

# **BAUM AI PREDICT**

**Руководство администратора**

Версия 1.0.0

Москва

2023

## СОДЕРЖАНИЕ:

<b>1. Общие положения</b>	<b>4</b>
1.1. Основные понятия и определения	4
<b>2. Введение</b>	<b>9</b>
2.1. Область применения	9
<b>3. Регистрация и авторизация пользователя в Системе</b>	<b>10</b>
<b>4. Личный кабинет пользователя</b>	<b>11</b>
<b>5. Интерфейс Платформы</b>	<b>13</b>
5.1. Меню интерфейса	13
5.2. Окно построения модели ИИ	15
5.3. Настройка внешнего вида интерфейса	17
5.4. Встроенные функции	18
<b>6. Загрузка данных в систему</b>	<b>21</b>
6.1. Создание новой папки	22
6.2. Загрузка файлов	22
6.3. Предпросмотр данных	24
6.4. Взаимодействие с данными	26
<b>7. Создание модели ИИ</b>	<b>26</b>
7.1. Создание новой и открытие сохраненной рабочей области	26
7.1.1. Создание новой рабочей области	26
7.1.2. Открытие сохраненной рабочей области	28
7.2. Построение блок-схемы	28
7.2.1. Блок «Запуск»	28
7.2.2. Блок «Источник данных»	29
7.2.3. Блок «Процесс»	33
7.3. Запуск блок-схемы на рабочей области	34
<b>8. Сохранение модели ИИ</b>	<b>35</b>
<b>9. Графическое представление информации на рабочей области</b>	<b>37</b>
9.1. Графики	37
9.2. Таблицы	39
9.4. Описание модели	40
<b>10. Работа с Дашбордами. Раздел «Визуализация».</b>	<b>41</b>
10.1. Таблица	43
<b>11. Создание отчета с результатами анализа данных</b>	<b>44</b>
<b>12. Конвейер приложений</b>	<b>45</b>
<b>13. Работа с проектом</b>	<b>47</b>
13.1. Создание нового проекта	47
13.2. Редактирование проекта	48
13.3. Наполнение проекта	49
<b>14. Настройка подключения к источникам данных</b>	<b>53</b>
14.1. Типы коннекторов	54
14.2. Порядок работы с коннекторами	55

14.2.1	Создание коннектора	56
14.2.2	Запуск коннектора	57
14.2.3	Подключение к коннектору на дашборде	57
14.3	Настройка подключения на примере ClickHouse	57
<b>15.</b>	<b>Примеры рабочих областей</b>	<b>62</b>
<b>15.1</b>	<b>Вводная часть</b>	<b>62</b>
15.1.1.	Описание проблемы	62
15.1.2.	Исходные данные	62
15.1.2.1.	Формат данных.	62
	Формат входных данных алгоритма - поток данных с СХД по SNMP 123 .	62
15.1.2.2.	Пример данных из СХД	62
15.1.3.	Постановка задачи	63
<b>15.2</b>	<b>Практическая часть</b>	<b>63</b>
15.1.1.	Функциональные требования	63
15.1.2.	Проект	63
15.1.2.1.	Рабочие области	64
15.1.2.1.1.	Пайплайн для прогнозирования загрузки	64
15.1.2.1.2.	Визуализация результатов	67
15.1.2.2.	Дашборды	68
15.1.2.3.	Файлы	69
15.1.2.4.	Коннекторы	69
15.1.2.5.1.	Источник данных	69
15.1.2.5.2.	ETL	70
15.1.2.5.3.	Коннектор	71
15.1.3.	Результаты	72
15.1.4.	Сохраняемые объекты	72
<b>16.</b>	<b>Администрирование Платформы</b>	<b>73</b>
16.1	Пользователи и группы	73
16.2	Настройка отправки уведомлений	76
<b>17.</b>	<b>Дополнительные возможности Платформы</b>	<b>79</b>
17.1.	Обращение в службу поддержки	79
17.2.	История изменений	79
<b>18.</b>	<b>Приложения</b>	<b>80</b>
	Приложение 1. Автоматизированные функции	80
<b>19.</b>	<b>Лист изменений</b>	<b>141</b>

## 1. Общие положения

### - 1.1. Основные понятия и определения

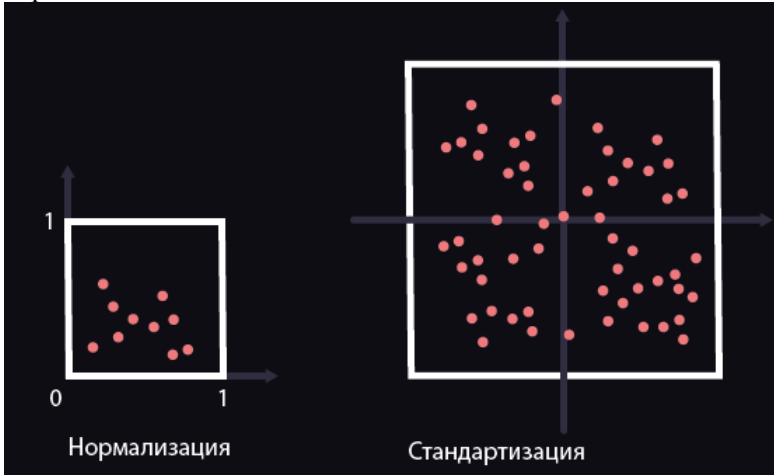
Таблица 1 – Основные термины, используемые в документе, и их определения

Термин	Определение
Apache Kafka	Брокер сообщений, реализующий паттерн Producer-Consumer. Данные из одного и того же топика могут считываться множеством консьюмер-групп одновременно.
Apache Spark	Фреймворк с открытым исходным кодом для реализации распределенной обработки неструктурированных и слабоструктурированных данных.
BPMN	С англ. Business Process Model and Notation. Нотация определяющая способ визуализации процессов в виде диаграмм с определенным набором блоков и взаимосвязей.
Временной ряд	Совокупность наблюдений, собранных за определенный временной интервал. Этот тип данных используется для поиска долгосрочного тренда, прогнозирования будущего и прочих видов анализа. Анализ временных рядов позволяет обнаруживать тенденции и закономерности в исследуемых процессах, строить прогнозы и предсказывать будущие изменения в бизнесе, на производстве, и в других областях.
Выборка	Случайное подмножество генеральной совокупности.
Датасет	С англ. <i>Data set</i> , набор данных. Коллекция из логических записей, хранящихся в виде <i>кортежа</i> . Набор данных можно сравнить с файлом, но в отличие от файла он является одновременно и каталогом, и файлом файловой системы, и не может содержать в себе другие наборы. Файловая система ориентирована на хранение записей, которые являются неделимыми единицами хранения. Множества записей объединяются в группы, которые и называются наборами данных. Записи в наборах данных используются приложениями, например, как входные данные. Так, записями набора данных могут быть как текстовые данные, так и изображения. К набору данных можно обратиться, указав точное место его хранения, или, если ранее для набора было зазервировано имя в файловой системе, по имени (второй вариант не реализован). Также <i>датасетом</i> называются данные, которые пользователь загрузил с локального компьютера, а Система при загрузке выполнила их предварительную обработку. Например, в систему загружается временной ряд в формате csv, и при загрузке он преобразуется в <i>структурированные данные</i> или по-другому датасет. Такой датасет пригоден для использования в моделях машинного обучения. В качестве <i>датасета</i> в настоящем документе также выступает набор изображений или файлов, которые используются для решения одной задачи.

<p>Дерево решений</p>	<p>Инструмент прогнозного моделирования. Строится с помощью алгоритмического подхода, который разделяет набор данных на основе различных условий. Относится к классу обучения с учителем, используется для задач классификации и регрессии. Цель в том, чтобы создать модель, которая предсказывает значение целевой переменной, изучая правила принятия решений, выведенные из характеристик данных.</p> <p>Алгоритм:</p> <ol style="list-style-type: none"> <li>1. Выбрать лучший атрибут, который разделяет наблюдения на группы.</li> <li>2. Задать соответствующий вопрос.</li> <li>3. Следовать по путям ответов.</li> <li>4. Вернуться к шагу 1.</li> </ol>
<p>Дисперсия</p>	<p>Мера удаленности того или иного значения выборки от среднего значения.</p> <p>Рассчитывается по формуле:</p> $\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N}$ <p>Стандартное отклонение рассчитывается как квадратный корень из полученной цифры.</p>
<p>Категориальная переменная (качественные данные)</p>	<p>Это данные с ограниченным числом уникальных значений или <i>категорий</i> (например, пол, страна проживания, номер группы, категория товаров, и т.п.). Категориальные поля могут быть как текстовыми, так и числовыми, в которых категории закодированы <i>числовыми кодами</i> (например, 0=женский, а 1=мужской). Номинальные поля, порядковые поля и флаги являются категориальными полями.</p> <p>-<i>Набор</i> (номинальная переменная). Поле, значения которого представляют категории без естественного упорядочивания (например, подразделение компании, в котором работает сотрудник).</p> <p>-<i>Упорядоченный набор</i> (порядковая переменная). Поле, значения которого представляют категории с некоторым естественным для них упорядочением (например, оценки, представляющие степень удовлетворенности или уверенности, или баллы, оценивающие предпочтение).</p> <p>-<i>Флаг</i>. Поле или переменная с двумя отдельными значениями, например Да и Нет.</p>
<p>Классификация</p>	<p>Задача машинного обучения, которая ставит своей целью назначить метку класса наблюдениям из предметной области.</p> <p>Основные типы классификации:</p> <ul style="list-style-type: none"> <li>• бинарная классификация;</li> <li>• мультиклассовая классификация;</li> <li>• классификация по нескольким меткам;</li> <li>• несбалансированная классификация.</li> </ul>
<p>Кластеризация</p>	<p>Техника обучения без учителя, которая включает в себя группирование или кластеризацию точек данных. Чаще всего она</p>

