--|--|
| <b>Описание</b>         | Ручной ввод<br>Ограничений на значение нет       | Описание узла  |
| <b>Очередь</b>          | Ручной ввод<br>Числовое неотрицательное значение | Порядок отображения при тестировании решения   |
| <b>Настройка модели</b> | Кнопка   | При выборе кнопки открывается окно <b>Настройка модели</b> для выбора необходимой публикации и мэппинга переменных |

**Таблица 56 Параметры узла "Модель"**

#### 4.2.2.1.4. Узел "Ветвление"

**Узел "Ветвление"** используется для обработки различных сценариев принятия решения.

**Список параметров узла** представлен в таблице ниже.

| Параметр                 | Возможные значения и ограничения                 | Описание  |
|--------------------------|--|---|
| <b>Название</b>          | Ручной ввод<br>Ограничений на значение нет       | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b>          | Ручной ввод<br>Ограничений на значение нет       | Описание узла   |
| <b>Очередь</b>           | Ручной ввод<br>Числовое неотрицательное значение | Порядок отображения при тестировании решения  |
| <b>Правила ветвления</b> | Кнопка   | При выборе кнопки открывается окно <b>Правила ветвления</b> для создания нескольких веток сценария. |

**Таблица 57 Параметры узла "Ветвление"**

#### 4.2.2.1.5. Узел "SQL"

**Узел "SQL"** используется для написания sql-запросов к базе данных для обогащения данных в процессе принятия решений.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения           | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет | Название узла, которое будет отображаться в интерфейсе |
| <b>Описание</b> | Ручной ввод<br>Ограничений на значение нет | Описание узла  |
| <b>Очередь</b>  | Ручной ввод                                | Порядок отображения при тестировании решения           |

| Параметр             | Возможные значения и ограничения  | Описание   |
|----------------------|-----------------------------------|--|
|                      | Числовое неотрицательное значение |  |
| <b>Сопоставление</b> | Кнопка                            | При выборе кнопки открывается окно <b>Сопоставление</b> для написания sql-запросов (подробное описание ниже) |

Таблица 58 Параметры узла "SQL"

### Окно Сопоставление

В окне **Сопоставление** Пользователь может написать sql-запрос.

Основными элементами окна **Сопоставление** являются:

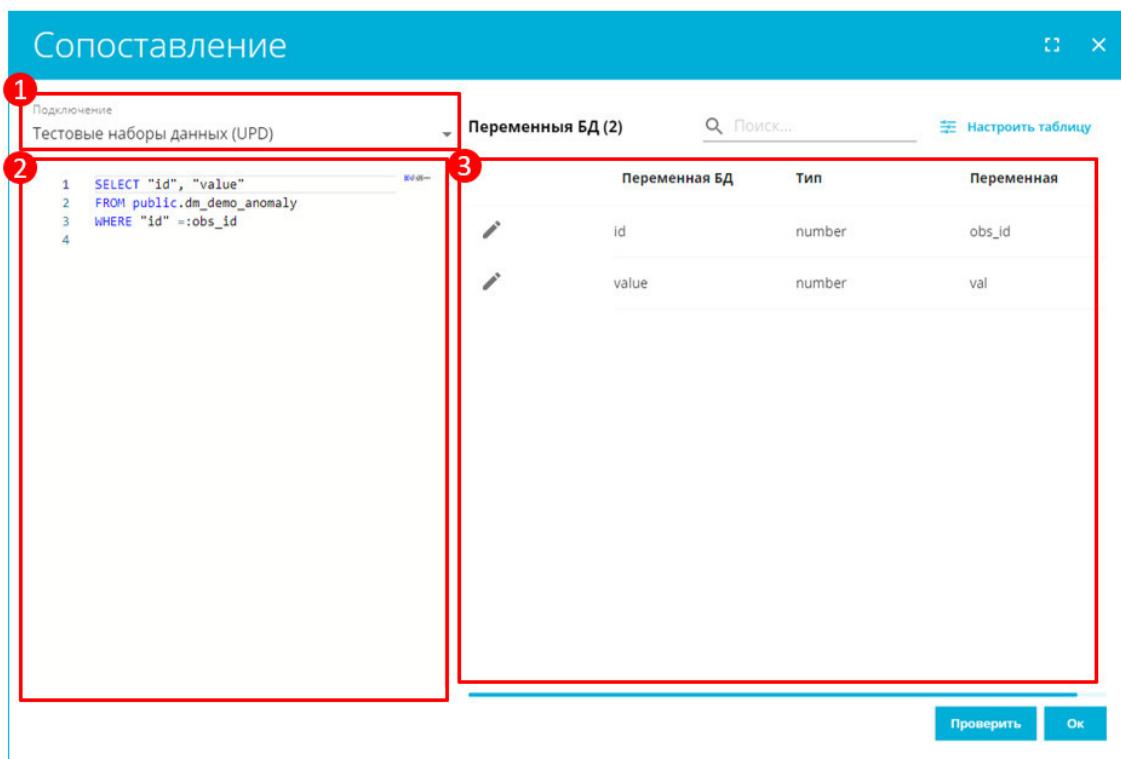


Рисунок 169 Окно Сопоставление

1. Список для выбора необходимого Подключения.
2. Поле ввода sql-запроса.
3. Мэппинг переменных БД и переменных сценария.

### Для задания sql-запроса необходимо:

1. Выбрать Подключение из списка доступных для сценария (задается на экране Конструктора сценария по кнопке Подключения).
2. Ввести необходимый запрос в поле ввода и нажать Проверить.
3. Если запрос корректный, то отобразится таблица, в которой необходимо задать мэппинг переменных БД и переменных сценария.
4. Сохранить.

#### 4.2.2.1.6. Узел "Решение"

**Узел "Решение"** используется для вызова другого решения и интеграции его в эту цепочку принятия решений.

**Список параметров узла** представлен в таблице ниже.

| Параметр             | Возможные значения и ограничения                 | Описание   |
|----------------------|--|--|
| <b>Название</b>      | Ручной ввод<br>Ограничений на значение нет       | Название узла, которое будет отображаться в интерфейсе   |
| <b>Описание</b>      | Ручной ввод<br>Ограничений на значение нет       | Описание узла  |
| <b>Очередь</b>       | Ручной ввод<br>Числовое неотрицательное значение | Порядок отображения при тестировании решения   |
| <b>Сопоставление</b> | Кнопка   | При выборе кнопки открывается окно <b>Сопоставление</b> для выбора встраиваемого решения из списка существующих и мэппинга переменных. |

**Таблица 59 Параметры узла "Решение"**

#### 4.2.2.1.7. Узел "Python"

**Узел "Python"** позволяет интегрировать произвольный python-код для реализации самых сложных нетиповых задач.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения                 | Описание  |
|-----------------|--|---|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет       | Название узла, которое будет отображаться в интерфейсе  |
| <b>Описание</b> | Ручной ввод<br>Ограничений на значение нет       | Описание узла   |
| <b>Очередь</b>  | Ручной ввод<br>Числовое неотрицательное значение | Порядок отображения при тестировании решения  |
| <b>Код</b>      | Кнопка   | При выборе кнопки открывается окно <b>Код</b> с редактором, куда необходимо добавить произвольный python-код. |

**Таблица 60 Параметры узла "Python"**

#### 4.2.2.1.8. Узел "Окончание"

**Узел «Окончание»** является базовым узлом цепочки решения.

**Список параметров узла** представлен в таблице ниже.

| Параметр        | Возможные значения и ограничения                 | Описание   |
|-----------------|--|--|
| <b>Название</b> | Ручной ввод<br>Ограничений на значение нет       | Название узла, которое будет отображаться в интерфейсе |
| <b>Описание</b> | Ручной ввод<br>Ограничений на значение нет       | Описание узла  |
| <b>Очередь</b>  | Ручной ввод<br>Числовое неотрицательное значение | Порядок отображения при тестировании решения           |

Таблица 61 Параметры узла "Окончание"

## 5. Компонент Управление моделями и решениями (Model Manager, ММ)

Основная задача, решаемая **компонентом Управление моделями и решениями** — ввод моделей в эксплуатацию и управление их жизненным циклом. В компоненте предусмотрен следующий функционал:

- Репозиторий моделей.
- Импорт сторонних python моделей.
- Публикация моделей.
- Согласование моделей.
- Мониторинг качества работы моделей.

### 5.1. Интерфейс ММ

Раздел **Управление моделями и решениями** открывается при выборе иконки



в левом верхнем меню.

| Наименование                                 | Описание   | Тип | Согласование | Дата обновления   | Создано           |
|--|--|-----|--------------|-------------------|-------------------|
| Ритеил С (канал коммуникации)                | Определение канала коммуникации:   |     |              | 02.09.2022, 19:23 | 17.05.2022, 15:48 |
| Внешняя модель Python - Отлик на предложение | Внешняя модель Python - отлик на предложение   |     |              | 13.10.2022, 15:45 | 20.07.2022, 16:41 |
| Отлик на предложение                         | Модель прогноза отклика на предложение в рамках кампании цевного маркетинга 'рост среднего чека клиента' |     |              | 25.07.2022, 17:04 | 17.05.2022, 15:50 |
| Ритеил С (контекст предложения)              | Определение контекста предложения:   |     |              | 18.05.2022, 18:39 | 18.05.2022, 18:38 |
| Определение контекста предложения 2          | Определение контекста предложения 2:   |     |              | 19.10.2022, 14:04 | 19.10.2022, 14:04 |
| Определение контекста предложения 1          | Определение контекста предложения 1:   |     |              | 18.10.2022, 11:55 | 29.09.2022, 8:42  |

**Рисунок 170** Интерфейс компонента Управление моделями

Основным элементом навигации в разделе **Управление моделями и решениями** является боковая панель, которая содержит следующие подразделы:

- **Управление моделями и решениями.**
  - **Репозиторий** — представляет собой централизованное хранилище моделей.
  - **Опубликованные модели** — представляет собой список моделей, которые были опубликованы.
  - **Согласование** — включает в себя функционал, необходимый для согласования моделей.

- **Настройки.**

- **Проекты** — представляет список проектов.
- **Библиотеки** — представляет собой список Python библиотек, используемых в моделях.
- **Шаблоны** — представляет собой список доступных для выполнения процессов согласования моделей.
- **Пользовательские атрибуты** — позволяют пользователю задать любую дополнительную информацию о версии модели.

Для удобства пользования левую боковую панель можно скрыть, выбрав кнопку



в нижней части панели. Раскрыть скрытую панель можно выбрав



также в нижней части панели.

Таблицы в разделе **Управление моделями и решениями** имеют гибкие настройки отображения. Так, пользователь может:

- изменить ширину любого столбца (для этого необходимо перетащить границу его заголовка до нужной ширины).
- сортировать таблицу (для этого необходимо выбрать иконку рядом с заголовком сортируемого столбца).
- скрывать/отображать столбцы и изменять их порядок в окне **Вид таблицы** (для открытия окна необходимо выбрать в правом верхнем углу таблицы; при выборе иконки столбец скроется, при наведении на иконку активируется возможность перемещения столбца).
- сбросить внесенные изменения также в окне **Вид таблицы** (для этого выбрать кнопку «**Сбросить**»).

Для быстрого поиска объекта в таблице предусмотрено поле  в правой верхней части таблицы.

Объекты таблицы можно выгрузить в формате Excel. Для этого нужно выбрать иконку **Экспорта в excel** .

## 5.2. Репозиторий

**Репозиторий** представляет собой централизованное хранилище моделей и их версий и позволяет управлять их жизненным циклом.

Перейти в **Репозиторий** можно, выбрав одноименный раздел боковой панели.

| Наименование...                                | Описание   | Тип | Согласование | Дата обновления   | Создано           |
|--|--|-----|--------------|-------------------|-------------------|
| Ритейл CI (канал коммуникаций)                 | Определение канала коммуникации  |     |              | 02.09.2022, 19:23 | 17.05.2022, 15:48 |
| Внешняя модель Python - отклики на предложение | Внешняя модель Python - отклики на предложение   |     |              | 13.10.2022, 15:45 | 20.07.2022, 16:41 |
| Отклики на предложение                         | Модель прогноза отклика на предложениях в рамках кампании ценового маркетинга "тест среднего чека клиента" |     |              | 25.07.2022, 17:04 | 17.05.2022, 15:50 |
| Ритейл CI (контекст предложения)               | Определение контекста предложения  |     |              | 18.05.2022, 18:39 | 18.05.2022, 18:38 |
| Определение контекста предложения 2            | Определение контекста предложения 2  |     |              | 19.10.2022, 14:04 | 19.10.2022, 14:04 |
| Определение контекста предложения              | Определение контекста предложения 1  |     |              | 18.10.2022, 11:55 | 29.09.2022, 8:42  |

**Рисунок 171 Выбор раздела Репозиторий боковой панели**

Репозиторий состоит из двух основных элементов:

1. Кольцевые диаграммы, отражающие характеристики моделей в Репозитории (открывается при выборе кнопки **Показать графики** в верхней части раздела):
  - Статусы согласования версии модели (возможные значения — Draft (черновик), Approved (согласовано), Rejected (отказано в согласовании), Awaiting (в процессе согласования)).
  - Используемые алгоритмы.
  - Библиотеки.
  - Статусы сборок для публикации (возможные значения — Ready (готово), Error (ошибка), InProgress (в процессе)).
2. Список доступных пользователю моделей в табличном виде.

| Наименование...                                | Описание   | Тип | Согласование | Дата обновления   | Создано           |
|--|--|-----|--------------|-------------------|-------------------|
| Ритейл CI (канал коммуникаций)                 | Определение канала коммуникации  |     |              | 02.09.2022, 19:23 | 17.05.2022, 15:48 |
| Внешняя модель Python - отклики на предложение | Внешняя модель Python - отклики на предложение   |     |              | 13.10.2022, 15:45 | 20.07.2022, 16:41 |
| Отклики на предложение                         | Модель прогноза отклика на предложениях в рамках кампании ценового маркетинга "тест среднего чека клиента" |     |              | 25.07.2022, 17:04 | 17.05.2022, 15:50 |

**Рисунок 172 Раздел Репозиторий**

В таблице отображается иерархическая структура Репозитория в соответствии с рисунком ниже. При нажатии на кнопку список раскрывается и сворачивается.

| Модели (23)              |                      | <a href="#">Импортировать</a> |
|--------------------------|----------------------|-------------------------------|
| <input type="checkbox"/> | Наименование         |                               |
| <input type="checkbox"/> | Задача Качество вина | <b>Проект</b>                 |
| <input type="checkbox"/> | Модель Качество вина | <b>Модель</b>                 |
| <input type="checkbox"/> | Модель Качество вина | <b>Версия модели</b>          |

**Рисунок 173 Пример иерархии объектов в Репозитории**

**Проект ММ** представляет собой логическое пространство, объединяющее несколько моделей по конкретной бизнес-задаче.

**Модель** — метаинформация, описывающая программное решение, позволяющее произвести трансформацию или анализ переменных, которые она ожидает на вход (Входные Переменные) и вернуть результаты (Выходные Переменные).

**Версия модели** — итеративная реализация модели.

В рамках **Проекта** может быть несколько **Моделей** и **Версий модели** (например, первоначальная версия модели и версия модели, полученная после дообучения на новой порции данных).

Для того, чтобы расширить таблицу с моделями и проектами, необходимо убрать

панель с графиками, выбрав кнопку [Скрыть Графики](#).

В Репозитории можно хранить как модели построенные и зарегистрированные из Model Designer и Decision Manager, так и сторонние python модели (и их версии).

## 5.3. Проект

Для импорта модели в **Репозиторий** необходимо первоначально создать **Проект**.

**Проект** представляет собой логическое пространство, объединяющее несколько моделей по конкретной бизнес-задаче. Параметры **Проекта** задают минимальные вычислительные ресурсы, с которыми запустится несогласованная модель (количество памяти (в Мб) и ядер), и шаблоны согласования версий модели и их публикаций.

При выборе пункта **Проекты** боковой панели открывается одноименный раздел со списком созданных Проектов ММ.

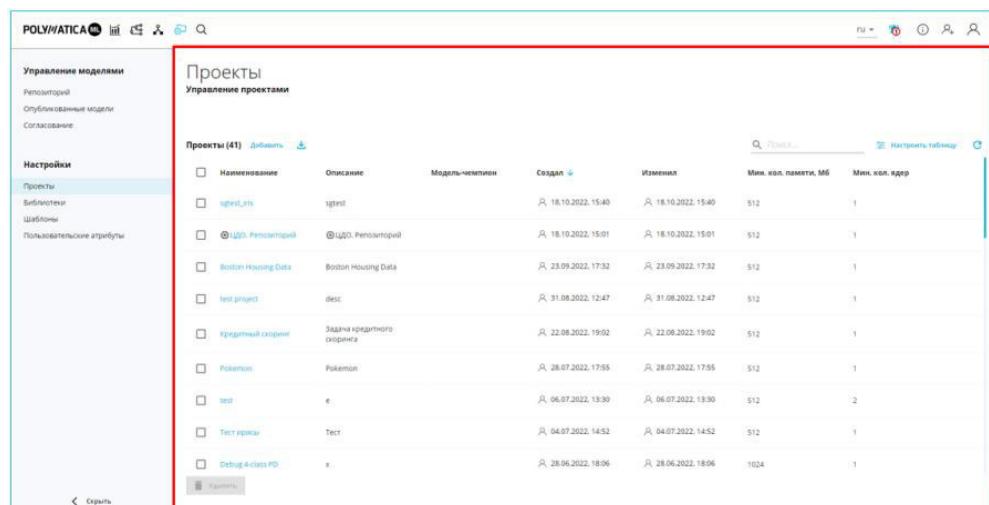


Рисунок 174 Раздел проекты

### 5.3.1. Создание проекта ММ

Для создания нового Проекта необходимо:

- Выбрать «**Добавить**» в верхней части таблицы раздела **Проекты**.
- В открывшемся окне **Создание проекта** задать:
  - параметры – **Имя проекта**, **Описание**, **Код** (Идентификатор), **Минимальное количество памяти**, **Минимальное количество ядер**;
  - шаблоны **согласования версии модели и публикаций** из списка
- Сохранить изменения.

Новый проект отобразится в списке созданных проектов раздела **Проекты**.

Созданный проект отобразится в **Репозитории** только после импорта в него Модели.

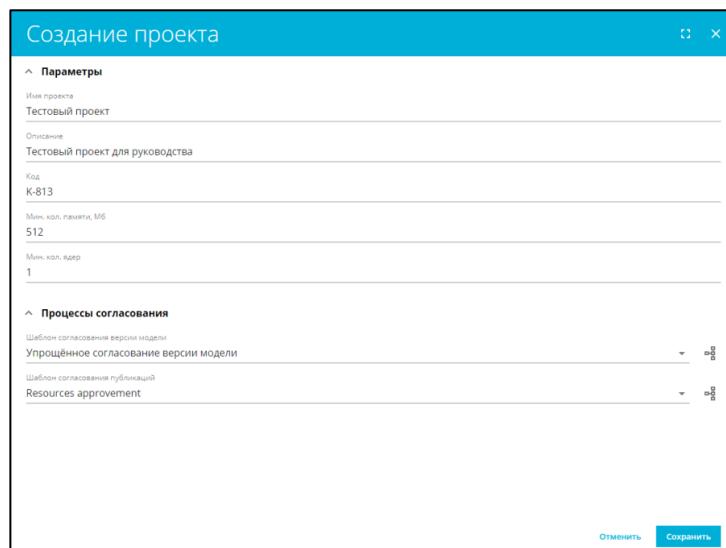


Рисунок 175 Окно Создание проекта



Ознакомиться с выбранным шаблоном согласования можно, выбрав кнопку рядом с полем выбора шаблона, а также в разделе **Шаблоны** боковой панели (подробнее в разделе [Шаблон согласования](#)).

### 5.3.2. Удаление проекта ММ

Для удаления Проекта необходимо:

- Выбрать чекбокс рядом с удаляемым проектом.
- Выбрать кнопку «**Удалить**».

| Наименование   | Описание                          | Модель-чемпион         | Создал            | Изменил           | Мин. кол. памяти, Мб | Мин. кол. ядер |
|--|-----------------------------------|------------------------|-------------------|-------------------|----------------------|----------------|
| <input checked="" type="checkbox"/> Ритейл CI (контекст предложения) | Определение контекста предложения |                        | 18.05.2022, 18:38 | 18.05.2022, 18:39 | 512                  | 1              |
| <input type="checkbox"/> Ритейл CI (канал коммуникаций)              | Определение канала коммуникаций   | Отклики на предложение | 17.05.2022, 15:48 | 02.09.2022, 19:23 | 512                  | 1              |

**Рисунок 176 Удаление проекта ММ**

В случае, если в Проекте имеются Модели, Пользователь не сможет его удалить. Для удаления Проекта сначала необходимо удалить Модели.

### 5.3.3. Редактирование проекта ММ

Параметры ранее созданного Проекта можно редактировать. Помимо изменения шаблонов согласования и минимальных вычислительных ресурсов, можно задать Модель-чемпион из соответствующего списка с моделями Проекта.

Для редактирования параметров проекта необходимо:

- Выбрать интересующий проект из списка.
- В открывшемся окне **Обновление проекта** внести необходимые изменения.
- Сохранить изменения.

Обновление проекта

▲ Параметры

Имя проекта  
Ритейл CI (контекст предложения)

Описание  
Определение контекста предложения

Код  
A-2

Мин. кол. памяти, Мб  
512

Мин. кол. ядер  
1

Модель-чемпион

▲ Процессы согласования

Шаблон согласования версии модели

Шаблон согласования публикаций

[Отменить](#) [Сохранить](#)

**Рисунок 177 Окно Обновление проекта**

## 5.4. Модель

**Модель** — метаинформация, описывающая программное решение, позволяющее произвести трансформацию или анализ переменных, которые она ожидает на вход (Входные Переменные) и вернуть результаты (Выходные Переменные).

### 5.4.1. Импорт модели

Зарегистрировать модель в Репозиторий можно из компонентов Построение моделей (MD) и Управление моделями и решениями (DM), а также импортировать извне.

**ВАЖНО!** Перед тем как импортировать модель необходимо создать **Проект ММ**.

#### 5.4.1.1. Регистрация модели из MD

Регистрация модели из компонента **Построение моделей** (Model Designer) производится при помощи узла **«Регистрация модели»** (подробнее раздел Узел «Регистрация модели»), при этом **Входные и Выходные переменные** подтягиваются автоматически.

Модель типа MD в Репозитории имеет иконку  в соответствующем столбце.

#### 5.4.1.2. Регистрация модели из DM

Регистрация модели из компонента **Разработка решений** (Decision Manager) производится при помощи кнопки **Регистрация** в Конструкторе решений.

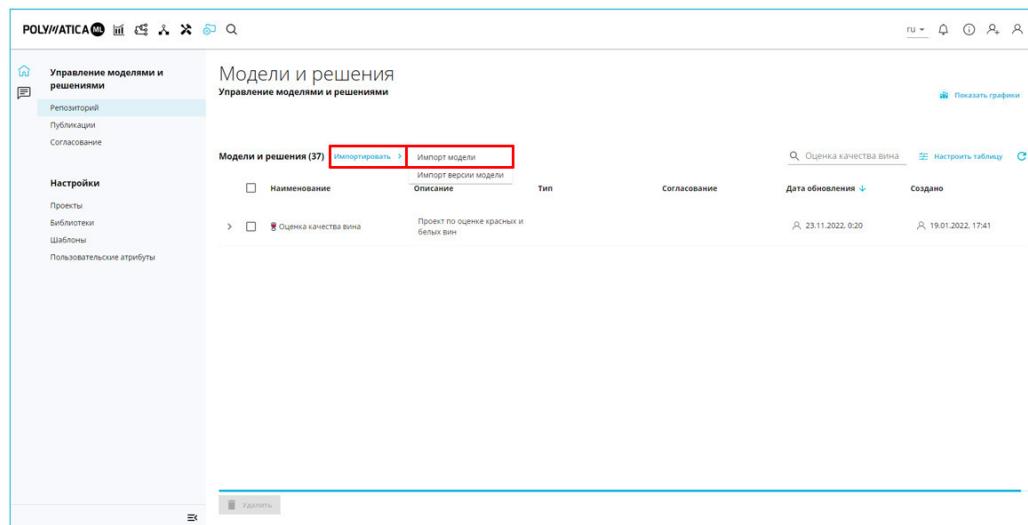
Модель типа DM в Репозитории имеет иконку  в соответствующем столбце.

#### 5.4.1.3. Импорт внешней модели

Импортировать в Модуль можно лишь модели, реализованные на Python.

Для импорта внешней модели необходимо:

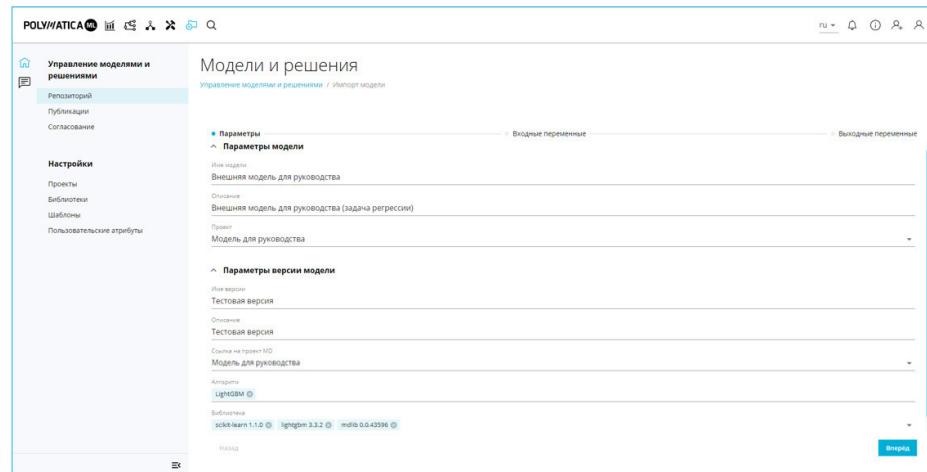
1. В **Репозитории** компонента **Управление моделями и решениями** выбрать кнопку **Импортировать** и в выпадающем меню выбрать пункт **Импорт модели**.



**Рисунок 178 Выбор кнопки Импорта модели**

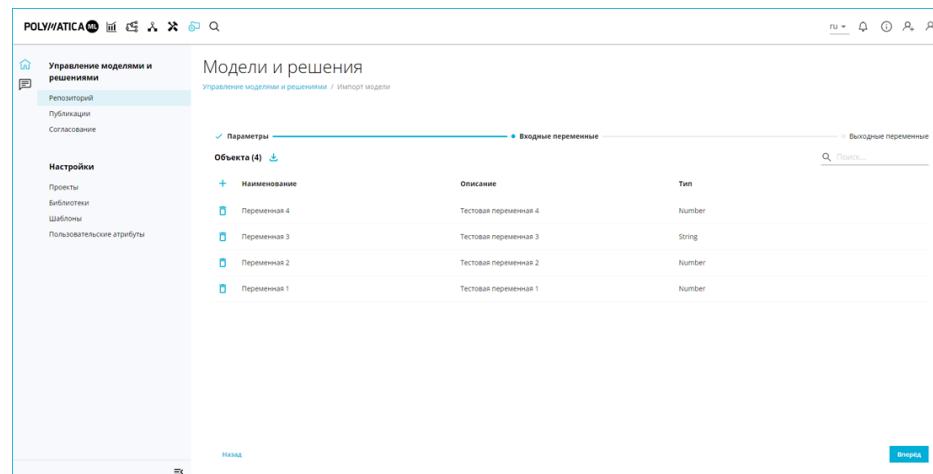
2. В открывшемся экране последовательно задать:

- Параметры модели – Имя и Описание модели, выбрать существующий Проект.
- Параметры версии проекта – Имя и Описание версии; проект MD; Алгоритм, лежащий в основе и библиотеки, в которых реализован используемый алгоритм.



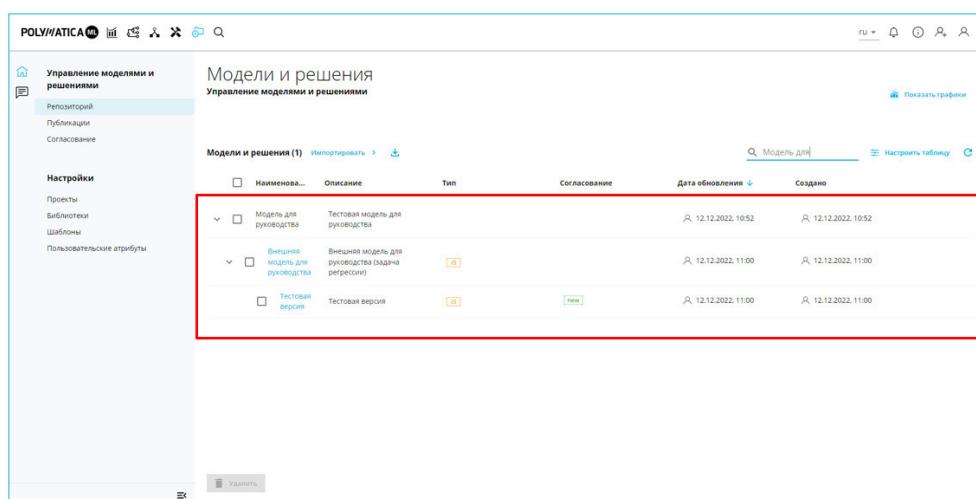
**Рисунок 179 Пример экрана импорта с параметрами модели**

- Входные и Выходные переменные.