	<p>используется для сегментации потребителей, выявления мошенничества и классификации документов.</p> <p>Кластеризация (или кластерный анализ) – это задача разбиения множества объектов на группы, называемые кластерами.</p> <p>Главное отличие от классификации состоит в том, что перечень групп четко не задан и определяется в процессе работы алгоритма.</p>
<p>Машинное обучение (ML – Machine learning)</p>	<p>Тренировка математической модели на исторических данных для того, чтобы прогнозировать какое-то событие или явление на новых данных. То есть попытка заставить алгоритмы программ совершать действия на основе предыдущего опыта, а не только на основе имеющихся данных.</p> <p>Для обучения нужны исторические данные (обучающая выборка) и значение целевой переменной (то, что прогнозируем), которое соответствует заданным историческим данным. Модель наблюдает и находит зависимости между данными и целевой переменной. Эти зависимости используются моделью для нового набора данных, чтобы прогнозировать целевую переменную, которая неизвестна.</p> <p>Машинное обучение включает в себя целый набор методов и алгоритмов, которые могут предсказать какой-то результат по входным данным.</p> <p>Алгоритмов машинного обучения большое множество: одни эффективны для решения одного типа задач, вторые – для другого.</p> <p><b>Суть технологии машинного обучения</b></p> <p>Говоря в общем, машинное обучение – это обучение компьютерной программы или алгоритма постепенному улучшению исполнения поставленной задачи.</p> <p>Машинное обучение обозначает множество математических, статистических и вычислительных методов для разработки алгоритмов, способных решить задачу не прямым способом, а на основе поиска закономерностей в разнообразных входных данных.</p> <p>Решение вычисляется не по четкой формуле, а по установленной зависимости результатов от конкретного набора признаков и их значений. Например, если каждый день в течение недели земля покрыта снегом и температура воздуха существенно ниже нуля, то вероятнее всего, наступила зима. Поэтому машинное обучение применяется для диагностики, прогнозирования, распознавания и принятия решений в различных прикладных сферах: от медицины до банковской деятельности.</p>
<p>Мониторинг</p>	<p>Процесс наблюдения и регистрации данных о каком-либо объекте на неразрывно примыкающих друг к другу интервалах времени, в течение которых значения данных существенно не изменяются .</p>
<p>Мониторинг состояния</p>	<p>Наблюдение за состоянием объекта для определения и предсказания момента перехода в предельное состояние.</p> <p>Результат мониторинга состояния объекта представляет собой совокупность диагнозов составляющих его субъектов, получаемых на неразрывно примыкающих друг к другу интервалах времени, в течение которых состояние объекта существенно не изменяется. Принципиальное отличие от</p>

	мониторинга параметров является наличие интерпретатора измеренных параметров в терминах состояния – экспертной системы поддержки принятия решений о состоянии объекта и дальнейшем управлении.
Наблюдение (строка, запись, точка, сущность)	<p>Ценные данные, собираемые во время исследования или эксперимента. Вместе с масштабом анализа определяет совокупность.</p> <p><i>Эмпирические исследования</i> – практические эксперименты с результатами на основе реального опыта, а не теории или убеждений. Основополагающим принципом <i>Науки о данных</i> является приоритет наблюдения над предположением.</p> <p>Типы наблюдений:</p> <ul style="list-style-type: none"> <li>● <i>Числовой</i>: целые (integer), вещественные (real number), числа с плавающей точкой (float).</li> <li>● <i>Булевый</i> (boolean) – принимает значения 1/0 (да/нет).</li> <li>● <i>Категориальный</i>. Например, жанры кино: комедия, ужасы, мелодрама.</li> <li>● <i>Текстовый</i>.</li> <li>● <i>Вектор</i>.</li> </ul>
Нейронная сеть (или Искусственная нейронная сеть)	<p>Представляет собой <i>математическую модель</i>, а также её программное или аппаратное воплощение, построенную по принципу организации и функционирования биологических нейронных сетей – сетей нервных клеток живого организма. Нейронные сети решают задачу: по точкам находят функцию. Происходит это путем <i>минимизации ошибки</i> – сводится к минимуму «расстояние» между значениями, предсказываемыми нейронной сетью, и значениями, которые наблюдаются. Под <i>архитектурой нейронной сети</i> понимается ее устройство – последовательность нейронов и связей между ними.</p>
Нормализация	Техника преобразования значений признака, масштабирующая значения таким образом, что они расположены в диапазоне от 0 до 1.
Обучение с учителем	<p>Контролируемое обучение – метод машинного обучения, при котором модель обучается на размеченных данных. Например, исследовав опухоли, установив их размер, плотность и другие метрики, мы передаем эти данные модели с обязательной пометкой, какое наблюдение к какому строению (доброкачественному или злокачественному) относится. Алгоритмы контролируемого обучения подразделяются на следующие модели: классификация, регрессия.</p>
Пайплайн	Последовательность стадий, внутри которых расположены задачи. Расположены они таким образом, что выход каждого элемента является входом следующего.
Признак	<p>Объективная характеристика, характерная черта или свойство, которое может быть определено или измерено.</p> <p>В статистике независимые переменные X используются для предсказания зависимого признака Y</p>
Регрессия (в математической статистике)	<p>Математическое выражение, отражающее связь между зависимой переменной y и независимыми переменными x. Алгоритмы регрессии используются для контролируемого обучения моделей искусственного интеллекта. Модели обучают прогнозировать числовые значения целевых переменных.</p>

<p>Стандартизация</p>	<p>Техника преобразования значений признака, адаптирующая признаки с разными диапазонами значений к моделям машинного обучения, использующих дистанцию для прогнозирования. Это разновидность нормализации с использованием стандартизированной оценки преобразует значения так, что из каждого наблюдения каждого признака вычитается среднее значение и результат делится на стандартное отклонение этого признака:</p> $x_{i, \text{станд.}} = \frac{x_i - \mu}{\sigma}$ <p><math>x_{i, \text{станд.}}</math></p> <p><math>x_i</math> – исходный элемент, <math>\mu</math> – среднее арифметическое, <math>\sigma</math> – стандартное отклонение.</p> <p>Преобразование необходимо, поскольку признаки датасета могут иметь большие различия между своими диапазонами, и для моделей машинного обучения это спровоцирует искаженное восприятие данных.</p> <p>Стандартизация, в отличие от нормализации, не имеет ограничивающего диапазона:</p> 
<p>Стандартизированная оценка</p>	<p>Метрика, характеризующая удаленность наблюдения от среднего значения совокупности данных.</p>
<p>Стандартное отклонение</p>	<p>Мера разброса в наборе числовых данных. Показывает, насколько далеко от среднего арифметического находятся точки данных. Чем меньше стандартное отклонение, тем более сгруппированы данные вокруг центра (среднего). Чем отклонение больше, тем больше разброс значений.</p>
<p>Тренировочные (обучающие) данные</p>	<p>Часть датасета, обучающая основа модели машинного обучения. Является одной из составляющих разделенного набора данных наряду с <i>тестовыми</i> и <i>валидационными</i> данными.</p>
<p>Валидационные (тестовые) данные</p>	<p>Часть датасета, основа для проверки работоспособности модели машинного обучения.</p>
<p>Числовая переменная</p>	<p>Numeric variable - переменная, выраженная различными видами чисел.</p>



Целевая (зависимая) переменная	Target variable – признак датасета, который предстоит предсказывать модели машинного обучения. Зависимой ее называют, поскольку в ходе разведочного анализа данных выявляется корреляция между одной или несколькими переменными-предикторами и рассматриваемым целевым признаком.
--------------------------------	--

## 2. Введение

### - 2.1. Область применения

.ПО предиктивного анализа систем хранения данных (СХД) **BAUM AI PREDICT** представляет собой удобный инструмент для системного аналитика или администратора СХД для решения задачи прогнозирования загрузки СХД и превентивного ремонта составных частей СХД в режиме визуального программирования (моделирования).

ПО предназначено для решения задачи прогнозирования загрузки СХД и превентивного ремонта составных частей СХД для снижения рисков недоступности и отказа СХД.

### 2.2. Уровни подготовки пользователей

В Системе предполагается наличие ролей пользователей – *аналитик*, моделирующий с помощью конструктора искусственного интеллекта работу бизнес-процесса, и *оператор*, обладающий возможностью только просматривать *результаты запуска конструктора*. Платформа позволяет настраивать уровни доступа к разделам системы для разных пользователей, подробнее об этом в разделе [Администрирование](#).

### 3. Регистрация и авторизация пользователя в Системе

Для начала работы с Системой необходимо пройти процедуру **регистрации**, для этого:

1. Перейдите на страницу регистрации:

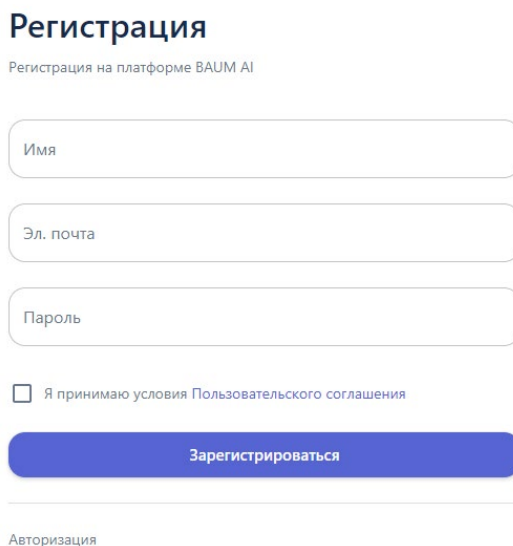


Рисунок 3.1 – Страница регистрации пользователя в Системе

По ссылке «Авторизация» доступен переход на страницу авторизации.

2. Заполните следующие поля:
  - Имя пользователя;
  - Адрес электронной почты;
  - Пароль для входа в Систему.
3. Ознакомьтесь с условиями пользовательского соглашения. Регистрация возможна только при условии их принятия.
4. Нажмите кнопку «Зарегистрироваться».

После этого шага выполняется регистрация пользователя в Системе – в БД **MongoDB** создается новая запись с уникальным идентификатором пользователя.

5. Вы можете сохранить связку логин и пароль для автоматического заполнения при последующей авторизации на текущем устройстве.

В случае, если регистрация была пройдена ранее, при входе в систему осуществляется процедура **авторизации**. Для этого:

1. Перейдите на страницу авторизации:

## Авторизация

Вход на платформу BAUM AI

[Регистрация](#)

[Восстановить пароль](#)

Рисунок 3.2 – Страница авторизации пользователя в Системе

*По ссылке «Регистрация» доступен переход на страницу регистрации.*



*По ссылке «Восстановить пароль» доступен переход на страницу восстановления пароля.*

2. Введите следующую информацию:
  - Адрес электронной почты;
  - Пароль для входа в Систему.
3. Нажмите кнопку «Войти».

## 4. Личный кабинет пользователя

**Личный кабинет** – это персональная страница на Платформе, доступ к которой есть только у одного пользователя (его владельца).