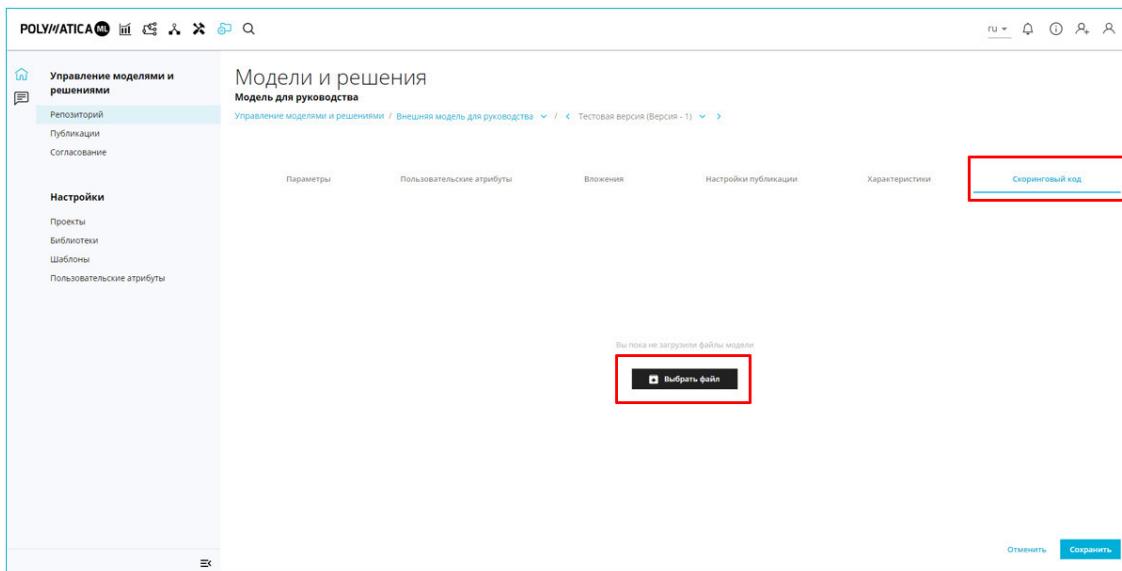
**Рисунок 180 Пример экрана с входными переменными**

3. Завершить импорт. Модель появится в списке Репозитория в соответствующем Проекте с типом MD.



**Рисунок 181 Импортированная в репозиторий внешняя модель**

4. Выбрать созданную версию модели, раскрыв иерархию Проекта и созданной модели.
5. В открывшемся окне перейти во вкладку **Скоринговый код** и загрузить архив. Архив должен содержать:
  - pkl-файлы модели.
  - файл со скоринговым кодом score.py.
  - файл инициализации модели \_\_init\_\_.py.



**Рисунок 182 Вкладка скоринговый код внешней импортированной модели**

### 5.4.2. Удаление модели

Для удаления Модели необходимо:

- Убедиться в том, что у удаляемой Модели нет Публикаций. Если Публикации есть – первоначально удалить их.
- Выбрать чекбокс рядом с удаляемой Моделью.
- Выбрать кнопку «**Удалить**».

Выбранная Модель совместно с ее Версиями пропадут из списка Репозитория.

**ВАЖНО!** Вместе с Моделью удаляются и ее версии.

### 5.4.3. Работа с моделью

Для редактирования и работы с моделью необходимо выбрать ее из списка Репозитория.

| Наименование                          | Описание  | Тип           | Согласование | Дата обновления   | Создано           |
|---------------------------------------|---|---------------|--------------|-------------------|-------------------|
| Модель для руководства                | Тестовая модель для руководства                   | Проект ММ     |              | 12.12.2022, 10:52 | 12.12.2022, 10:52 |
| Внешняя модель для руководства        | Внешняя модель для руководства (задача регрессии) | Модель        |              | 12.12.2022, 11:00 | 12.12.2022, 11:00 |
| Тестовая версия                       | Тестовая версия                                   | Версия модели |              | 12.12.2022, 11:00 | 12.12.2022, 11:00 |
| Внешняя версия модели для руководства | Тестовая модель для руководства                   | Версия модели |              | 12.12.2022, 10:56 | 12.12.2022, 10:56 |

Рисунок 183 Пример выбора модели

Открывшийся экран представляет собой набор вкладок, которые позволяют настраивать Модель.

Настройки разных типов моделей (MD или DM) несколько отличаются - у **моделей DM** отсутствует вкладка **Дополнительные функции**.

Первая вкладка **Параметры** отображает тип (MD или DM), имя и описание Модели, а также имеющиеся Версии модели.

Настройки Модели и решения:

- Модель для руководства
- Параметры:
  - Тип: MD
  - Имя модели: Внешняя модель для руководства
  - Описание: Внешняя модель для руководства (задача регрессии)
- Версии (2):
 

| Название                              | Описание                        | Версия | Создал            | Изменил           |
|---------------------------------------|---------------------------------|--------|-------------------|-------------------|
| Тестовая версия                       | Тестовая версия                 | 1      | 12.12.2022, 11:00 | 12.12.2022, 11:00 |
| Внешняя версия модели для руководства | Тестовая модель для руководства | 2      | 12.12.2022, 10:56 | 12.12.2022, 10:56 |

Рисунок 184 Пример экрана с настройками Модели типа MD

Вкладка **Пользовательские атрибуты** позволяет указать список атрибутов для версий текущей модели. Добавить необходимые атрибуты для отображения в данной вкладке можно в разделе **Пользовательские атрибуты** бокового экрана MD (подробнее в разделе [Пользовательские атрибуты](#)).

The screenshot shows the POLYMATICA web interface. In the top navigation bar, there are icons for home, search, and user profile, along with language selection (ru). The main menu on the left includes 'Управление моделями и решениями' (Management of Models and Solutions), 'Репозиторий' (Repository), 'Публикации' (Publications), and 'Согласование' (Approval). Under 'Настройки' (Settings), there are links for 'Проекты' (Projects), 'Библиотеки' (Libraries), 'Шаблоны' (Templates), and 'Пользовательские атрибуты' (User Attributes). The central content area is titled 'Модели и решения' (Models and Solutions) and 'Модель для руководства' (Model for Guidance). A breadcrumb navigation shows 'Управление моделями и решениями / Внешняя модель для руководства'. The main table has tabs at the top: 'Параметры' (Parameters), 'Пользовательские атрибуты' (User Attributes) which is highlighted with a red box, 'Входные переменные' (Input Variables), 'Выходные переменные' (Output Variables), and 'Дополнительные функции' (Additional Functions). The 'Параметры' tab displays three columns: 'Название' (Name), 'Описание' (Description), and 'Тип' (Type). The entries are: 'Приоритет' (Priority) with type 'numeric'; 'Дата ревизии' (Review Date) with type 'datetime'; and 'Разработчик' (Developer) with type 'string'.

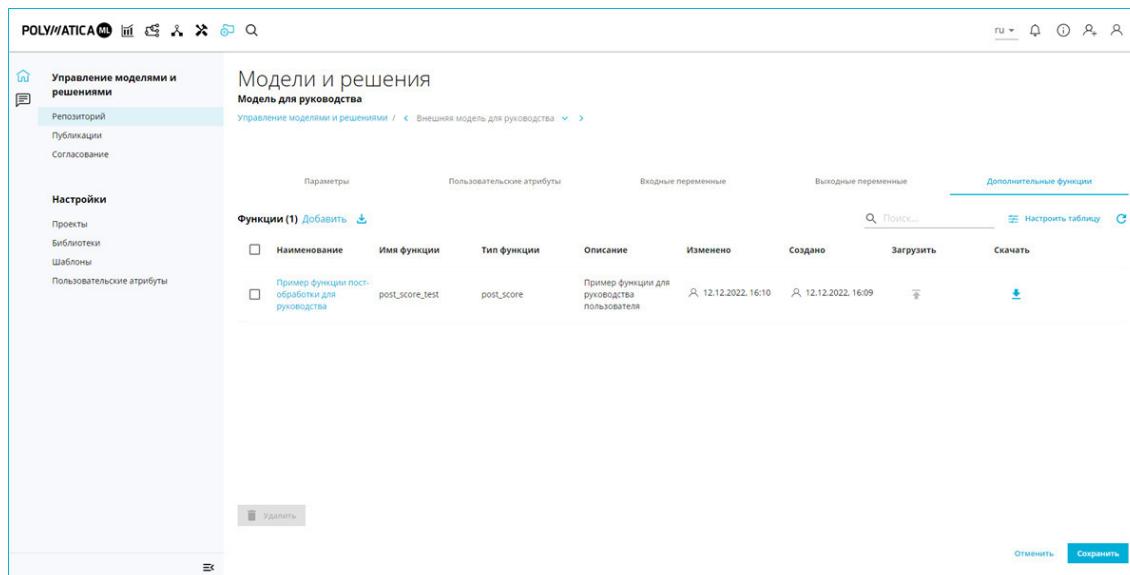
**Рисунок 185 Пример вкладки Пользовательские атрибуты**

Вкладки **Входные** и **Выходные переменные** имеют одинаковый интерфейс, в котором отображаются автоматически подтянутые из MD или DM входные/выходные переменные или заданные при импорте сторонней модели. Они будут одинаковы для всех версий текущей модели.

This screenshot shows the same POLYMATICA interface as the previous one, but with the 'Входные переменные' (Input Variables) tab selected in the 'Parameters' section. The table structure is identical, with columns for 'Наименование' (Name), 'Описание' (Description), and 'Тип' (Type). The entries listed are: 'Переменная 4' (Variable 4) with description 'Тестовая переменная 4' (Test variable 4) and type 'Number'; 'Переменная 2' (Variable 2) with description 'Тестовая переменная 2' (Test variable 2) and type 'Number'; 'Переменная 3' (Variable 3) with description 'Тестовая переменная 3' (Test variable 3) and type 'String'; and 'Переменная 1' (Variable 1) with description 'Тестовая переменная 1' (Test variable 1) and type 'Number'.

**Рисунок 186 Пример вкладки Входные переменные**

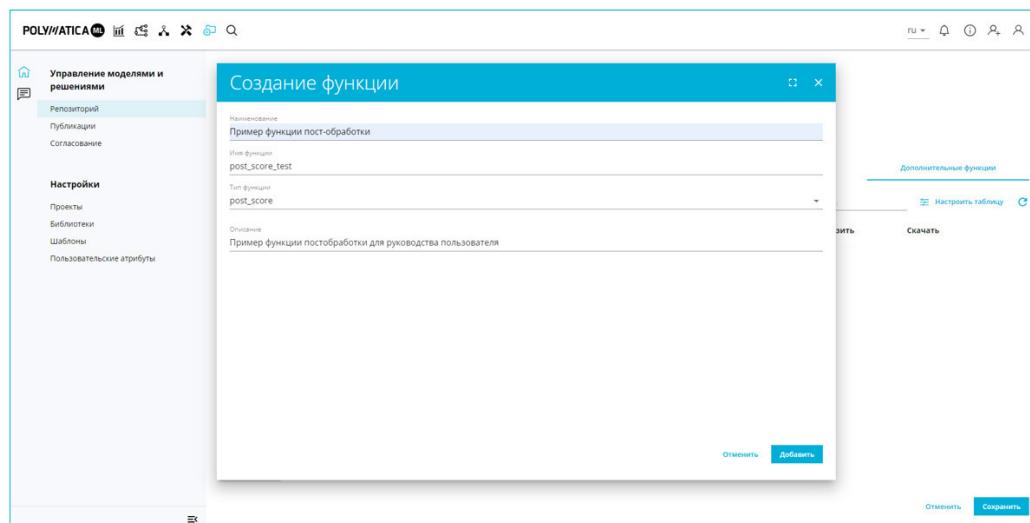
Последняя вкладка **Дополнительные функции** позволяет добавить pre- и post-score функции к исходному score коду версий модели MD.



**Рисунок 187 Пример вкладки Дополнительные функции**

Для добавления функции необходимо:

- Выбрать кнопку **Добавить** (ниже панели с вкладками).
- В открывшемся окне **Создание функции** задать Наименование и Описание функции, выбрать ее Тип (pre-score или post-score), а также задать ее имя. **Добавить изменения.**



**Рисунок 188 Окно Создание функции**

- В списке вкладки **Дополнительные функции** отобразится созданная функция. В столбце **Загрузить** необходимо выбрать иконку , выбрать необходимый файл с расширением .ru и вновь нажать на эту же иконку.

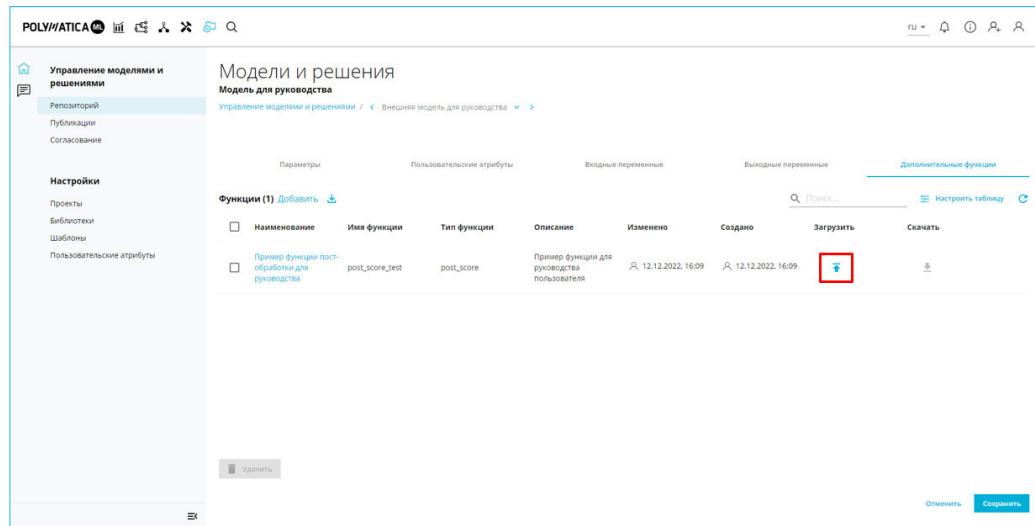


Рисунок 189 Пример загрузки функции пост-обработки

- После загрузки файла с функцией активизируется иконка в столбце Скачать.

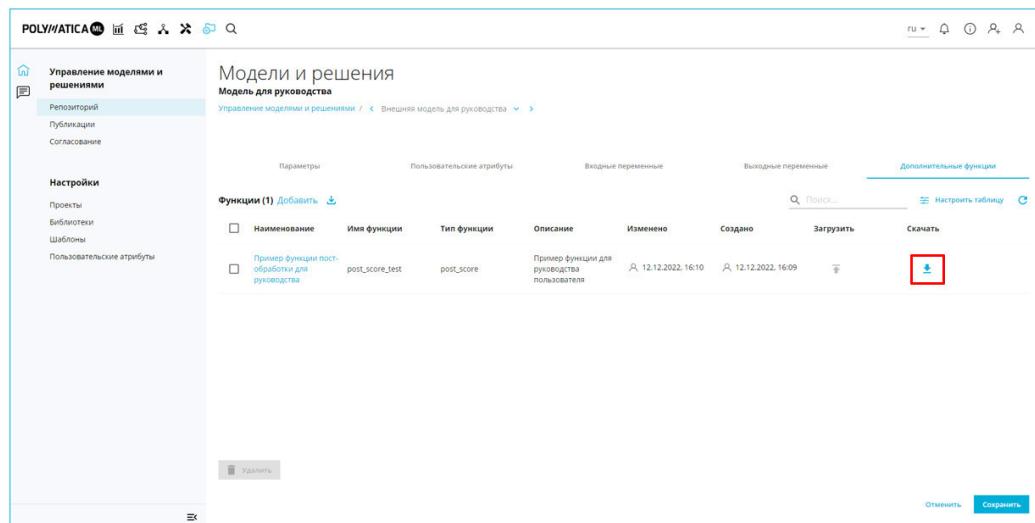


Рисунок 190 Пример выгрузки функции пост-обработки

## 5.5. Версия модели

**Версия модели** является итеративной реализацией модели.

### 5.5.1. Импорт версии модели

Зарегистрировать **Версию модели** в **Репозиторий** можно аналогично **Модели** - из компонентов **Построение моделей (MD)** и **Управление моделями и решениями (DM)**, а также импортировать извне.