Как работать с личным кабинетом:

1. Перейдите на страницу личного кабинета одним из способов:
  - 1.1. На панели инструментов нажмите на кнопку с инициалами/аватаркой пользователя (кнопка ) -> кнопку «Профиль» (кнопка ):

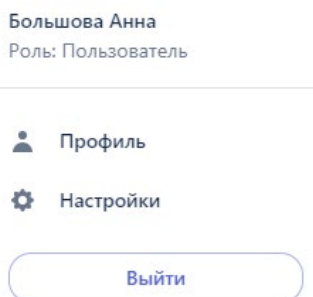


Рисунок 4.1 – Переход в личный кабинет

- 1.2. В левом верхнем углу Главного окна в разделе с информацией о пользователе нажмите кнопку с инициалами/аватаркой пользователя:

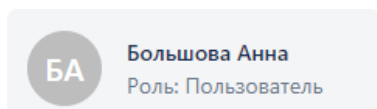


Рисунок 4.2 – Отображение данных пользователя

2. Откроется страница «Профиль» по умолчанию на вкладке «Главная»:

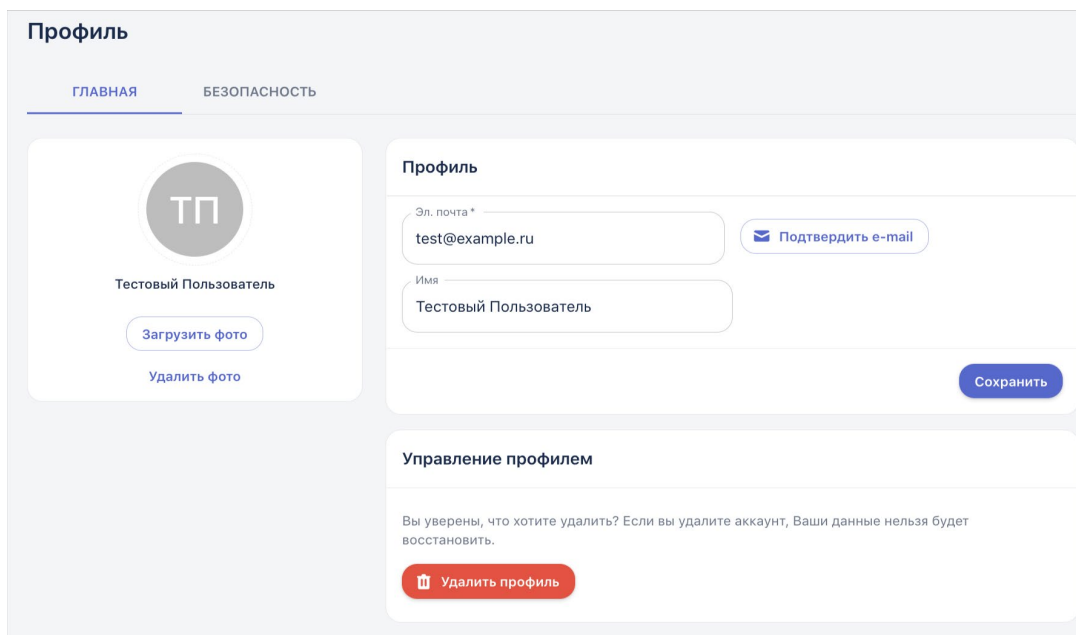


Рисунок 4.3 – Страница «Профиль»

Вкладка «Главная» состоит из:

- аватарки пользователя (по умолчанию это инициалы имени пользователя, загрузка фото не реализована в текущей версии);
- блока для ввода личной информации о пользователе – имени пользователя (логина), адреса электронной почты.
- блока «Управление профилем», в котором реализована возможность удаления своей учетной записи. Перед попыткой выполнить это действие, система выдает предупреждение о последствиях.

Чтобы изменить адрес электронной почты в поле «Электронная почта» укажите новый адрес и нажмите кнопку «Сохранить» (далее сохранение новых настроек в профиле предполагается по умолчанию). Подтверждение e-mail не является обязательным действием при смене адреса.

Профиль

---

Эл. почта\*  
user23@example.com

Подтвердить e-mail

Имя  
Большова Анна

---

Сохранить

Рисунок 4.4 – Изменение адреса электронной почты в настройках профиля

Чтобы изменить имя пользователя, в поле «Имя» укажите свой новый логин. Тогда следующая авторизация на Платформе будет выполняться уже с новым логином.

## 5. Интерфейс Платформы

### - 5.1. Меню интерфейса

Пункты меню имеют древовидную структуру и представлены в виде вложенных папок:

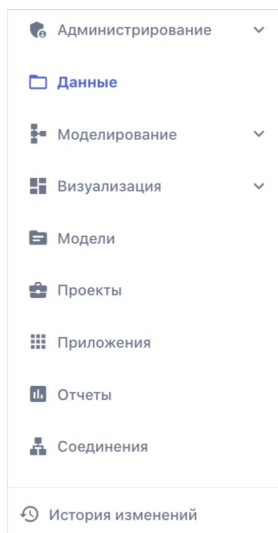




Рисунок 5.1 – Пункты меню программы **BAUM AI PREDICT**

**Состав:**

Таблица 2 – Описание пунктов меню

	<p><b>Администрирование</b></p>	<p>Данный раздел позволяет создавать ролевые модели пользователей - <i>группы</i>, согласно которым разделяются уровни доступа к разным модулям системы. Тут же осуществляется управление и назначение ролей всем пользователям системы.</p> <p>В разделе также реализована возможность настройки <i>каналов уведомлений</i> (это может быть email или telegram), которые могут быть использованы для получения автоматических оповещений от системы в случае выполнения заданных условий.</p>
	<p><b>Данные</b></p>	<p>Данный блок предназначен для загрузки в Систему <i>входных данных</i>, это могут быть файлы в различных форматах:</p> <ul style="list-style-type: none"> <li>● <b>Таблицы в форматах csv, txt, xlsx, xls</b> - табличные данные для создания интеллектуальной базы знаний, в том числе временные ряды. <i>Временной ряд</i> – это собранный в разные моменты времени статистический материал о значении каких-либо параметров исследуемого процесса. Возникают временные ряды в результате измерения некоторого показателя. Это могут быть как показатели технических систем, так и показатели природных (погодные условия), социальных, экономических и других систем. Пример временного ряда – показания датчиков на производстве, анализ которых позволяет прогнозировать выход за критические отметки целевых переменных, принять меры и предотвратить возможную поломку оборудования или аварию.</li> </ul> <p>Уже загруженные в систему файлы с данными, преобразованные в датасеты, отображаются в формате таблицы (реализовано только для файлов с расширением <i>.csv</i>). Эти датасеты используются в качестве входных данных для обучения модели искусственного интеллекта.</p> <ul style="list-style-type: none"> <li>● <b>Изображения в форматах: jpeg, jpg, png.</b></li> <li>● <b>Видео в форматах: avi, mp4.</b></li> </ul>

		<ul style="list-style-type: none"> <li>● <b>Текстовые файлы в форматах: txt, doc, docx.</b> Примеры задач для работы с текстом – классификация текста (например, определение авторства), распознавание и оцифровка рукописного текста, чтение и анализ новостного фона, определение тенденций в науке по научным статьям «Скопус», и т.д.</li> </ul>
	<b>Моделирование</b>	<p>Сердце программы, которое позволяет описывать бизнес-процессы и выполнять целевые действия.</p> <ul style="list-style-type: none"> <li>● <b>Рабочая область.</b> Для реализации бизнес-процессов заказчика на рабочей области создаются <i>блок-схемы из элементарных блоков, соединенных стрелками</i>, которые позволяют выстраивать цепочки, взаимосвязи, условия и т.д. Блок-схема отвечает на вопрос «Что делает процесс».</li> <li>● <b>Сохраненные рабочие области.</b> Рабочие области с ранее созданными блок-схемами. Есть возможность в любое время вернуться и продолжить работу над блок-схемой.</li> </ul>
	<b>Визуализация</b>	<ul style="list-style-type: none"> <li>● <b>Дашборды.</b> Динамически настраиваемые дашборды. Содержат загруженные и преобразованные массивы данных, представленные в виде таблиц, графиков, гистограмм.</li> <li>● <b>Сохранённые дашборды.</b> Дашборды, созданные пользователем или группой пользователей.</li> </ul>
	<b>Модели</b>	<p><i>Модель на основе машинного обучения</i> – это абстракция, которая обучена распознаванию определенного типа закономерностей. Хранится в виде файла.</p> <p>Для обучения модели нужны исторические данные (обучающая выборка) и значение целевой переменной (то, что прогнозируем), которое соответствует заданным историческим данным. С помощью алгоритмов машинного обучения модель наблюдает и находит зависимости между данными и целевой переменной. Эти зависимости используются моделью для нового набора данных (тестовой выборки), чтобы прогнозировать целевую переменную, которая неизвестна.</p> <p>Все модели разделяются на обучение с учителем и без учителя. <i>Обучение с учителем</i> подразделяется на две подкатегории: регрессия и классификация. В <i>регрессионных моделях</i> вывод является непрерывным. Наиболее распространенные типы регрессионных моделей – линейная регрессия, деревья решений, случайный лес, нейронная сеть. В <i>классификационных моделях</i> вывод является дискретным. Наиболее распространенные типы классификационных моделей – логистическая регрессия, метод опорных векторов.</p> <p>В отличие от обучения с учителем, <i>обучение без учителя</i> используется для того, чтобы сделать выводы из входных данных без отсылок на отмеченные результаты. Два основных метода, используемых в обучении без учителя, включают <i>кластеризацию</i> и <i>снижение размерности</i>.</p> <p><b>Сохранённые модели</b> – обученные модели, которые используются в моделировании бизнес-процессов с данными в режиме реального времени.</p>
	<b>Проекты</b>	<p>Сущность «Проект» объединяет в себе: загруженные в Систему файлы, созданные рабочие области, обученные модели ИИ, визуализацию результатов на дашбордах и сформированные отчеты. Так Система позволяет консолидировать всю информацию по проекту в одном разделе для систематизации и удобного доступа пользователей.</p>
	<b>Приложения</b>	<p>Приложение содержит «упакованную» обученную модель ИИ – без необходимости разработки отдельного кода для создания приложения с моделью. Такое приложение можно скачать и развернуть за пределами программы, интегрировать с внешними системами, настроить получение входных данных и выполнить прогнозы.</p>



	<p><b>Отчеты</b></p>	<p>Подготовленные по установленной форме отчеты. Могут храниться результаты ИИ исследований, результаты спроектированного бизнес-процесса.</p>
	<p><b>Соединения</b></p>	<p>Раздел «Соединения» предназначен для настройки подключения Платформы к внешним источникам данных с целью получения данных из других систем. В этом разделе пользователь уже может создавать новые и просматривать уже созданные коннекторы. <i>Коннектор</i> – это сущность, которая объединяют в себе источник подключения и запрос на получение данных из него.</p>

**5.2. Окно построения модели ИИ**

На примере страницы Платформы, на которой строится конструктор модели искусственного интеллекта, рассматривается, из каких частей состоит интерфейс:

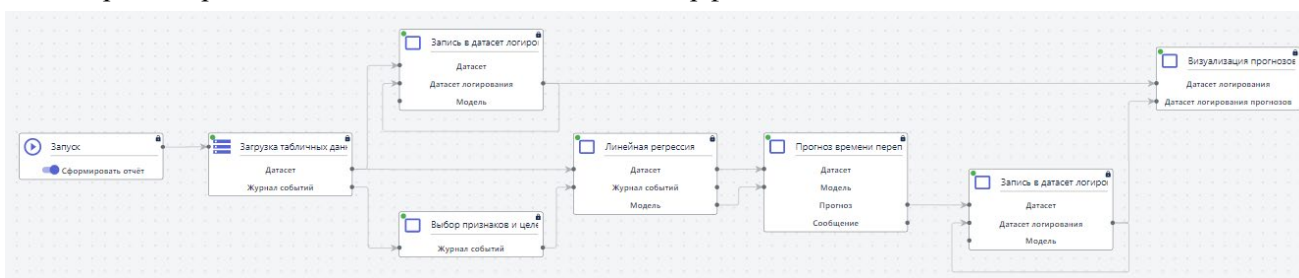
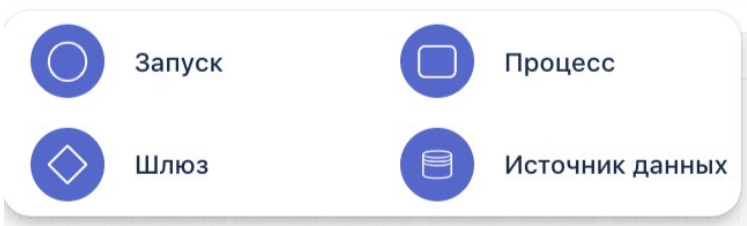






Рисунок 5.2 – Интерфейс BAUM AI PREDICT

**На верхней панели инструментов рабочей области представлены следующие кнопки:**

Таблица 4 – Кнопки панели инструментов рабочей области

<p><b>Создание рабочей области</b></p>	<p>При нажатии на кнопку открывается форма, в которой нужно указать название рабочей области:</p> <div style="text-align: center;"> <p>Введите имя рабочей области</p> <input type="text" value="Имя рабочей области"/> <p><b>Создать</b></p> </div> <p>Одна рабочая область может содержать несколько блок-схем. Реализована возможность запуска как отдельной блок-схемы, так и всех блок-схем на рабочей области.</p> <p>Рабочую область можно добавить в проект, таким образом организовав доступ к рабочей области для всех пользователей этого проекта.</p>
--	---

	<p><b>Добавить элемент</b></p>	<p style="text-align: center;">BPMN</p> <div style="text-align: center;">  </div> <p>Чтобы добавить элемент на блок-схему достаточно нажать на кнопку с элементом.</p> <p>Типы элементов нотации BPMN 2.0:</p> <table border="1" data-bbox="555 701 1497 1326"> <tr> <td data-bbox="555 701 756 801">Запуск</td> <td data-bbox="756 701 1497 801"><i>Иницилирующее событие</i> – главный элемент, обозначающий начало блок-схемы.</td> </tr> <tr> <td data-bbox="555 801 756 902">Процесс</td> <td data-bbox="756 801 1497 902">Действие, выполняемое в ходе бизнес-процесса. Является основным элементом.</td> </tr> <tr> <td data-bbox="555 902 756 1126">Шлюз</td> <td data-bbox="756 902 1497 1126">Действие предназначено для разветвления алгоритма по веткам – прохождение сценария по каждой из веток выполняется при определенных условиях. Когда выполняется раздвоение потока операций, прописывается одно условие, и при его выполнении сценарий проходит по одной из веток.</td> </tr> <tr> <td data-bbox="555 1126 756 1326">Источник данных</td> <td data-bbox="756 1126 1497 1326"><i>Объект данных</i> – это информационный объект (датасет, файл, модель и т.д.), который обрабатывается и передается в ходе выполнения бизнес-процесса. Действие предназначено для выбора уже загруженной и/или сохраненной в Системе сущности.</td> </tr> </table>	Запуск	<i>Иницилирующее событие</i> – главный элемент, обозначающий начало блок-схемы.	Процесс	Действие, выполняемое в ходе бизнес-процесса. Является основным элементом.	Шлюз	Действие предназначено для разветвления алгоритма по веткам – прохождение сценария по каждой из веток выполняется при определенных условиях. Когда выполняется раздвоение потока операций, прописывается одно условие, и при его выполнении сценарий проходит по одной из веток.	Источник данных	<i>Объект данных</i> – это информационный объект (датасет, файл, модель и т.д.), который обрабатывается и передается в ходе выполнения бизнес-процесса. Действие предназначено для выбора уже загруженной и/или сохраненной в Системе сущности.
Запуск	<i>Иницилирующее событие</i> – главный элемент, обозначающий начало блок-схемы.									
Процесс	Действие, выполняемое в ходе бизнес-процесса. Является основным элементом.									
Шлюз	Действие предназначено для разветвления алгоритма по веткам – прохождение сценария по каждой из веток выполняется при определенных условиях. Когда выполняется раздвоение потока операций, прописывается одно условие, и при его выполнении сценарий проходит по одной из веток.									
Источник данных	<i>Объект данных</i> – это информационный объект (датасет, файл, модель и т.д.), который обрабатывается и передается в ходе выполнения бизнес-процесса. Действие предназначено для выбора уже загруженной и/или сохраненной в Системе сущности.									
	<b>Графики</b>	После успешной отработки пайплайна здесь будет отображаться список доступных для визуализации графиков.								
	<b>Таблицы</b>	После успешной отработки пайплайна здесь будет отображаться список доступных для визуализации таблиц.								
	<b>Изображения</b>	После успешной отработки пайплайна здесь будет отображаться список доступных для визуализации изображений.								
	<b>Описание</b>	Отображается описание обученной модели ИИ.								

- **Рабочая область.** Предназначена для построения графической бизнес-модели. С помощью блок-схемы выполняется настройка последовательных операций обработки данных, обучение моделей, выстраивается «pipeline» (цепочка процессов преобразования).

В правом нижнем углу рабочей области отображается миниатюра бизнес-процесса:

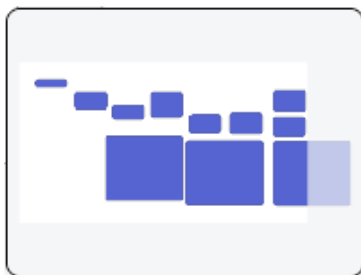


Рисунок 5.3 – Миниатюра моделируемого бизнес-процесса

- **Карточка элемента блок-схемы.** В карточке прописываются условия обработки данных – откуда берутся данные, какие операции над ними выполняются. Для перехода в карточку нажмите кнопку «Настройки» на элементе. Чтобы скрыть карточку элемента, используйте кнопку «Свернуть».

Карточка содержит *программный код* с основными блоками:

- входные данные (источник входных данных);
  - преобразование данных;
  - выходные данные – куда и в каком виде передаются преобразованные данные.
- **Вкладки с рабочими областями** - вы можете работать с несколькими рабочими областями одновременно, для удобного перехода между ними используйте вкладки.

### - 5.3. Настройка внешнего вида интерфейса

Чтобы изменить настройки интерфейса, кликните на свою аватарку и перейдите в пункт меню «Настройки». Откроется окно, в котором можно изменить следующее:

- **Тема** – светлая или темная.
- **Адаптивные размеры шрифтов** – адаптация размеров шрифтов программы для работы с ней в мобильных устройствах.
- **Компактный вид** - фиксированная ширина на некоторых экранах.
- **Закругленные углы.** Если выбрать настройку, то окна в программе будут иметь закругленные углы. Иначе окна будут иметь прямые углы.

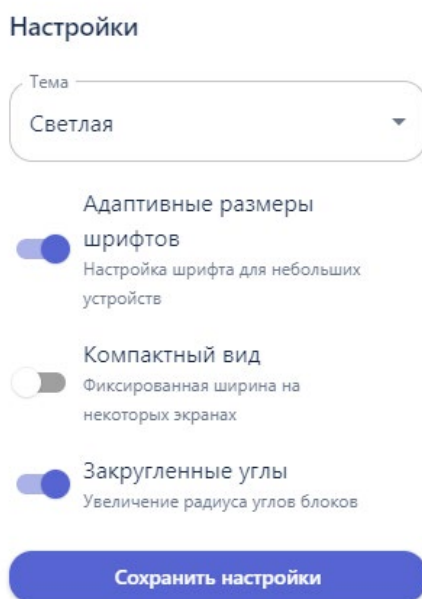


Рисунок 5.4 – Настройки внешнего вида интерфейса программы

Чтобы применить выбранные настройки к интерфейсу нажмите кнопку «Сохранить настройки».

### - 5.4. Встроенные функции

Каждая функция относится к одному из модулей Системы: *модулю препроцессинга входных данных, модулю машинного обучения, модулю нейронных сетей/глубокого обучения, модулю анализа данных или модулю визуализации.*

Описание всех функций представлено в Приложении 1 настоящего документа.

За каждым типом блока закреплен определенный набор функций:

#### 5.4.1 Функции элемента «Источник данных»

Таблица 5 – Набор функций элемента «Источник данных»

Группа	Функция
Загрузка данных	1. Загрузка табличных данных из коннектора
	2. Загрузка табличных данных
	3. Загрузка текстовых файлов
	4. Загрузка модели
Spark	5. Загрузка табличных данных из файла CSV (Spark)
	6. Загрузка табличных данных из папки CSV (Spark)

Группа	Функция
	7. Загрузка модели
	8. Загрузка табличных данных из коннектора (Spark)

#### 5.4.2 Функции элемента «Процесс»

Таблица 6 – Набор функций элемента «Процесс»

Группа	Подгруппа	Функция
Анализ данных	Загрузка данных	1. Преобразование данных во временной ряд
	Препроцессинг	2. Стабилизация дисперсии
		3. Стандартизация
		4. Дифференцирование временного ряда
		5. One-Hot Encoding
		6. Создание признаков для временного ряда
		7. Кодирование целевого признака
		8. Порядковое кодирование категориальных признаков
		Тесты на нормальность распределения
	Тесты на стационарность временного ряда	10. Тест Дики-Фуллера
	-	11. Выбор признаков и целевых признаков
		12. Матрица корреляции
		13. Косинусное расстояние
		14. Поиск пропущенных значений

		15. Анализ временных рядов
		16. Визуализация Real Time
		17. Запись в датасет логирования
<b>Машинное обучение</b>	<b>Классификация</b>	18. Логистическая регрессия
		19. Модель XGBClassifier
		20. Дерево решений для классификации
		21. Случайный лес для классификации
		22. Categorical Naive Bayes
		23. Multinomial Naive Bayes
		24. Complement Naive Bayes
		25. Gaussian Naive Bayes
		26. Bernoulli Naive Bayes
	<b>Обучение без учителя</b>	27. Кластеризация DBSCAN
		28. Метод локтя K-Means
		29. Кластеризация K-Means
		30. Агломеративная иерархическая кластеризация
		31. Изоляционный лес
	<b>Регрессия</b>	32. Линейная регрессия
33. Полиномиальная регрессия		
34. Дерево решений для регрессии		

		35. Случайный лес для регрессии
		36. Метод опорных векторов для регрессии
		37. Байесовская гребневая регрессия
		38. Метод k-ближайших соседей для регрессии
	<b>Авторегрессия</b>	39. ARIMA/SARIMAX
<b>Глубокое обучение</b>	<b>Классификация</b>	40. Классификация изображений (Валидация)
		41. Классификация (табличные данные)
	<b>Регрессия</b>	42. Регрессия (табличные данные)
<b>Предобработка данных</b>	-	43. Заполнение пропусков
		44. Сглаживание временного ряда
		45. Срез временного ряда по индексу
<b>Отправка уведомлений</b>	-	46. Отправка уведомлений
<b>Spark</b>	-	47. Выбор признаков и целевых признаков
		48. Разделение датасета на обучающую и тестовую выборки
		49. Валидация модели
		50. Сохранение датасета Spark в CSV
		51. Прогноз модели
		52. Косинусное расстояние
	<b>Преппроцессинг</b>	53. Порядковое кодирование признаков

		54. Нормализация признаков
	<b>Классификация</b>	55. Модель градиентного бустинга Spark для бинарной классификации
	<b>Кластеризация</b>	56. Кластеризация Spark DBSCAN

“–” в таблице означает, что у функции нет подгруппы, и она напрямую относится к группе функций.