### 5.5.2. Удаление версии модели

Для удаления Версии модели необходимо:

- Убедиться в том, что у удаляемой Версии модели нет Публикаций. Если Публикации есть – снять их с публикации и удалить.
- Выбрать чекбокс рядом с удаляемой Версию модели.
- Выбрать кнопку «Удалить».

### 5.5.3. Работа с версией модели

Для редактирования и работы с моделью необходимо выбрать ее из списка Репозитория.

| <input type="checkbox"/> Наименование                          | Описание  | Тип       | Согласование  | Дата обновления   | Создано           |
|--|---|-----------|---------------|-------------------|-------------------|
| <input type="checkbox"/> Модель для руководства                | Тестовая модель для руководства                   | Проект ММ |               | 12.12.2022, 10:52 | 12.12.2022, 10:52 |
| <input type="checkbox"/> Внешняя модель для руководства        | Внешняя модель для руководства (задача регрессии) |           | Модель        | 12.12.2022, 11:00 | 12.12.2022, 11:00 |
| <input type="checkbox"/> Тестовая версия                       | Тестовая версия                                   |           | Версия модели | 12.12.2022, 11:00 | 12.12.2022, 11:00 |
| <input type="checkbox"/> Внешняя версия модели для руководства | Тестовая модель для руководства                   |           | Версия модели | 12.12.2022, 10:56 | 12.12.2022, 10:56 |

**Рисунок 191 Пример выбора версии модели**

Открывшийся экран представляет собой набор вкладок, которые позволяют настраивать Версию модели. Они различаются для Версий модели MD и DM.

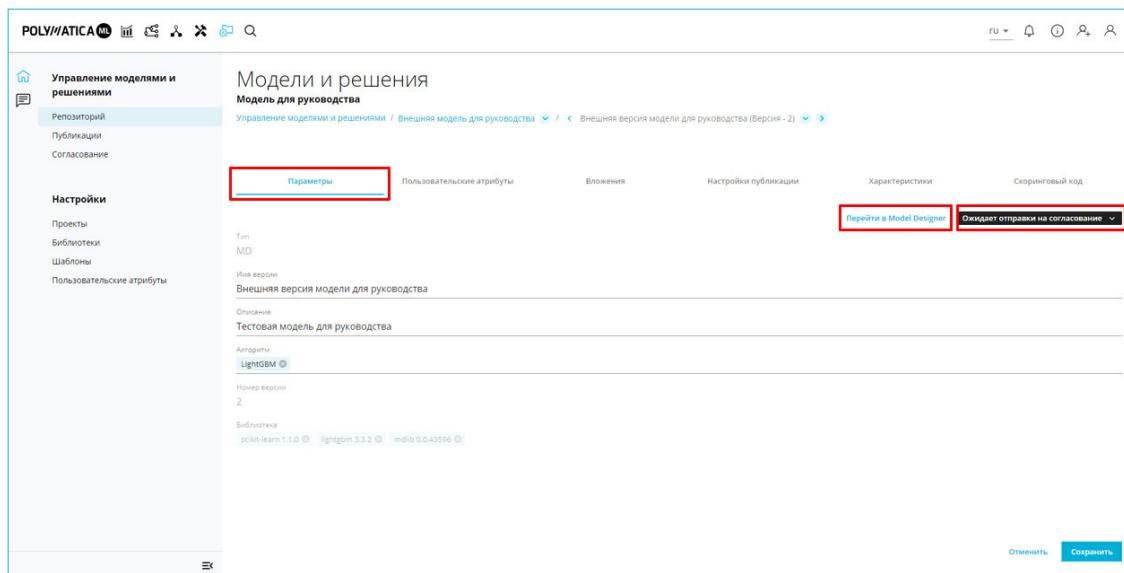
Набор вкладок для Версии модели MD:

- Параметры
- Пользовательские атрибуты
- Вложения
- Настройки публикации
- Характеристики
- Скоринговый код

Набор вкладок для Версии модели DM:

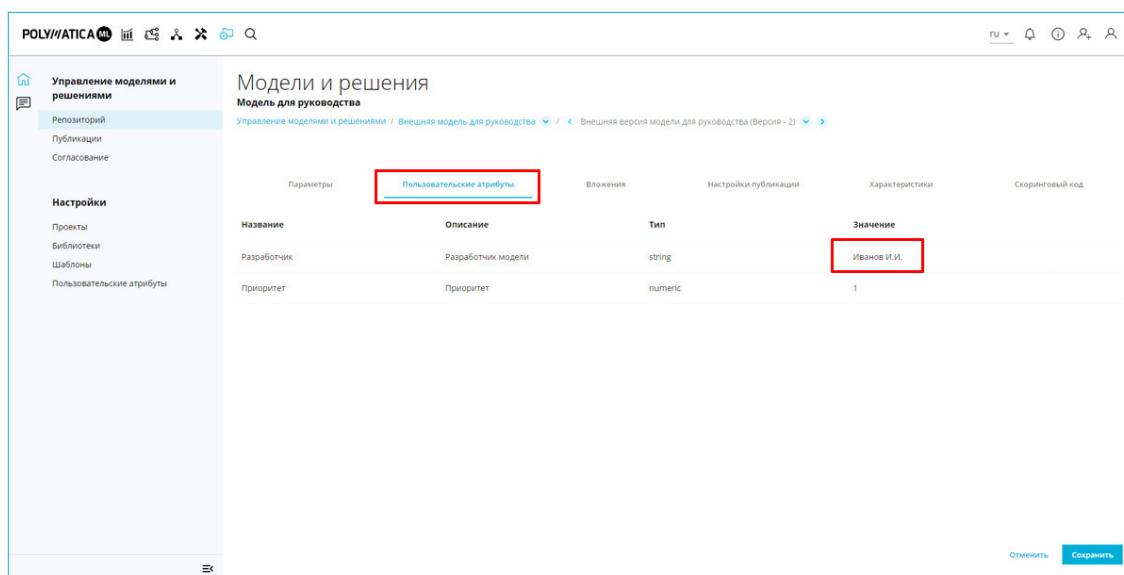
- Параметры
- Пользовательские атрибуты
- Подключения
- Вложения
- Настройки публикации
- Скоринговый код

На вкладке **Параметры** отображены Тип, Имя, Описание и номер версии модели, а также используемые алгоритмы и библиотеки. Помимо этого здесь есть возможность перейти в проект MD, из которого была зарегистрирована версия (кнопка Перейти в **Model Designer**), а также отправить версию модели на согласование (кнопка **Ожидает отправки на согласование**, подробнее о процессе согласования в разделе [Согласование версий модели и публикаций](#)).



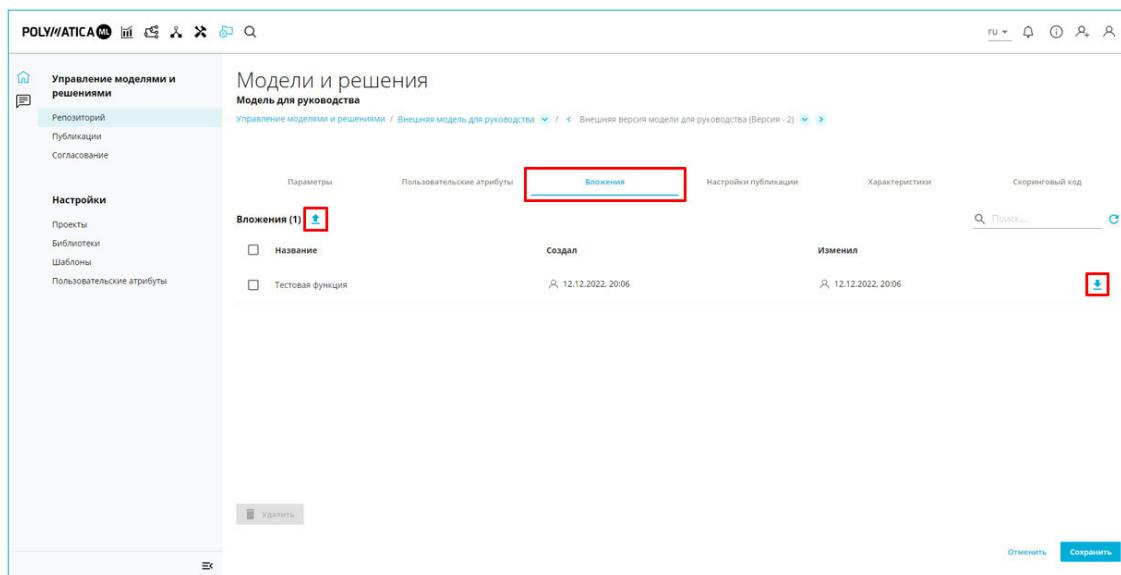
**Рисунок 192 Вкладка Параметры**

Вкладка **Пользовательские атрибуты** позволяет задать конкретные значения для выбранных атрибутов. Для этого нужно нажать на пустое поле столбца **Значение** и ввести необходимое значение.



**Рисунок 193 Вкладка Пользовательские атрибуты**

Вкладка **Вложения** позволяет добавить необходимые артефакты (документацию, инструкции, список и объяснение атрибутов и т.д.) к версии модели. Для этого необходимо выбрать иконку и загрузить необходимые файлы. Далее любой Пользователь, имеющий доступ к данной версии модели, сможет скачать необходимый файл, выбрав иконку рядом.



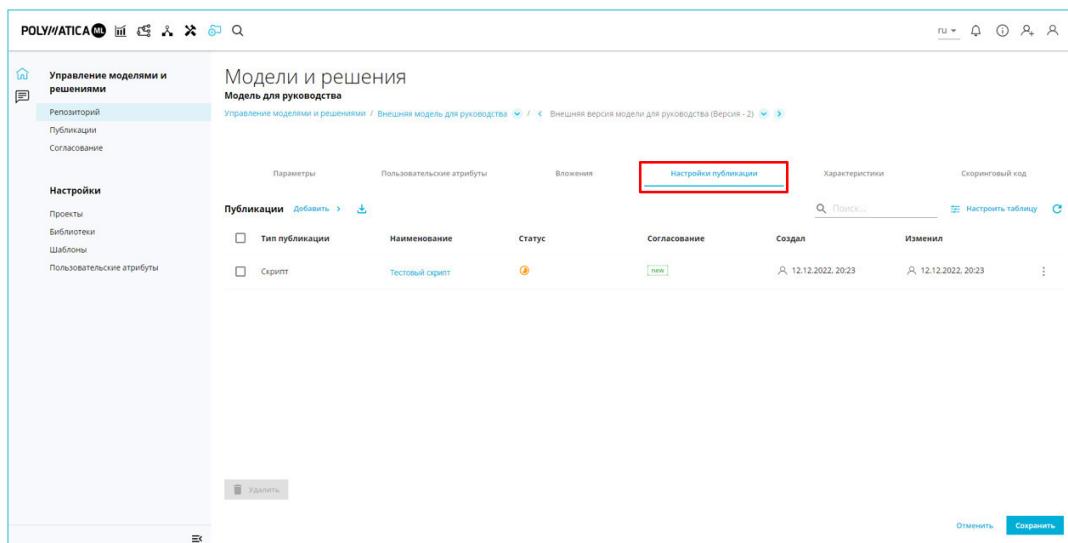
**Рисунок 194 Вкладка Вложения**

Вкладка **Настройки публикации** позволяет опубликовать версию модели в виде:

- Сервиса.
- Пакетного обработчика (Пакетная публикация).
- Скрипта (только для типа модели MD).
- Оценить качество модели (Оценка характеристик) (только для типа модели MD).
- Polymatica Analytics (только для типа модели MD).

Подробнее про публикацию модели в разделе [Публикация версии модели](#).

Также в меню рядом с публикацией можно инициировать процесс согласования публикации версии модели (подробнее в разделе [Согласование версий модели и публикаций](#)).



**Рисунок 195 Вкладка Настройки публикации**

Вкладка **Характеристики** актуальна только для моделей типа MD и содержит в себе результаты оценки качества модели при расчетах в MD (в столбце **Тип выборки** - Обучение, Тест, Валидация), а также при публикации модели типа Оценка характеристик (в столбце **Тип выборки** - Скоринг).

Управление моделями и решениями

## Модели и решения

Оценка качества вина

Управление моделями и решениями / Регрессионная модель качества вина (Версия - 6) / Регрессионная модель качества вина (Версия - 6)

Настройки

Характеристики (4)

|                           | Тип выборки | Подключение                  | Имя таблицы            | Дата создания        | MSE   | RMSE  | MAE   | MAPE  | R2    |
|---------------------------|-------------|------------------------------|------------------------|----------------------|-------|-------|-------|-------|-------|
| Пользовательские атрибуты | Обучение    | Тестовые наборы данных (UPD) | Качество красного вина | 18.01.2022, 07:36:45 | 0.413 | 0.643 | 0.488 | 0.091 | 0.362 |
|                           | Тест        | Тестовые наборы данных (UPD) | Качество красного вина | 18.01.2022, 07:36:45 | 0.471 | 0.686 | 0.53  | 0.096 | 0.294 |
|                           | Валидация   | Тестовые наборы данных (UPD) | Качество красного вина | 18.01.2022, 07:36:45 | 0.528 | 0.726 | 0.567 | 0.103 | 0.178 |
|                           | Скоринг     | Тестовые наборы данных (UPD) | winequality_white      | 18.05.2022, 21:05:43 | 0.743 | 0.862 | 0.657 | 0.112 | 0.052 |

Характеристики

Скоринговый код

### **Рисунок 196 Вкладка Характеристики**

Во вкладке **Скоринговый код** отображен код модели. Для раскрытия архива с моделью необходимо выбрать иконку  **core**. В раскрывшемся списке можно увидеть pk1-файлы (просмотр данных файлов не предусмотрен), файл инициализации модели, а также файл со скоринговым кодом score.py. Два последних файла можно открыть и посмотреть их содержимое (изменить содержимое в данном редакторе невозможно).

В файле `score.py` можно заменить функции `pre_score` и `post_score`. Именно сюда добавляются пользовательские функции, загруженные в модель (подробнее в разделе [Работа с моделью](#)).