## 6. Загрузка данных в систему

Сразу после авторизации в Системе открывается её начальная страница – раздел «Данные».

Раздел «Данные» имеет внешний вид аналогичный проводнику файлов в операционной системе компьютера. Пользователи имеют возможность создавать папки, загружать файлы и делать структуры, удобные для личного использования. По умолчанию для нового пользователя раздел Данные выглядит следующим образом:

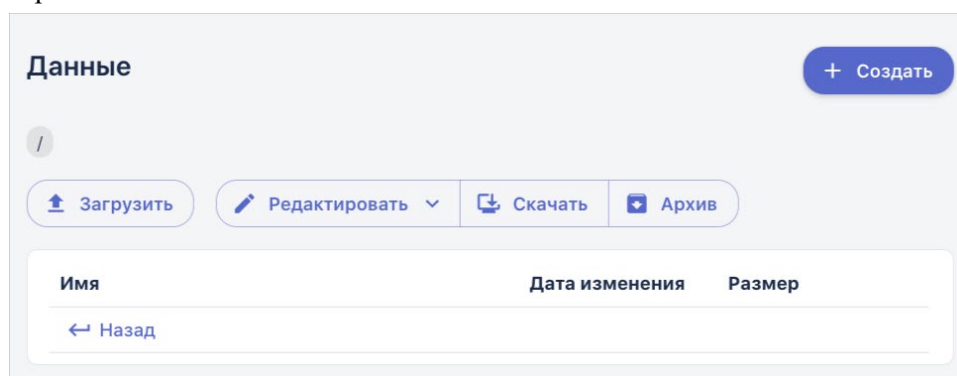


Рисунок 6.1 – Вид раздела Данные

Основные кнопки для работы с разделом:

1. «Создать» - позволяет создавать папки и разметки.
2. «Загрузить» - осуществление непосредственно загрузки файлов в систему
3. «Редактировать» - не работает в данной версии системы
4. «Скачать» - не работает в данной версии системы
5. «Архив» - не работает в данной версии системы

В данном разделе вы можете создавать папки, добавлять в них файлы, создавая удобную и понятную структуру. Также в разделе осуществляется создание разметки для видео и изображений, и добавление папок классификации для дальнейшего использования в обучении искусственного интеллекта.

Вы можете загружать в систему файлы разных форматов и типов. Данные можно загружать напрямую в основную директорию, или создавая новые папки внутри.

### - 6.1. Создание новой папки



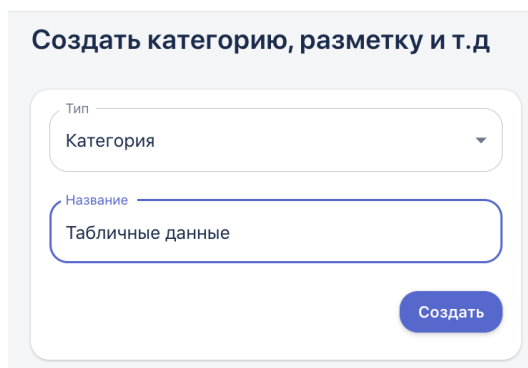


Рисунок 6.2 – Создание новой папки в разделе Данные

После этого папка появится в разделе Данные:

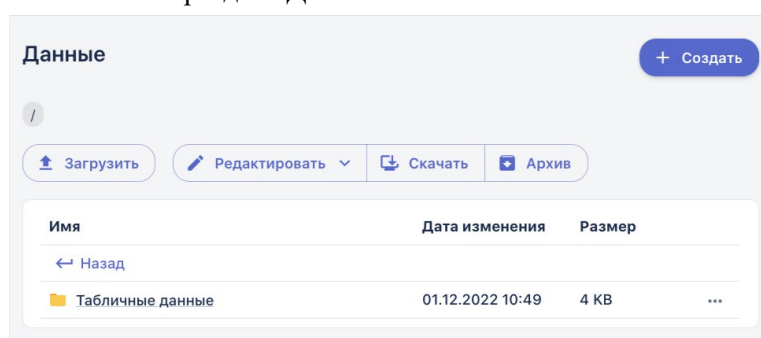


Рисунок 6.3 – Новая папка в разделе Данные

Обратите внимание, что папка будет добавлена в том разделе, из которого вы нажали кнопку «Создать». Т.е. далее вы можете перейти в папку «Табличные данные» и создать внутри еще одну категорию-папку.

## - 6.2. Загрузка файлов

Для того чтобы загрузить файлы в папку, кликните на неё и перейдите в ее содержимое. Далее нажмите кнопку «Загрузить»:

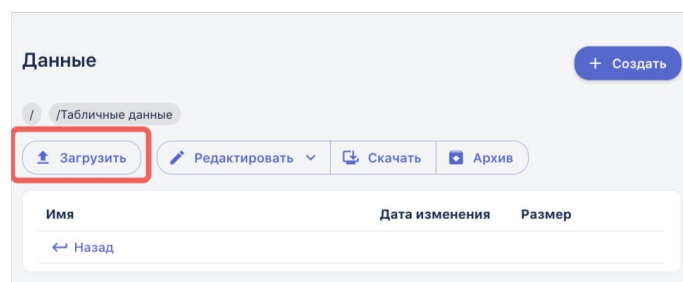


Рисунок 6.4 – Загрузка файлов в папку

В открывшемся окне левой кнопкой мыши нажмите на ссылку выбора файла. Указать путь к файлу для загрузки на вашем ПК. Второй вариант – перенести файлы с локального компьютера в этот раздел по технологии «drag n drop».

Выбранные файлы отобразятся в нижней части окна загрузки:

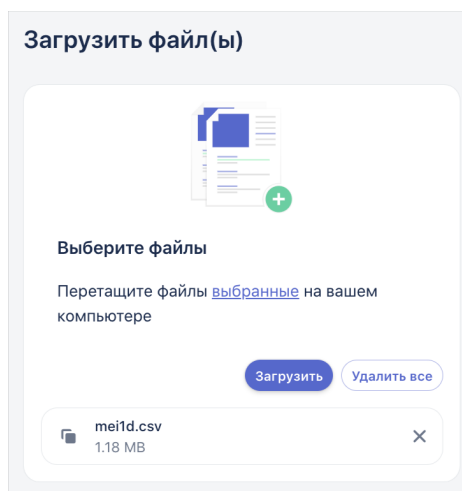


Рисунок 6.5 – Отображение выбранного файла

При необходимости выбранные файлы можно удалить по одному, нажав на крестик, или все вместе, нажав кнопку «Удалить все».

Для того чтобы загрузить выбранные файлы, нажмите на кнопку «Загрузить». Файлы отобразятся в папке:

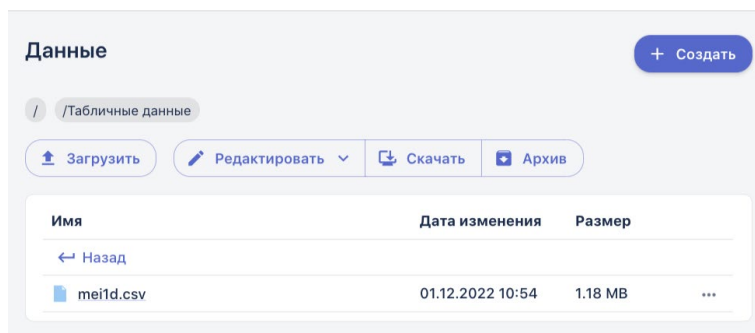


Рисунок 6.6 – Загруженный файл в разделе Данные

### - 6.3. Предпросмотр данных

В системе реализована функция предпросмотра загруженных файлов следующих форматов:

- табличные данные формата .csv
- изображения в форматах: jpeg, jpg, png.
- видео в форматах: avi, mp4.

Если кликнуть на файл с изображением или видео - его предпросмотр начнется прямо в окне, чтобы выйти из режима предпросмотра - просто кликните в любую область экрана за пределами медиафайла.

При предпросмотре табличных данных в формате csv, в нижней части экрана отобразится окно с тремя вкладками, содержащими описательно статистический анализ.

На вкладке «Подробнее» (показано на примере файла с временным рядом) пользователь имеет возможность увидеть:

- размер датасета (количество строк, столбцов)

- гистограммы распределения каждого признака (числового)

vol\_size\_used.csv ×

[Подробнее](#)   [Компактно](#)   [Столбцы](#)

Количество строк: 73349. Количество столбцов: 4

date	value	name	clock
2022-05-14 10:05:00	0.02	Logicalused by volume: pool39/ISCSI39	1652522700
2022-05-14 10:10:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523000
2022-05-14 10:15:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523300
2022-05-14 10:20:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523600
2022-05-14 10:25:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523900
2022-05-14 10:30:00	0.02	Logicalused by volume: pool39/ISCSI39	1652524200
2022-05-14 10:35:00	0.02	Logicalused by volume: pool39/ISCSI39	1652524500
2022-05-14 10:40:00	0.02	Logicalused by volume: pool39/ISCSI39	1652524800

Рисунок 6.7– Вкладка «Подробнее» окна отображения данных о датасете

На вкладке «Компактно» отображается состав строк и столбцов файлов:

vol\_size\_used.csv ×

[Подробнее](#)   [Компактно](#)   [Столбцы](#)

date	value	name	clock
2022-05-14 09:35:00	0.02	Logicalused by volume: pool39/ISCSI39	1652520900
2022-05-14 09:40:00	0.02	Logicalused by volume: pool39/ISCSI39	1652521200
2022-05-14 09:45:00	0.02	Logicalused by volume: pool39/ISCSI39	1652521500
2022-05-14 09:50:00	0.02	Logicalused by volume: pool39/ISCSI39	1652521800
2022-05-14 09:55:00	0.02	Logicalused by volume: pool39/ISCSI39	1652522100
2022-05-14 10:00:00	0.02	Logicalused by volume: pool39/ISCSI39	1652522400
2022-05-14 10:05:00	0.02	Logicalused by volume: pool39/ISCSI39	1652522700
2022-05-14 10:10:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523000
2022-05-14 10:15:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523300
2022-05-14 10:20:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523600
2022-05-14 10:25:00	0.02	Logicalused by volume: pool39/ISCSI39	1652523900
2022-05-14 10:30:00	0.02	Logicalused by volume: pool39/ISCSI39	1652524200

Рисунок 6.8 – Вкладка «Компактно» окна отображения данных о датасете

На вкладке «Столбцы» для каждого признака отображается:

- количество пропусков (Missing) в абсолютном и в процентном выражении, под пропуском понимается пустая ячейка в таблице
- среднее значение (Mean)
- стандартное отклонение (Std. Deviation)
- квантили (quantiles):

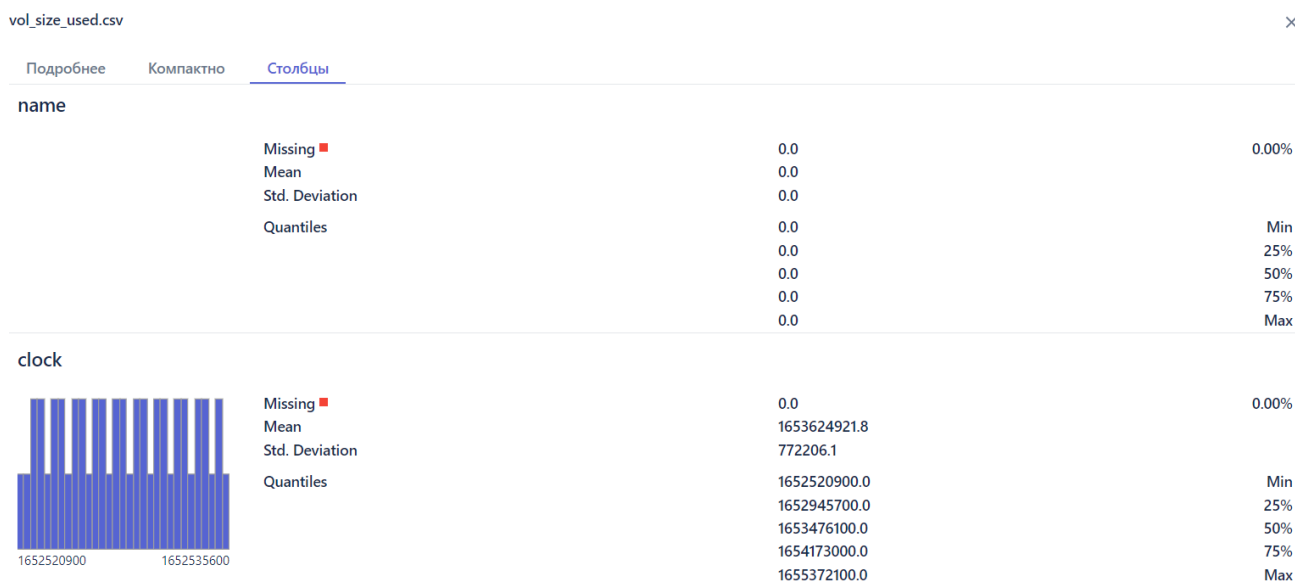


Рисунок 6.7– Вкладка «Столбцы» окна отображения данных о датасете

### 6.4. Взаимодействие с данными

Загруженный файл или созданную папку можно скачать, скачать архивом, переименовать или удалить. Для этого нажмите на три точки в правой части раздела и выберите соответствующую кнопку:

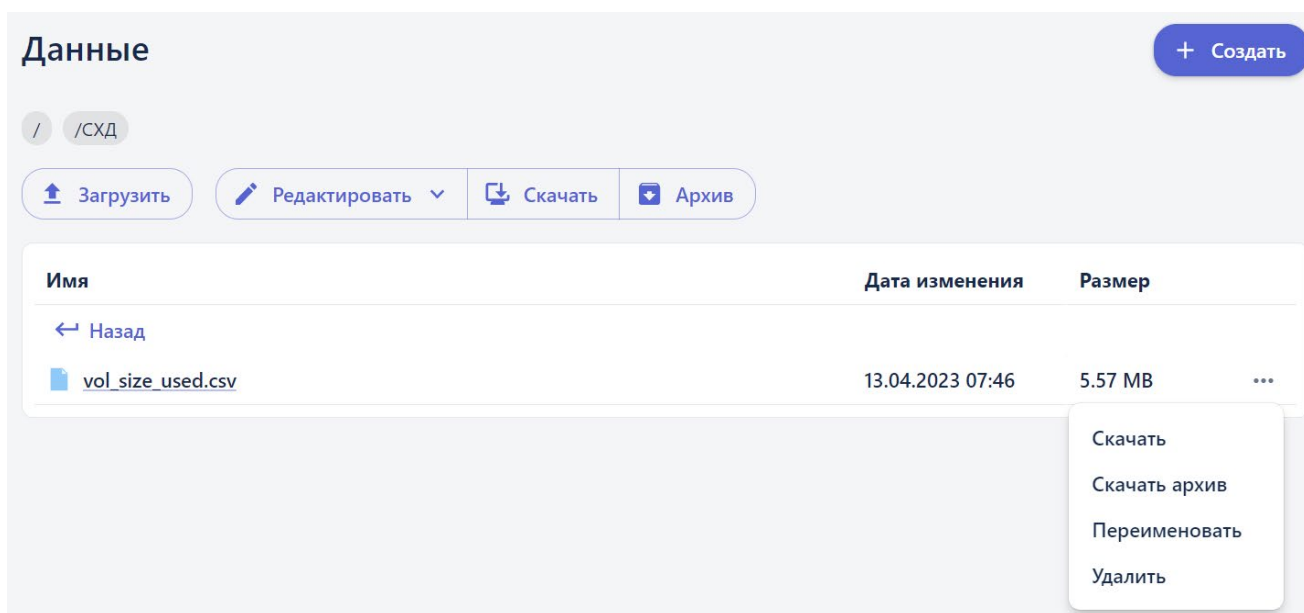


Рисунок 6.8 – Удаление или скачивание файла из раздела Данные

Вы можете удалять файлы по отдельности или сразу целую папку.



## 7. Создание модели ИИ

В разделе «Моделирование» осуществляется процесс построения блок-схем - соединение последовательных функции процессов анализа, обработки и преобразования исходных данных для построения и обучения моделей искусственного интеллекта. Построение таких блок-схем осуществляется на рабочих областях.

### - 7.1 Создание новой и открытие сохраненной рабочей области

#### ■7.1.1. Создание новой рабочей области

1. Перейдите в пункт меню системы Моделирование -> Рабочая область:

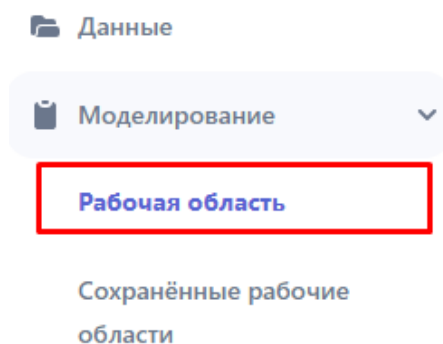


Рисунок 7.1.1 – Переход на рабочую область

Откроется страница с пустой рабочей областью:

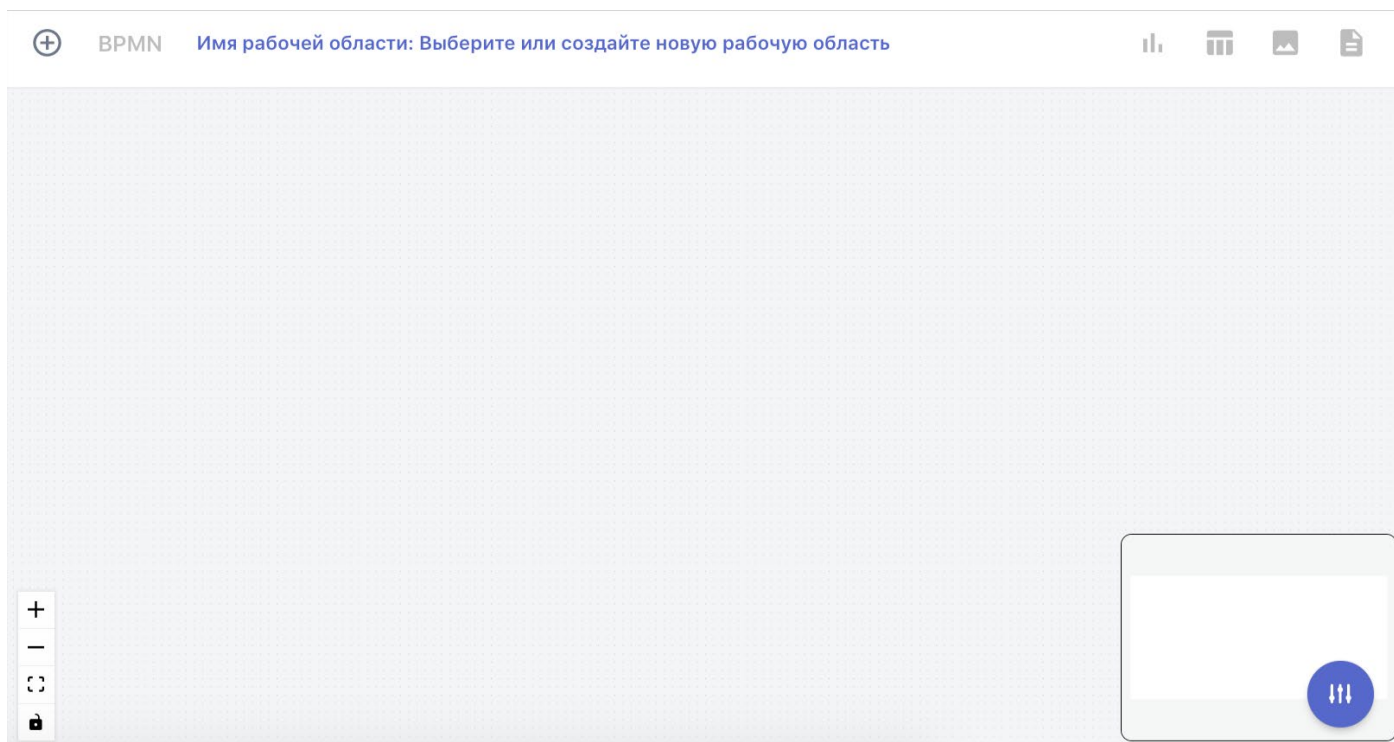


Рисунок 7.1.2 – Пустая рабочая область



В открывшейся форме введите название новой рабочей области и нажмите кнопку «Создать»:

Введите имя рабочей области

Имя рабочей области

Создать

Рисунок 7.1.3 – Создание новой рабочей области

На панели инструментов отобразится название созданной рабочей области:

 BPMN    Имя рабочей области: Animals

Рисунок 7.1.4 – Отображение названия рабочей области

После присвоения названия рабочая область автоматически сохраняется в раздел «Сохраненные рабочие области».