Управление моделями и решениями

## Модели и решения

Оценка качества вина

управление моделями и решениями / Регрессионная модель качества вина < / > Регрессионная модель качества вина (Версия - 6) < / >

Настройки Параметры Пользовательские атрибуты Вложения Настройки публикации Характеристики Скоринговый код

Настройки

Проекты

Библиотеки

Шаблоны

Пользовательские атрибуты

core

core.py

```
1 import os
2 from typing import Dict, Any
3
4 import joblib
5 import numpy
6 import pandas
7
8 def on_init() -> Dict[Any, Any]:
9     script_dir = os.path.dirname(__file__)
10    models = {}
11    pkls = [ "85d10b82-6ac9-4e6c-aa60-a525c964285e", "1250d275-23d0-427b-b160-b337149b7ef3.pkl" ]
12
13    for pkl in pkls:
14        models[pkl] = joblib.load(os.path.join(script_dir, f'{pkl}.pkl'))
15
16    return models
17
18 #def pre_score(df): # optional
19 #    print("pre score stage")
20 #    #print(df)
21 #    pass
22
23
24    def score(df, models):
25        _df = df.copy()
26
27
```

Скачать архив с версией

Score.py

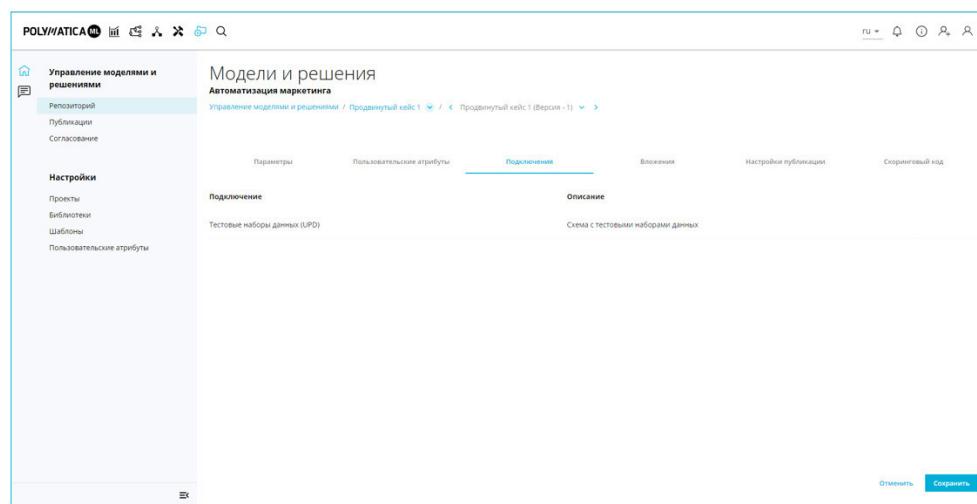
Размер: 256|Кб

Дата обновления: 23.01.2021 19:02:04

Отменить Сохранить

### Рисунок 197 Вкладка Скоринговый код

Вкладка **Подключения** актуальна только для моделей типа DM. В ней отображен список всех Подключений (к БД), которые используются в решении.



**Рисунок 198 Вкладка Подключения**

## 5.6. Публикация версии модели

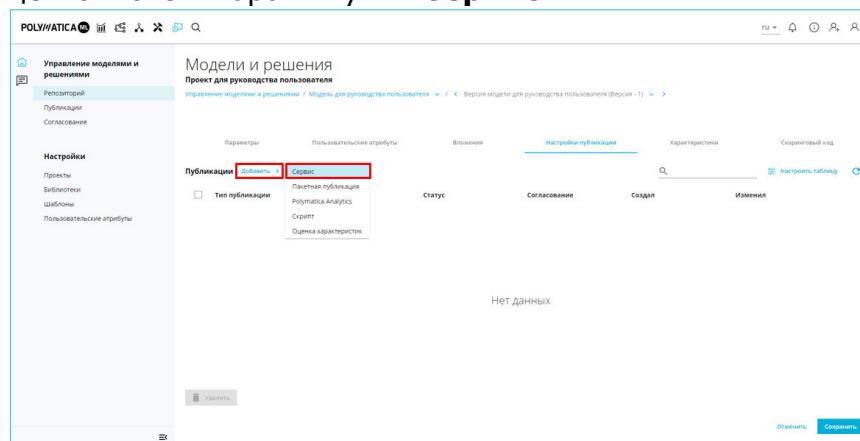
Модель может быть опубликована в виде:

- Сервиса.
- Пакетной публикации.
- Скрипта (только для типа модели MD).
- Оценить качество модели (Оценка характеристик) (только для типа модели MD).
- Polymatica Analytics (только для типа модели MD).

### 5.6.1. Сервис

Для публикации модели в виде Сервиса необходимо:

1. Выбрать наименование интересующей Версии модели в **Репозитории**.
2. Во вкладке **Настройки публикации** нажать на кнопку **Добавить** и в выпадающем списке выбрать пункт **Сервис**.



**Рисунок 199 Создание публикации типа Сервис**

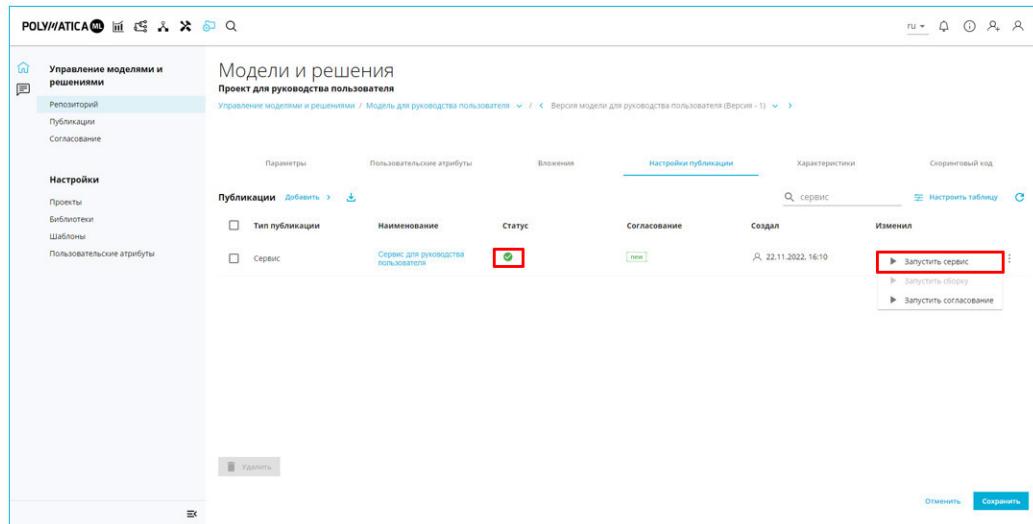
3. В открывшемся окне **Создание экземпляра** во вкладке **Параметры** задать следующие настройки:

- **Наименование.**
- **Уровень логирования** — какие события добавлять в лог (All, Debug, Info, Warn, Error, Fatal, Off).
- **Количество реплик.**
- Название эндпоинта (заглавными латинскими буквами без пробелов)
- **Максимальное количество ядер.**
- **Максимальное количество оперативной памяти** (в Мб).
- **Pre score** — функция предобработки (при необходимости).
- **Post score** — функция постобработки (при необходимости).



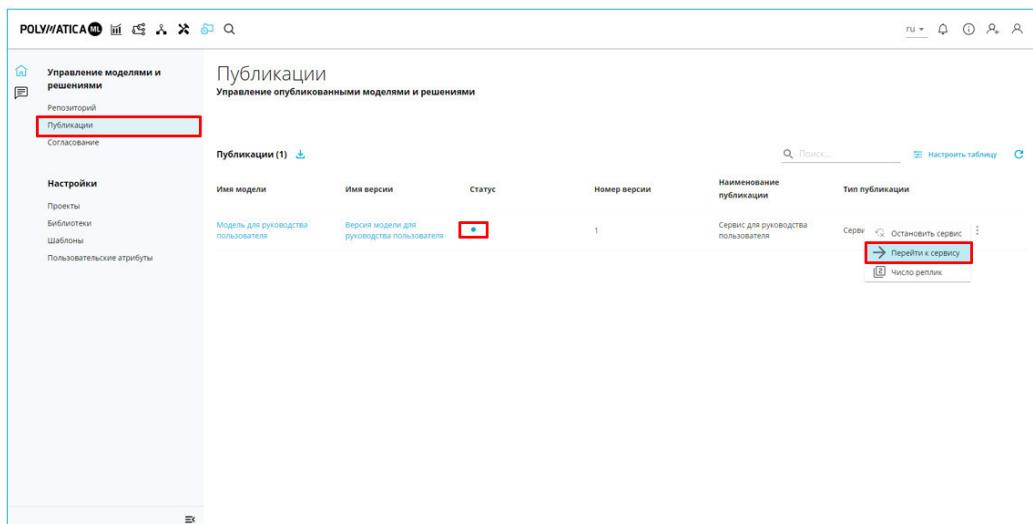
**Рисунок 200 Окно Создание экземпляра**

4. При необходимости изменить имена переменным (столбец **Имя переменной мэппинга**) во вкладке **Мэппинг**.
5. Сохранить изменения.
6. После успешной сборки модели (иконка в столбце Статус) выбрать меню и в выпадающем списке выбрать пункт **Запустить сервис**.



**Рисунок 201 Запуск сервиса**

7. Опубликованный сервис появится в разделе **Публикации** боковой панели компонента **Управление моделями и решениями (ММ)**. После изменения статуса публикации на Работает (иконка ●) выбрать меню ⏮ и в выпадающем списке выбрать пункт **Перейти к сервису**.



**Рисунок 202 Переход к сервису**

**Для снятия с публикации модели в виде Сервиса необходимо:**

- В разделе **Публикации** боковой панели компонента **Управление моделями и решениями (ММ)** выбрать меню рядом с созданной публикацией и в выпадающем списке нажать на **Остановить сервис**.

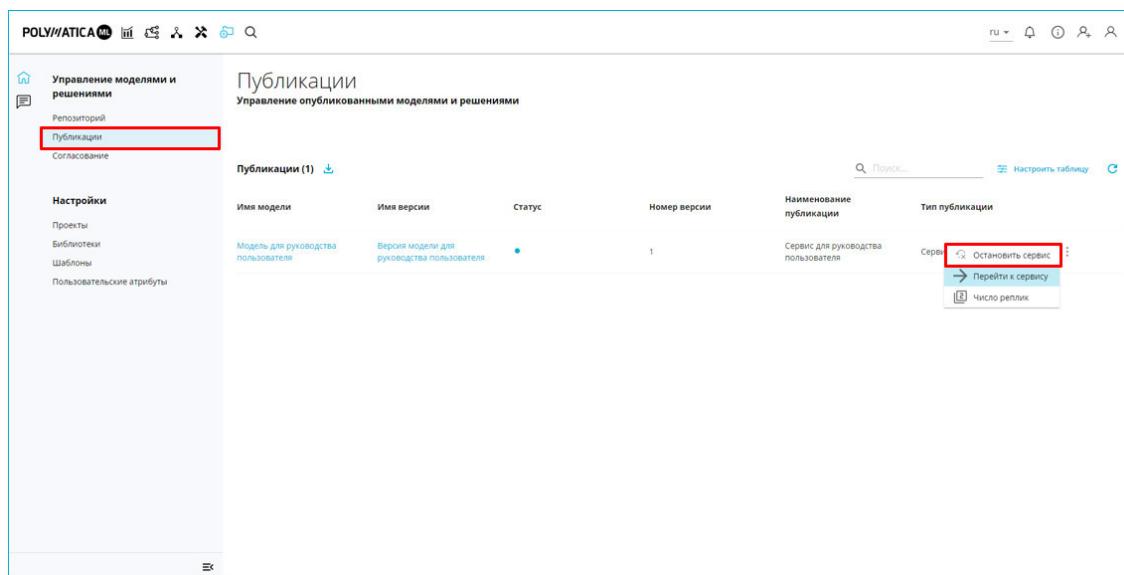


Рисунок 203 Снятие сервиса с публикации

## 5.6.2. Пакетная публикация

Для создания **пакетной публикации** модели необходимо:

1. Выбрать наименование интересующей Версии модели в **Репозитории**.
2. Во вкладке **Настройки публикации** нажать на кнопку **Добавить** и в выпадающем списке выбрать пункт **Пакетная публикация**.

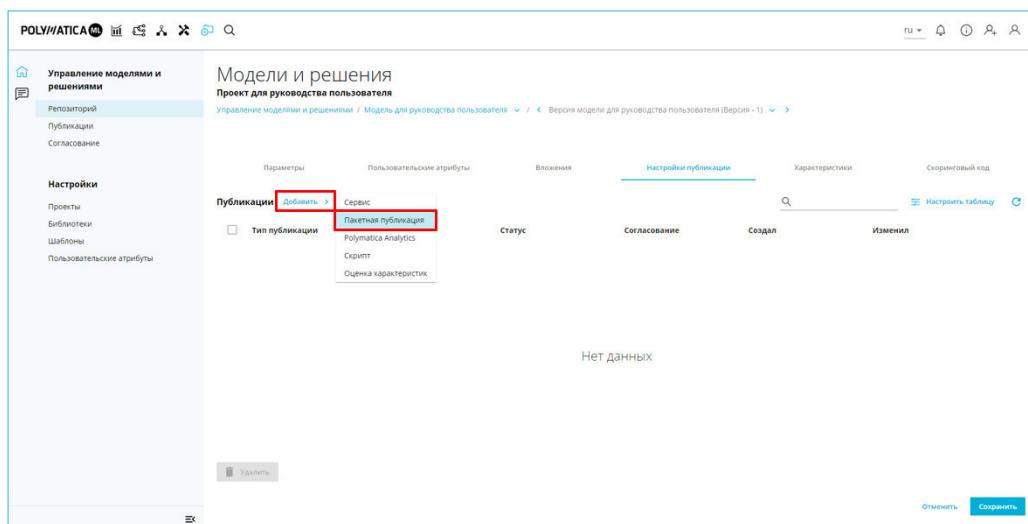
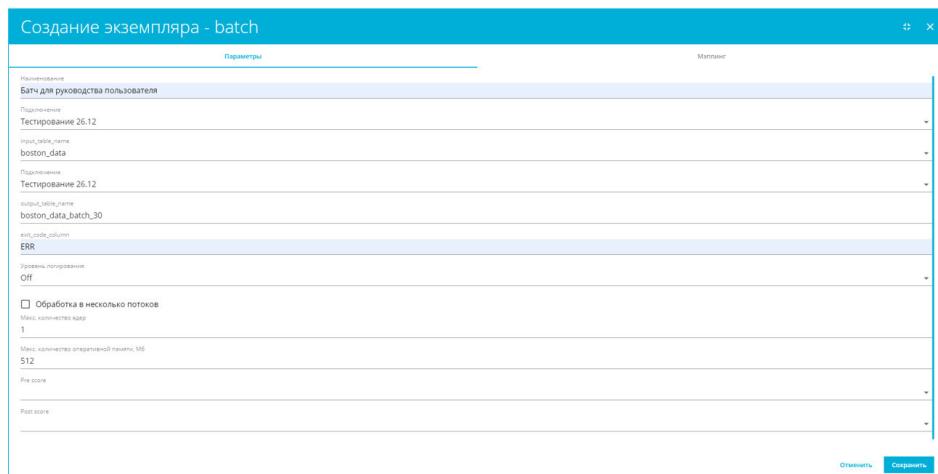


Рисунок 204 Создание пакетной публикации

3. В открывшемся окне **Создание экземпляра** во вкладке **Параметры** задать следующие настройки:
  - **Наименование.**
  - **Подключение** — указать подключение, в котором находится необходимый для скринга набор данных.
  - **input\_table\_name** — название входного набора данных.