### ■ 7.1.2. Открытие сохраненной рабочей области

Для того чтобы открыть ранее созданную рабочую область, перейдите в пункт меню системы Моделирование -> Сохраненные рабочие области. Откроется страница «Рабочие области»:











Рабочие области		
🔍 Поиск		
Название	Создан	
Прогнозирование загрузки СХД (Лаба 1)	17.03.2023 12:12	 
Информирование оператора СХД	28.02.2023 18:03	 
Превентивный ремонт дисков	22.02.2023 11:58	 
Диагностика СХД	23.03.2023 12:16	 
Прогнозирование загрузки СХД (Энерго)	27.02.2023 05:07	 

Рисунок 7.1.2 – Страница со списком созданных в Системе рабочих областей

На странице отображаются рабочие области. Чтобы открыть сохраненную рабочую область кликните на её название. Чтобы удалить ненужную рабочую область - кликните кнопку удалить в правой части строки с названием области.

## - 7.2. Построение блок схемы

Предварительным условием для построения блок схемы является:

- 1) Загруженные в систему данные (файлы или датасеты)

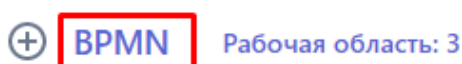
- 2) Созданная рабочая область (на одной рабочей области можно создать неограниченное количество блок-схем).

Построение блок-схемы осуществляется путем добавления на рабочую область элементов (блоков) и соединение их между собой.

### ■7.2.1 Блок «Запуск»

Блок «Запуск» обозначает начало блок-схемы, и всегда является её первым элементом. Так как на рабочей области может быть несколько блок-схем, именно по блоку «Запуск» определяется их количество, и идентифицируется принадлежность блоков к той или иной блок-схеме.

#### BPМN



Откроется меню выбора блоков – библиотека графических элементов нотации BPМN 2.0:

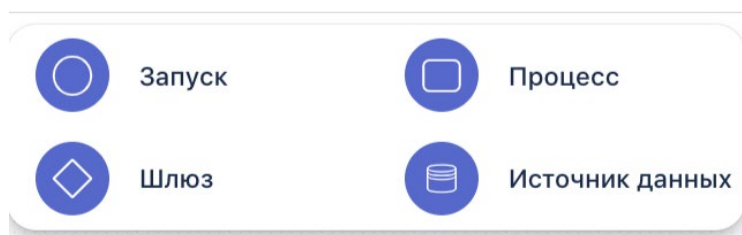
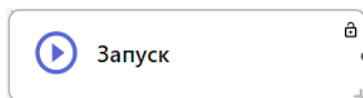


Рисунок 7.2.1 – Меню блоков

Выберите элемент «Запуск». Выбранный элемент будет добавлен на рабочую область:



При необходимости размер блока можно увеличить или уменьшить, потянув за уголок в правой нижней части элемента:

Обратите внимание, что в правом нижнем углу рабочей области отображается уменьшенное изображение создаваемой блок-схемы, эта функция упрощает навигацию по полотну.

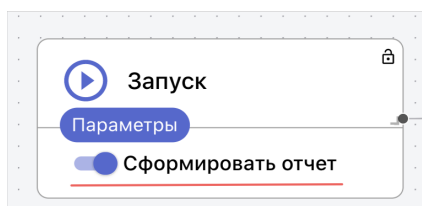


Рисунок 7.1.2 – Опция формирования отчета на блоке «Запуск»



### ■ 7.2.2. Блок «Источник данных»

Следующим элементом в блок-схеме после «Запуска» всегда является «Источник данных» - блок, который определяет какие данные будут использоваться в сценарии.



Для объединения элементов в блок-схему, их требуется соединить между собой. Для этого нажмите на точку выхода блока, которая отображается в виде круглой точки на правой грани блока, и перетащите мышью появившуюся стрелку в сторону нужного блока, как показано на рисунке:

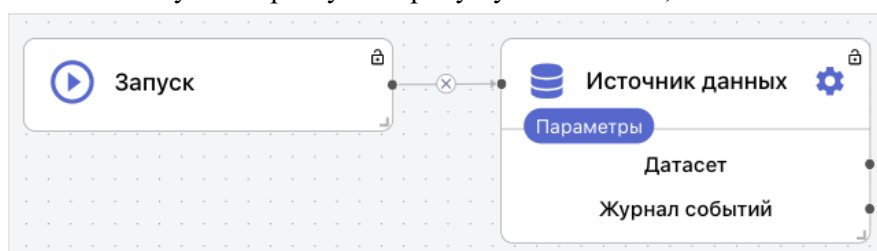


Рисунок 7.2.1 – Соединение блоков между собой



На блоке «Источник данных» отображаются два компонента:

- 1) «Датасет» - это непосредственно сами данные.
- 2) «Журнал событий» содержит информацию обо всех преобразованиях с данными, которые выполняются в текущем блоке пайплайна. Ведение журнала позволяет сохранить историю преобразований над данными, и при необходимости выполнить обратное преобразование.



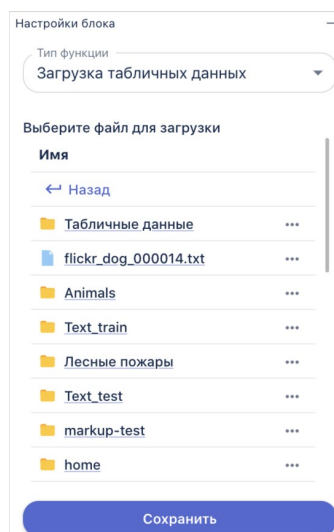


Рисунок 7.2.2 – Настройка параметров блока «Источник данных»

\* Дальнейшие действия по настройке блоков указаны для примера. При работе с Системой пользователь должен выбрать необходимые параметры исходя из своей задачи и загруженных данных.

После добавления на рабочую область, для элемента «Источник данных» по умолчанию выбрана функция: тип функции «Загрузка данных» -> функция «Загрузка табличных данных». Это можно увидеть в верхней части окна настройки параметров элемента:

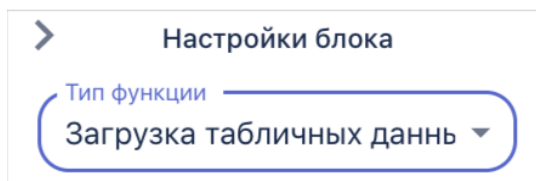


Рисунок 7.2.3 – Отображение типа функции в параметрах блока

Для того чтобы посмотреть другие доступные функции, нужно нажать на выпадающий список. Весь список доступных функций и их описания доступны в [Таблице 18.1 – Перечень автоматизированных функций элемента «Источник данных»](#).

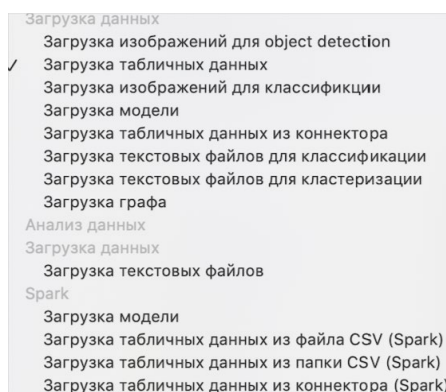


Рисунок 7.2.4 – Список возможных функций элемента «Источник данных»

В разделе «Выберите файл» отображается структура папок из разделе «Данные», чтобы выбрать файл достаточно перейти в нужную папку и кликнуть на три точки в правой части строки с названием файла и нажать «Выбрать»:

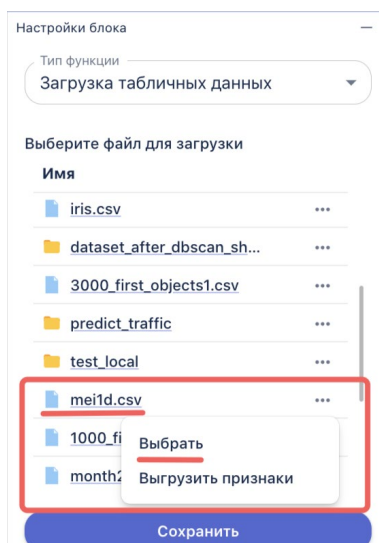


Рисунок 7.2.5 – Отображение папок и файлов из раздела «Данные» в параметрах блока «Источник данных»

После этого в нижней части окна отобразится его название:

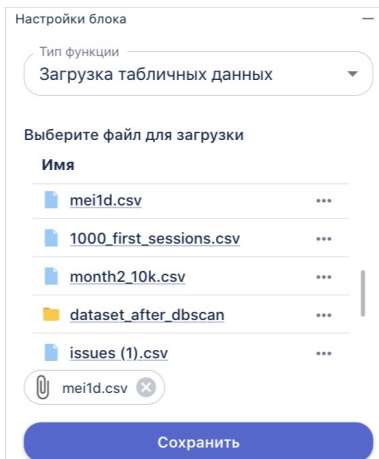


Рисунок 7.2.6 – Отображение выбранного файла, предназначенного для загрузки в блок-схему

Кнопка «Выгрузить признаки» используется для других блоков, где в настройках необходимо указать целевые признаки для конкретной функции.

Если кликнуть на название файла на рабочей области отобразится его предпросмотр:

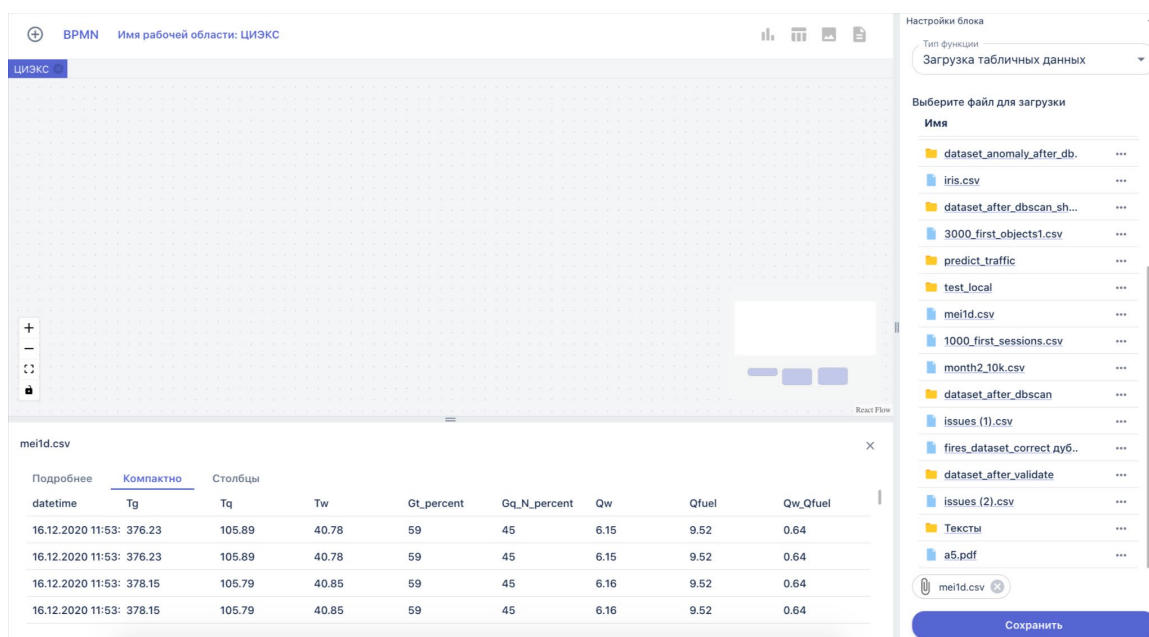


Рисунок 7.2.6.1 – Предпросмотр файла на рабочей области

Для сохранения выбранных настроек нажмите на панели параметров кнопку «Сохранить» (далее сохранение настроек элемента предполагается по умолчанию). Для удаления неверно добавленного на рабочую область элемента предусмотрена кнопка «Удалить блок».