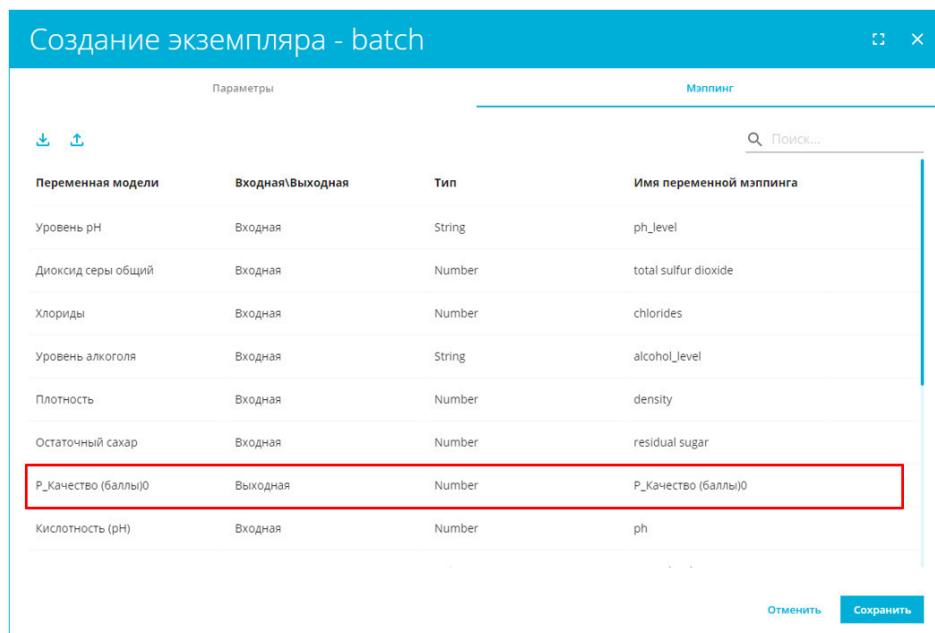
- **Подключение** — указать подключение, в которое необходимо положить результаты скоринга.
- **output\_table\_name** — название выходного набора данных.
- **exit\_code\_column** (при необходимости записи ошибок и информации об обработке в отдельную колонку).
- **Уровень логирования** — какие события добавлять в лог (All, Debug, Info, Warn, Error, Fatal, Off).
- Чекбокс **Обработка в несколько потоков** (при выборе чекбокса отобразится параметр **Количество реплик**, в котором можно указать необходимое число потоков)
- **Максимальное количество ядер.**
- **Максимальное количество оперативной памяти** (в Мб).
- **Pre score** — функция предобработки (при необходимости).
- **Post score** — функция постобработки (при необходимости).



**Рисунок 205 Окно Создание экземпляра**

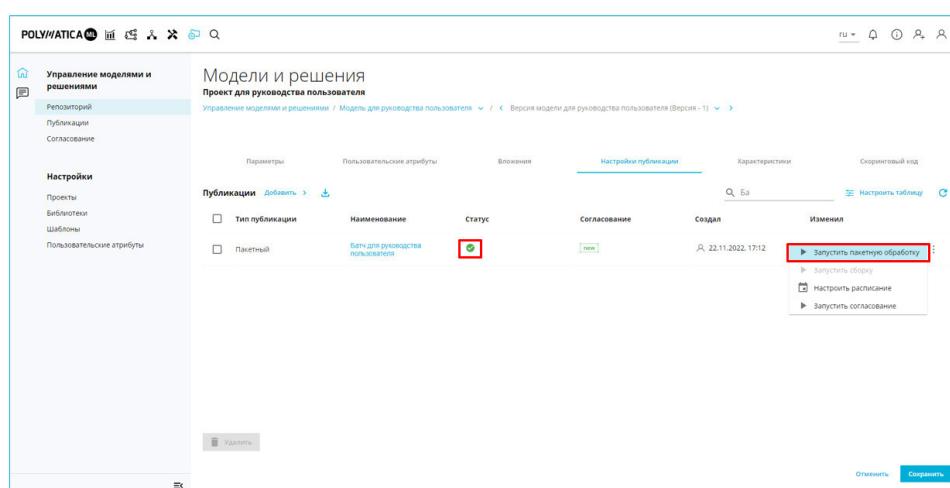
4. Задать имена переменным во вкладке **Мэппинг**.

Значения столбца **Имя переменной мэппинга** для входных переменных должно совпадать с наименованием соответствующего столбца в таблице. Для выходных переменных можно оставить определенное системой имя.



**Рисунок 206 Мэппинг переменных**

5. Сохранить изменения.
6. После успешной сборки модели (иконка в столбце Статус) выбрать меню и в выпадающем списке выбрать пункт **Запустить пакетную обработку**.



**Рисунок 207 Запуск пакетной обработки**

7. Публикация появится в разделе **Публикации** боковой панели компонента **Управление моделями и решениями (ММ)**. После завершения пакетной обработки публикация пропадет из раздела **Публикации**. С результатами расчетов Пользователь может ознакомиться в выходном файле (параметр **output\_table\_name**) в соответствующей БД (параметр **Подключение**). В выходной таблице помимо исходных столбцов появится столбец **exit\_code\_column** (значение которого в случае отсутствия ошибок должно быть равным 1), а также выходные столбцы модели (количество зависит от решаемой задачи).

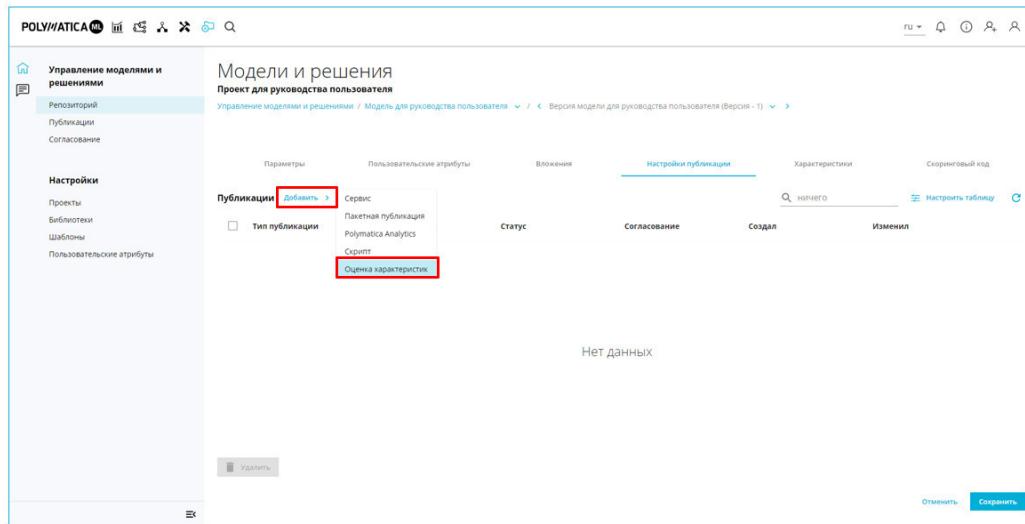
**Запуск пакетной обработки можно поставить на расписание.** Для этого необходимо:

- Рядом с интересующей публикацией в меню выбрать пункт **Настроить расписание**.
- Выбрать **Добавить Правило**.
- В открывшемся окне **Настройка расписания** задать:
  - Название, которое далее будет отображаться в системе.
  - Выбрать дату и время первого запуска.
  - Настроить повторение (период — ежедневно, еженедельно, ежемесячно и время повторного запуска).
  - Окончание расписания (Всегда запускать/Число выполнений (указать число запусков)/Дата окончания (указать последнюю дату запуска)).
  - Сохранить настройки.

### 5.6.3. Оценка характеристик

Для оценки качества модели необходимо:

1. Выбрать наименование интересующей Версии модели в **Репозитории**.
2. Во вкладке **Настройки публикации** нажать на кнопку **Добавить** и в выпадающем списке выбрать пункт **Оценка характеристик**.



**Рисунок 208 Создание публикации типа Оценка**

3. В открывшемся окне **Создание экземпляра** во вкладке **Параметры** задать следующие настройки:
  - **Наименование**.
  - **Подключение** — указать подключение, в котором находится необходимый для скринга набор данных.

- **input\_table\_name** — название входного набора данных.
- **Подключение** — указать подключение, в которое необходимо положить результаты скоринга.
- **output\_table\_name** — название выходного набора данных.
- **Уровень логирования** — какие события добавлять в лог (All, Debug, Info, Warn, Error, Fatal, Off).
- **Максимальное количество ядер.**
- **Максимальное количество оперативной памяти** (в Мб).
- **Pre score** — функция предобработки (при необходимости).
- **Post score** — функция постобработки (при необходимости).
- **Целевая переменная** — переменная в исходном наборе данных, относительно которой будет рассчитываться оценка.

**Задание параметра Целевая переменная**  
зависит от решаемой задачи.  
**Для задачи регрессии:** Для единственной  
выходной переменной необходимо выбрать **Тип  
метрики = interval**.

| + | Целевая переменная | Переменная Веса | Тип Метрики | Переменные Модели |
|---|--------------------|-----------------|-------------|-------------------|
| □ | medv               |                 | interval    | P_medv            |

**Рисунок 209 Пример задания Целевой переменной  
для задачи регрессии**

**Для задачи бинарной классификации:** Для двух переменных с префиксом P\_ необходимо выбрать **Тип метрики = binary\_prob** (т.к. они отражают вероятность принадлежности наблюдения к классу), для переменной с префиксом C\_ необходимо выбрать **Тип метрики = nominal\_class** (т.к. она отражает принадлежность к одному из классов).

| + | Целевая переменная | Переменные Веса | Тип Метрики  | Переменные Модели |
|---|--------------------|-----------------|--------------|-------------------|
| □ | C1_disease         |                 | nominal_0101 | P_C1_disease      |
| □ | C1_disease         |                 | binary_0101  | P_C1_disease_0    |

**Рисунок 210 Пример задания Целевой переменной  
для задачи бинарной классификации**

**Для задачи многоклассовой классификации:**  
Для всех переменных (количество совпадает с количеством классов) с префиксом P\_ необходимо выбрать **Тип метрики = nominal\_prob** (т.к. они

отражают вероятность принадлежности наблюдения к классу), для переменной с префиксом C\_ необходимо выбрать **Тип метрики = nominal\_class** (т.к. она отражает принадлежность к одному из классов).

| Целевая переменная | Переменная Веса | Тип Метрики   | Переменные Модели                                |
|--------------------|-----------------|---------------|--|
| C_sqm              |                 | nominal_class | C_sqm  |
| C_sqm              |                 | nominal_prob  | P_C_sqm_posprob P_C_sqm_posclass P_C_sqm_negprob |

**Рисунок 211 Пример задания Целевой переменной для задачи многоклассовой классификации**



**Рисунок 212 Окно Создание экземпляра**

4. Задать имена переменным во вкладке **Мэппинг**.

Значения столбца **Имя переменной мэппинга** для входных переменных должно совпадать с наименованием соответствующего столбца в таблице. Для выходных переменных можно оставить определенное системой имя.

The screenshot shows a configuration interface for a machine learning model. The top bar has tabs for 'Параметры' (Parameters) and 'Мэппинг' (Mapping). The 'Мэппинг' tab is active. A search bar at the top right contains the placeholder 'Поиск...'. Below the search bar is a table with columns: 'Переменная модели' (Model Variable), 'Входная\Выходная' (Input\Output), 'Тип' (Type), and 'Имя переменной мэппинга' (Mapping Variable Name). The table lists several variables: Уровень pH (Input, String, ph\_level), Диоксид серы общий (Input, Number, total sulfur dioxide), Хлориды (Input, Number, chlorides), Уровень алкоголя (Input, String, alcohol\_level), Плотность (Input, Number, density), Остаточный сахар (Input, Number, residual sugar), and R\_Качество (баллы)0 (Output, Number, R\_Kachество (баллы)0). The last row, which includes 'Кислотность (рН)' (Input, Number, ph), is not highlighted. At the bottom right are 'Отменить' (Cancel) and 'Сохранить' (Save) buttons.

**Рисунок 213 Задание мэппинга переменных**

5. Сохранить изменения.
6. После успешной сборки модели (иконка в столбце Статус) выбрать меню и в выпадающем списке выбрать пункт **Запустить оценку метрик**.
7. Публикация появится в разделе **Публикации** боковой панели компонента **Управление моделями и решениями (ММ)**. После завершения оценки публикация пропадет из раздела **Публикации**. С результатами расчетов Пользователь может ознакомиться в во вкладке **Характеристики модели**.

The screenshot shows the 'Models and solutions' section of the 'Management of models and solutions' component. The left sidebar has sections for 'Управление моделями и решениями' (Management of models and solutions), 'Репозиторий' (Repository), 'Публикации' (Publications), and 'Согласование' (Approval). The main area displays a table of model characteristics. The table has columns: 'Параметры' (Parameters), 'Пользовательские атрибуты' (User attributes), 'Вложения' (Attachments), 'Настройки публикации' (Publication settings), 'Характеристики' (Characteristics), and 'Скоринговый код' (Scoring code). The 'Характеристики' column is highlighted with a red border. The table rows represent different model runs: 'Обучение' (Training), 'Тест' (Test), 'Валидация' (Validation), and 'Скоринг' (Scoring). The 'Скоринг' row is also highlighted with a red border. At the bottom right are 'Отменить' (Cancel) and 'Сохранить' (Save) buttons.

**Рисунок 214 Результаты расчета оценки модели**

Запуск оценки характеристик модели можно поставить на расписание. Для этого необходимо:

- Рядом с интересующей публикацией в меню  выбрать пункт **Настроить расписание**.
- Выбрать **Добавить Правило**.
- В открывшемся окне **Настройка расписания** задать:
  - Название, которое далее будет отображаться в системе.
  - Выбрать дату и время первого запуска.
  - Настроить повторение (период — ежедневно, еженедельно, ежемесячно и время повторного запуска).
  - Окончание расписания (Всегда запускать/Число выполнений (указать число запусков)/Дата окончания (указать последнюю дату запуска)).
  - Сохранить настройки.

#### **5.6.4. Скрипт**

Для публикации модели в виде **Скрипта** необходимо:

1. Выбрать наименование интересующей Версии модели в **Репозитории**.
2. Во вкладке **Настройки публикации** нажать на кнопку **Добавить** и в выпадающем списке выбрать пункт **Скрипт**.
3. В открывшемся окне **Создание экземпляра** во вкладке **Параметры** задать **Наименование**.
4. Сохранить изменения.
5. После успешной сборки модели (иконка



в столбце Статус) выбрать меню



и в выпадающем списке выбрать пункт **Скачать**. В скачанном архиве можно найти все необходимые для запуска файлы, а также инструкцию по использованию (файл **README.md**).

#### **5.6.5. Публикация для скоринга мультисфер (**Polymatica Analytics**)**

Для публикации модели для скоринга мультисфер (**Polymatica Analytics**) необходимо:

1. Выбрать наименование интересующей Версии модели в **Репозитории**.

2. Во вкладке **Настройки публикации** нажать на кнопку **Добавить** и в выпадающем списке выбрать пункт **Polymatica Analytics**.
3. В открывшемся окне **Создание экземпляра** во вкладке **Параметры** задать следующие настройки:
  - **Наименование.**
  - **Уровень логирования** — какие события добавлять в лог (All, Debug, Info, Warn, Error, Fatal, Off).
  - **Количество реплик.**
  - Название эндпоинта (заглавными латинскими буквами без пробелов)
  - **Максимальное количество ядер.**
  - **Максимальное количество оперативной памяти** (в Мб).
  - **Pre score** — функция предобработки (при необходимости).
  - **Post score** — функция постобработки (при необходимости).
1. При необходимости поправить **Имя переменной мэппинга** во вкладке **Мэппинг**.
2. Сохранить изменения.
3. После успешной сборки модели (иконка в столбце Статус) выбрать меню и в выпадающем списке выбрать пункт **Запустить**.
4. Публикация появится в разделе **Публикации** боковой панели компонента **Управление моделями и решениями (ММ)**. После изменения статуса публикации на Работает (иконка ) выбрать меню и в выпадающем списке выбрать пункт **Перейти к скорингу мультисфер**.

## 5.7. Согласование версий модели и публикаций

### 5.7.1. Шаблон согласования

**Шаблон** представляет собой BPMN-процесс, в соответствии с которым будет выполняться согласование версии модели или публикации.

Шаблоны могут быть произвольными и базируются на ролевых моделях пользователей. Создаются шаблоны Администратором.

При выборе пункта **Шаблоны** боковой панели открывается одноименный раздел со списком Шаблонов согласования.

**Рисунок 215 Раздел Шаблоны согласования**

При нажатии на наименование шаблона открывается окно просмотра шаблона.

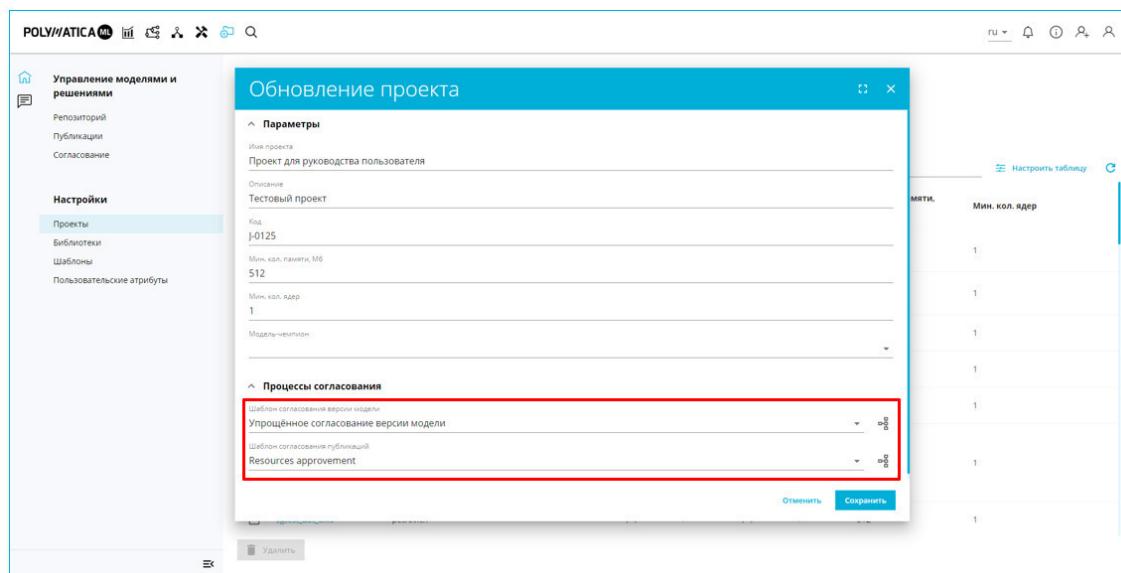
```

graph LR
    Start((Начало согласования)) --> Task[Согласование версии модели администратором]
    Task --> Decision{Согласовано?}
    Decision -- Да --> Approved([Версия модели согласована])
    Decision -- Нет --> NotApproved([Версия модели не согласована])
  
```

**Рисунок 216 Окно просмотра шаблона согласования**

### 5.7.2. Процесс согласования

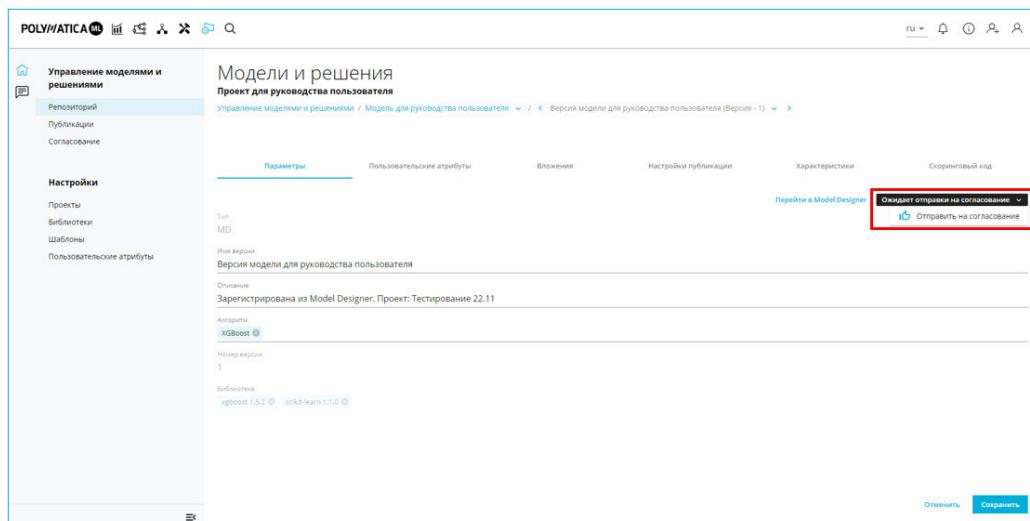
Процесс согласования происходит в соответствии с шаблонами, которые были указаны Пользователем в **Проекте ММ**.



**Рисунок 217 Задание шаблонов согласования в проекте ММ**

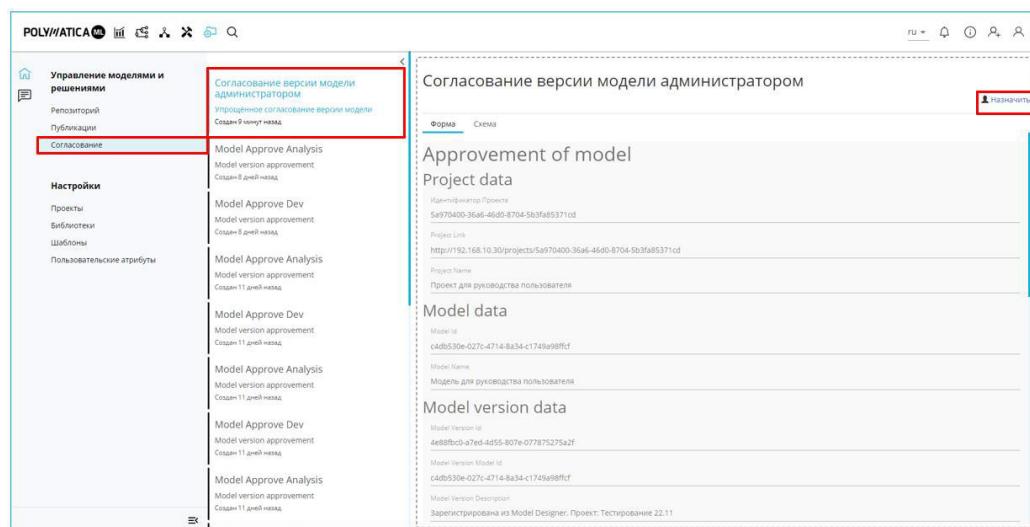
#### Для запуска согласования версии модели:

1. Выбрать наименование интересующей Версии модели в **Репозитории**.
2. Во вкладке **Параметры** нажать на кнопку **Ожидает отправки на согласование** и выпадающем списке выбрать пункт **Отправить на согласование**.

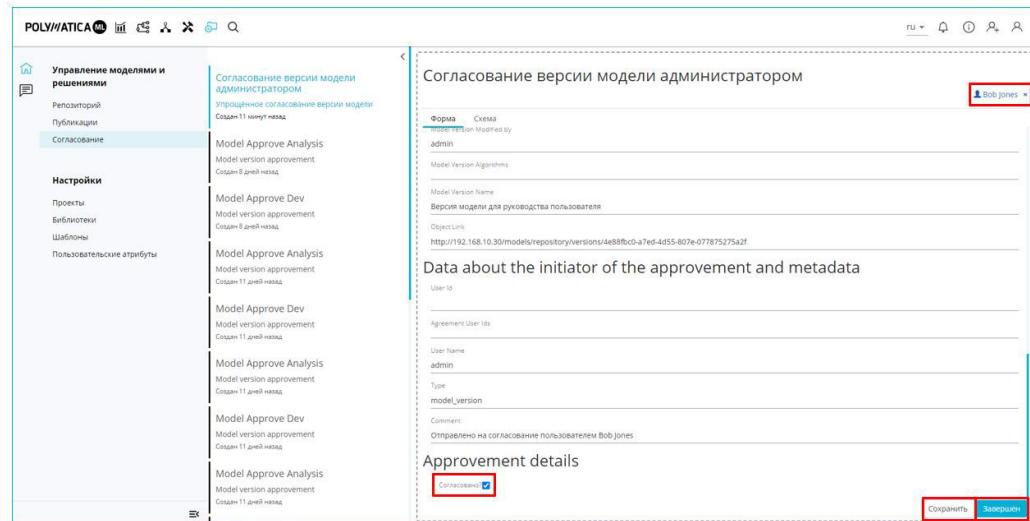


**Рисунок 218 Запуск процесса согласования**

3. Процесс согласования запущен. Статус версии модели изменился на **В процессе согласования**. В соответствии с шаблоном согласования всем заинтересованным пользователям придет сообщение о **Согласовании версии модели**. Для согласования (или отправки на доработку) Пользователю с соответствующими правами необходимо перейти в раздел **Согласование**, выбрать соответствующее сообщение, назначить себя в качестве Согласующего, проверить характеристики модели и согласовать (или отправить на доработку версию модели).

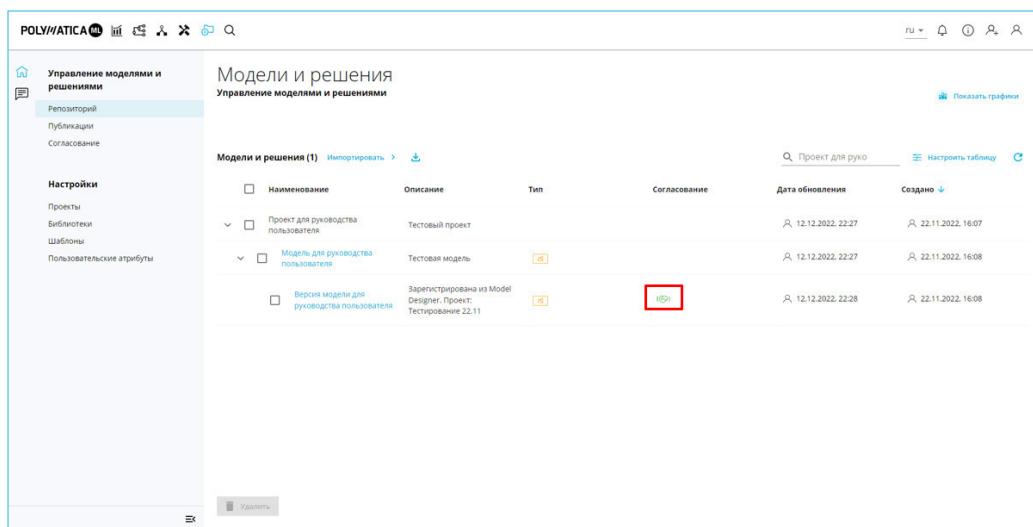


**Рисунок 219 Процесс согласования Версии модели**



**Рисунок 220 Процесс согласования Версии модели**

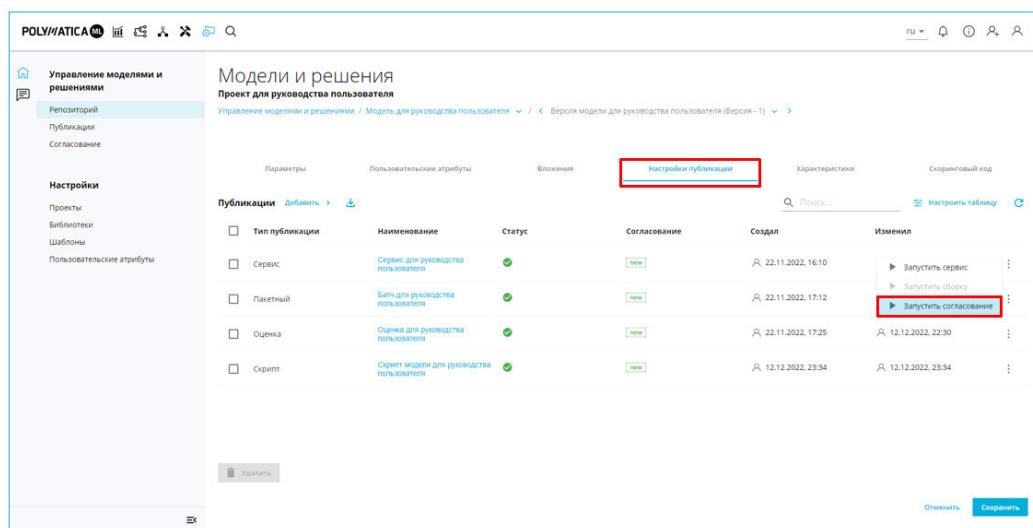
4. В зависимости от решения Согласующего в Репозитории в строке с Версией модели в столбце Согласование отобразится соответствующий статус (Согласовано/Не согласовано).



**Рисунок 221 Статус согласования Версии модели**

**Для запуска согласования публикации:**

1. Выбрать наименование интересующей Версии модели в **Репозитории**.
  2. Во вкладке **Настройки публикации** выбрать интересующую публикацию и в меню  выбрать пункт **Запустить согласование**.



## **Рисунок 222 Запуск процесса согласования публикации**

3. Процесс согласования запущен. Статус версии модели изменился на **В процессе согласования**. В соответствии с шаблоном согласования всем заинтересованным пользователям придет сообщение о **Согласование версии модели**. Для согласования (или отправки на доработку) Пользователю с соответствующими правами необходимо перейти в раздел **Согласование**, выбрать соответствующее сообщение, назначить себя в качестве Согласующего, проверить характеристики модели и согласовать (или отправить на доработку версию модели).

4. В зависимости от решения Согласующего во вкладке **Настройки публикации** в строке с Публикацией в столбце **Согласование** отобразится соответствующий статус (Согласовано/Не согласовано).

The screenshot shows a table of publications under the 'Настройки публикации' tab. One row is highlighted with a red box around the 'Согласование' column, which contains the status 'Не'. Other columns include 'Тип публикации', 'Назначение', 'Статус', 'Создан', and 'Изменил'.

| Тип публикации | Назначение                                 | Статус | Согласование | Создан            | Изменил           |
|----------------|--|--------|--------------|-------------------|-------------------|
| Сервис         | Сервис для руководства пользователя        | ✓      | Не           | 22.11.2022, 16:10 | 12.12.2022, 22:30 |
| Батч           | Батч для руководства пользователя          | ✓      | new          | 22.11.2022, 17:12 | 12.12.2022, 22:30 |
| Оценка         | Оценка для руководства пользователя        | ✓      | new          | 22.11.2022, 17:25 | 12.12.2022, 22:30 |
| Скрипт         | Скрипт модели для руководства пользователя | ✓      | new          | 12.12.2022, 23:34 | 12.12.2022, 23:34 |

**Рисунок 223 Статус согласования публикации**

## 5.8. Библиотеки

Библиотеки используются при сборке публикаций версии модели.

Для ознакомления со списком библиотек необходимо выбрать одноименный пункт боковой панели компонента **Управление моделями и решениями (ММ)**.

The screenshot shows a table of libraries under the 'Библиотеки' tab. One row is highlighted with a red box around the 'Библиотеки' item in the sidebar. Other columns include 'Наименование', 'Описание', 'Версия', and 'Ссылка'.

| Наименование | Описание            | Версия | Ссылка  |
|--------------|---------------------|--------|---|
| catboost     | catboost 1.0.4      | 1.0.4  | <a href="https://catboost.ai/en/docs/">https://catboost.ai/en/docs/</a>   |
| xgboost      | xgboost 1.5.2       | 1.5.2  | <a href="https://xgboost.readthedocs.io/en/stable/install...">https://xgboost.readthedocs.io/en/stable/install...</a>                             |
| scikit-learn | scikit-learn 1.0.2  | 1.0.2  | <a href="https://scikit-learn.org/stable/auto_examples/release_highligh...">https://scikit-learn.org/stable/auto_examples/release_highligh...</a> |
| scikit-learn | scikit-learn 1.0.1  | 1.0.1  | <a href="https://scikit-learn.org/stable/auto_examples/release_highligh...">https://scikit-learn.org/stable/auto_examples/release_highligh...</a> |
| scikit-learn | scikit-learn 1.1.0  | 1.1.0  | /link   |
| lightgbm     | lightgbm 3.3.2      | 3.3.2  | <a href="https://lightgbm.readthedocs.io/en/latest/install...Guide.html">https://lightgbm.readthedocs.io/en/latest/install...Guide.html</a>       |
| auto-sklearn | auto-sklearn 0.15.1 | 0.15.1 | /link   |
| skorch       | skorch 0.12.0       | 0.12.0 | /link   |
| torch        | torch 1.12.1        | 1.12.1 | /link   |

**Рисунок 224 Раздел Библиотеки бокового меню**

**Для создания Библиотеки необходимо:**

1. Выбрать кнопку **Добавить** в верхней части таблицы.
2. В открывшемся окне **Создание библиотеки** ввести следующие параметры:

- Наименование (в соответствии с реальным названием используемой библиотеки).
- Описание.
- Версию библиотеки.
- Ссылку.

3. Выбрать кнопку **Сохранить**.

**Для удаления Библиотеки необходимо:**

1. Выбрать чекбокс рядом с наименованием удаляемой библиотеки.
2. Выбрать кнопку **Удалить**.

**Для обновления параметров библиотеки необходимо:**

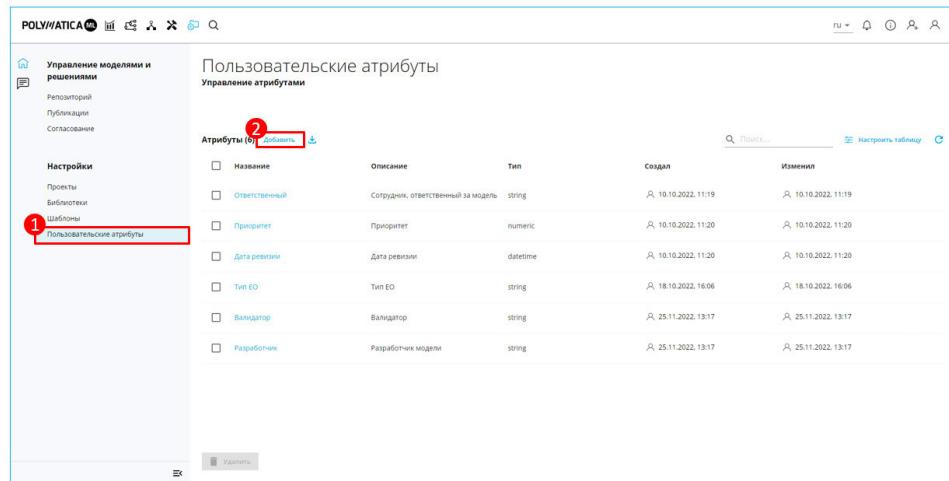
1. Нажать на наименование библиотеки.
2. В открывшемся окне **Просмотр библиотеки** ввести необходимые изменения.
3. Сохранить редактированные параметры, выбрав кнопку **Сохранить**.

## 5.9. Пользовательские атрибуты

Атрибуты предусмотрены для указания дополнительной информации о версии модели - инициалы автора, менеджера модели, дата разработки/важных изменений, приоритет версии и все то, что посчитает важным Пользователь.

**Для создания нового пользовательского атрибута необходимо:**

1. Выбрать одноименный пункт в боковой панели компонента **Управление моделями и решениями (ММ)**.
2. На открывшемся экране **Пользовательские атрибуты** выбрать кнопку **Добавить** в верхней части списка атрибутов.

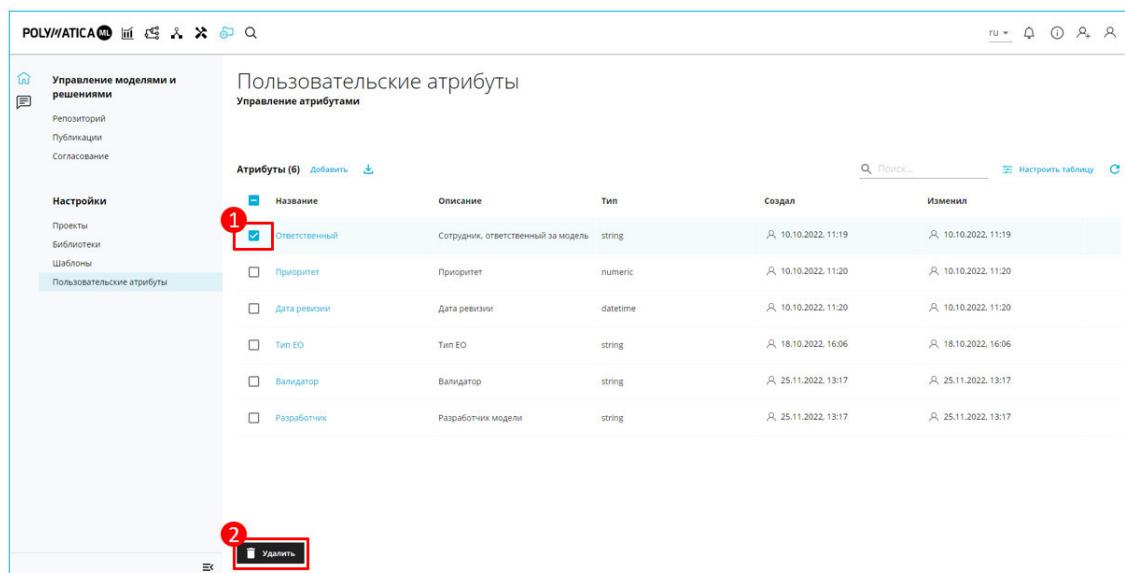


**Рисунок 225 Создание нового Пользовательского атрибута**

3. В открывшемся окне **Создание атрибута** задать название, описание и тип. Сохранить внесенные изменения
4. Созданный атрибут появится в списке экрана Пользовательские атрибуты и его можно будет указать в Модели в Репозитории и задать конкретное значение для Версий модели.

#### **Для удаления пользовательского атрибута необходимо:**

1. Выбрать чекбокс рядом с удаляемым атрибутом.
2. Нажать кнопку **Удалить** в нижней части экрана.



**Рисунок 226 Удаление пользовательского атрибута**

## 6. Глобальный поиск

### 6.1. Интерфейс экрана Глобального поиска

Экран Поиска открывается при выборе иконки в левом верхнем меню и представляет собой список доступных пользователю объектов системы.

The screenshot shows the 'Поиск' (Search) screen with the title 'Глобальный поиск'. It displays a table with 36 results, each row representing an object with columns for Name, Type, Version, Description, Creation Date, Author, and Last Modification Date. The table includes a search bar at the top right and a 'Настройка Таблицы' (Table Settings) button. The data rows are as follows:

| Наименование             | Тип       | Версия | Описание                                    | Создано              | Автор | Изменено             |
|--------------------------|-----------|--------|---|----------------------|-------|----------------------|
| Атамасов. Примеры        | MDProject | 1      | Атамасов. Примеры                           | 03.02.2022, 17:00:09 | admin | 03.02.2022, 17:00:09 |
| Таблица Менделеева       | MDProject | 1      | Модель плотности                            | 13.10.2021, 17:55:54 | admin | 31.01.2022, 09:22:06 |
| Test clustering          | MDProject | 1      |   | 07.02.2022, 09:28:31 | admin | 07.02.2022, 09:28:31 |
| POLYANALYTICA            | MDProject | 1      | Проект модели для мультифермы               | 11.11.2021, 12:55:30 | admin | 11.11.2021, 12:55:30 |
| Минтрез средняя зарплата | MDProject | 1      |   | 19.11.2021, 19:55:07 | admin | 19.11.2021, 19:55:07 |
| Узлы                     | MDProject | 1      | Для отражок группировки пользователи и т.д. | 24.01.2022, 13:25:52 | admin | 24.01.2022, 13:25:52 |
| Закупки                  | MDProject | 1      | Закупки                                     | 24.11.2021, 13:59:50 | admin | 24.11.2021, 13:59:50 |
| Кредитный скрипте для ЦБ | MDProject | 1      | Построение модели кредитного скрипта        | 02.12.2021, 10:14:50 | admin | 08.12.2021, 16:07:56 |
| Oracle mendeleev         | MDProject | 1      | Oracle mendeleev                            | 22.12.2021, 16:05:08 | admin | 22.12.2021, 16:05:08 |

Рисунок 227 Экран глобального поиска

При выборе наименования объекта откроется соответствующий раздел Модуля с выбранным объектом для дальнейшей работы. Помимо наименования самих объектов в таблице также отображены:

- Тип объекта – DataSet, MDProject, MMProject.
- Версия
- Описание объекта
- Дата создания и изменения объекта
- Автор
- GUID

Таблица с доступными пользователю объектами системы имеет гибкие настройки отображения. Пользователь может:

- изменить ширину любого столбца (для этого необходимо перетащить границу его заголовка до нужной ширины);
- сортировать таблицу (для этого необходимо выбрать иконку рядом с заголовком сортируемого столбца);
- скрывать/отображать столбцы и изменять их порядок в окне **Вид таблицы** (для открытия окна необходимо выбрать иконку в правом верхнем углу таблицы; при выборе иконки столбец скроется, при наведении на иконку активируется возможность перемещения столбца);
- сбросить внесенные изменения также в окне **Вид таблицы** (для этого выбрать кнопку «**Сбросить**»).

Для быстрого поиска объекта в таблице предусмотрено поле  Пойск... в правой верхней части таблицы.

Объекты таблицы можно выгрузить в формате Excel. Для этого нужно выбрать

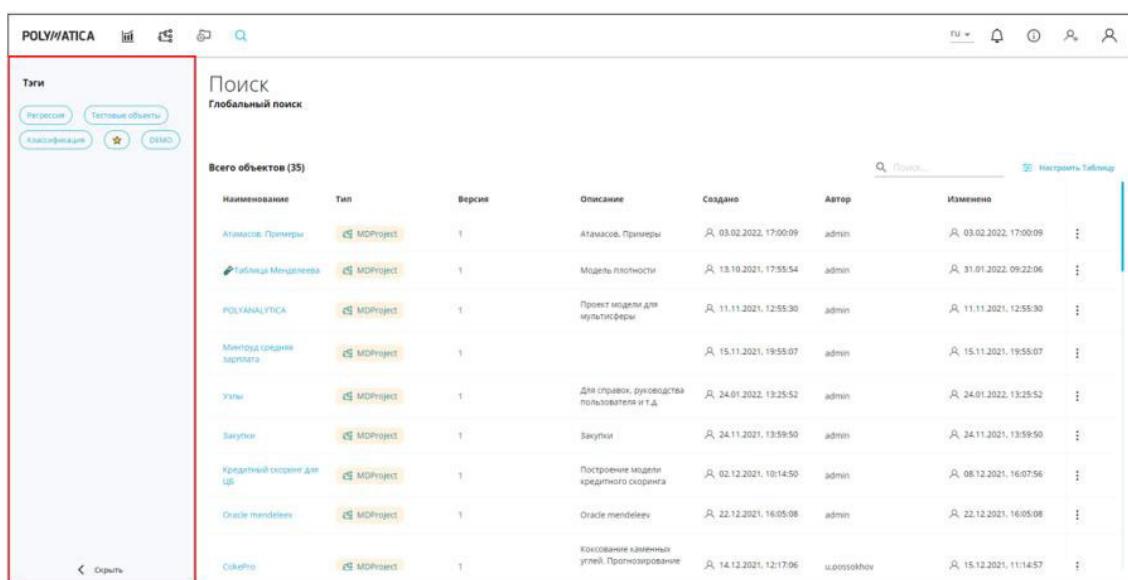
иконку **Экспорта в excel** .

Помимо этого, реализована работа с тэгами для быстрого поиска объектов.

## 6.2. Работа с тэгами

Скрытая левая боковая панель содержит тэги, при помощи которых можно быстро найти интересующий объект.

Для раскрытия панели необходимо выбрать  в нижней части интерфейса.



| Наименование              | Тип       | Версия | Описание                                     | Создано              | Автор     | Изменено             |
|---------------------------|-----------|--------|--|----------------------|-----------|----------------------|
| Атамасов. Примеры         | MDProject | 1      | Атамасов. Примеры                            | 03.02.2022, 17:00:09 | admin     | 03.02.2022, 17:00:09 |
| Таблица Менделеева        | MDProject | 1      | Модель плотности                             | 13.10.2021, 17:55:54 | admin     | 31.01.2022, 09:22:06 |
| POLYANALYTICA             | MDProject | 1      | Проект модели для мультисферы                | 11.11.2021, 12:55:30 | admin     | 11.11.2021, 12:55:30 |
| Минтрод греющие зарядки   | MDProject | 1      |  | 15.11.2021, 19:55:07 | admin     | 15.11.2021, 19:55:07 |
| Улицы                     | MDProject | 1      | Для справок, руководства пользователя и т.д. | 24.01.2022, 19:25:52 | admin     | 24.01.2022, 13:25:52 |
| Закупки                   | MDProject | 1      | Закупки                                      | 24.11.2021, 13:59:50 | admin     | 24.11.2021, 13:59:50 |
| Кредитный скрининг для ЦБ | MDProject | 1      | Построение модели кредитного скрининга       | 02.12.2021, 10:14:50 | admin     | 08.12.2021, 16:07:56 |
| Oracle mendeleev          | MDProject | 1      | Oracle mendeleev                             | 22.12.2021, 16:05:08 | admin     | 22.12.2021, 16:05:08 |
| CokePro                   | MDProject | 1      | Кассование канимных услуг. Прогнозирование   | 14.12.2021, 12:17:06 | ш.посохов | 15.12.2021, 11:14:57 |

Рисунок 228 Боковая панель со списком тэгов

### 6.2.1. Добавление тэга

Для задания объекту тэга необходимо:

- В строке с интересующим объектом выбрать меню в виде трех вертикальных точек  и в раскрывшейся панели выбрать пункт **Работа с Тэгами**.
- В открывшемся окне **Работа с Тэгами** в строку ввода ввести тэг и выбрать кнопку **Добавить**.

### 6.2.2. Фильтрация при помощи тэгов

Для фильтрации объектов при помощи тэгов необходимо:

- Раскрыть левую боковую панель.
- В списке с тэгами найти интересующий и выбрать его (можно выбрать несколько тэгов у одного объекта).

- Таблица с доступными объектами отсортируется в соответствии с выбранными тэгами.

## 6.3. Связанные объекты

Помимо прочего на экране глобального поиска можно узнать связь между объектами. Для этого необходимо:

- В строке с интересующим объектом выбрать меню  и в раскрывшейся панели выбрать пункт **Связанные Объекты**.
- В открывшемся окне содержится список со связанными объектами. При выборе наименования объекта из данного списка откроется соответствующий раздел Модуля с выбранным объектом для дальнейшей работы.