Любой блок можно переименовать, чтобы дать ему понятное название, отображающее суть происходящего процесса. Для этого дважды щелкните левой кнопкой мыши на текущее название элемента в рабочей области и измените его. Чтобы новое название сохранилось достаточно щелкнуть мышью в любом месте на рабочей области.

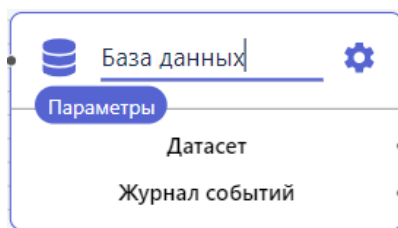


Рисунок 7.2.7 – Ввод названия элемента

Для того чтобы удалить блок, кликните по нему правой кнопкой мыши и нажмите «Удалить».

### ■ 7.2.3. Блок «Процесс»

Блок «Процесс» предназначен для выполнения операций над данными. Блок-схема может содержать несколько элементов «Процесс», настроенных пользователем для выполнения определенных задач. Весь список доступных функций и их описания доступны в [Таблице 18.2 – Перечень автоматизированных функций элемента «Процесс»](#).



Далее будет показан принцип настройки свойств блока на примере одной функции. Выбор функции определяется типом решаемой задачи.

Для примера выберите для элемента функцию: раздел «Машинное обучение» -> функция «Разделение датасета на обучающую и тестовую выборки»:

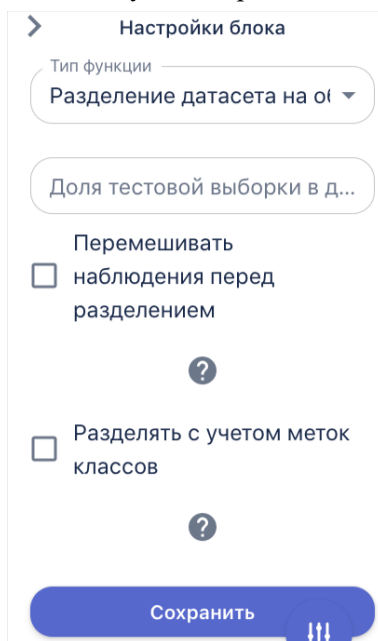


Рисунок 7.2.8 – Панель свойств блока «Процесс»

Далее осуществляется настройка параметров следующим образом:

- В разделе «Параметры» -> в поле «Доля тестовой выборки в датасете» введите значение 0.2;
- Оставьте пустым поле «Перемешивать наблюдения перед разделением». Рядом с полем есть подсказка, что не рекомендуется перемешивать наблюдения во временных рядах (выбирайте действие в зависимости от типа входных данных);
- Установите галочку в поле «Разделять с учетом меток классов» – применяется для задач классификации (выбирайте действие в зависимости от решаемой задачи).

Измените название элемента на «Сплит датасета»:

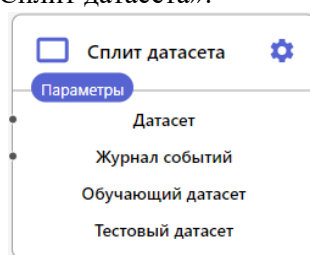


Рисунок 7.2.9 – Отображение блока на рабочей области после настройки его параметров и ввода названия

Обратите внимание, что соединение элемента «Процесс» с другими элементами блок-схемы выполняется только после настройки и сохранения его параметров. Это связано с тем, что каждая функция имеет свой набор компонентов, который отображается на элементе после сохранения его настроек.

Соединить элемент «Процесс» с предыдущими элементами блок-схемы можно следующим образом:

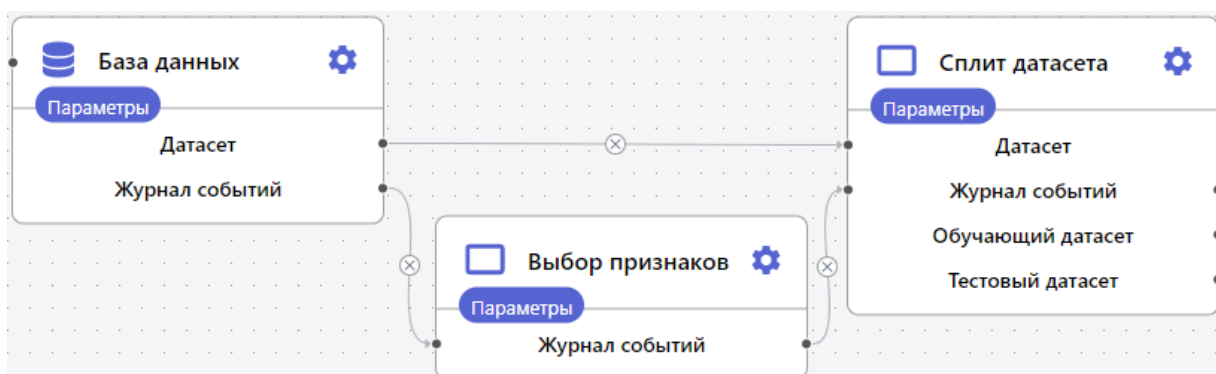


Рисунок 7.2.10 – Соединения между несколькими блоками (организация связей/пробрасывание данных)

Обратите внимание, что соединять можно только идентичные (одноименные) компоненты блоков.

Больше примеров построения блок схем, сценарии сохранения моделей, отображения на рабочих областях таблиц, графиков и изображений можно прочитать в разделе [Примеры работы с Платформой](#).

### - 7.3. Запуск блок-схемы на рабочей области

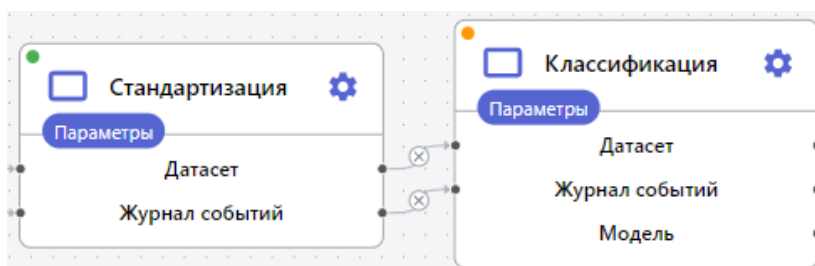


Рисунок 7.3.1 – Последовательная обработка блоков пайплайна

В случае, если блок не отработал (например, вследствие того, что был неправильно настроен или из-за неверных входных данных), на нем появится индикатор красного цвета. Такой блок необходимо проверить, изменить настройки и запустить заново:

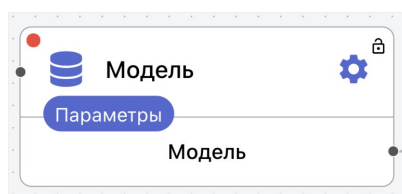


Рисунок 7.3.2 – Индикатор ошибки при обработке блока

## 8. Сохранение модели ИИ

Для сохранения модели ИИ, обученной выполнению какой-либо задачи, на блок-схему необходимо добавить специальный блок. Этот блок настраивается для элемента «Процесс» с помощью функции «Управление моделями» → функция «Сохранение модели»:

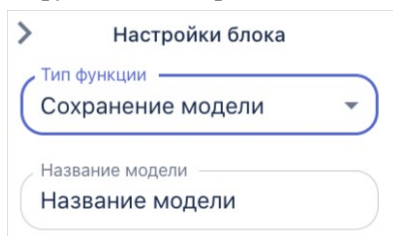


Рисунок 7.3.1 – Настройка блока для сохранения обученной модели ИИ

После успешной отработки блок схемы, которая содержит такой блок, название модели будет добавлено в список сохраненных моделей в меню «Модели»:

Модели		
Поиск		
Название	Создан	
CLF_spark_model_2022-10-18_08:10:40.664728.UTC	18.10.2022 11:10	+ ↓ 🗑
CLF_spark_model_2022-10-18_08:05:43.798686.UTC	18.10.2022 11:05	+ ↓ 🗑
DBSCAN_spark_model_2022-10-18_06:41:21.483562.UTC	18.10.2022 09:41	+ ↓ 🗑
mei_2022-10-17_11:22:53.663173.UTC	17.10.2022 14:22	+ ↓ 🗑
animals_2022-10-14_13:50:05.011675.UTC	14.10.2022 16:50	+ ↓ 🗑
object_detection_2022-10-12_12:02:36.819048.UTC	12.10.2022 15:02	+ ↓ 🗑
Модель прогнозирования пожаров_2022-10-07_13:23:27.446806.UTC	07.10.2022 16:23	+ ↓ 🗑
Модель прогнозирования пожаров_2022-10-07_12:46:27.554901.UTC	07.10.2022 15:46	+ ↓ 🗑

Рисунок 7.3.2 – Вкладка меню «Модели»

Обратите внимание, что модель будет сохраняться столько раз, сколько будет запущена блок-схема с элементом «Сохранение модели», при этом в разделе будет меняться временная отметка создания записи.

Модель можно использовать следующими способами:



- **model.pkl** - сама модель.

- **vars\_dict.pkl** - словарь преобразований. Преобразования необходимо сохранять, чтобы при анализе новой порции данных над ними выполнялись все те же преобразования, что и при обучении модели.
  - **info.json** - служебный файл, куда прописывается тип модели.
3. Использовать при построении новой блок схемы в качестве источника данных (например, для целей прогнозирования). Пример можно посмотреть в разделе [14.4 Работа с данными в режиме реального времени](#).
  4. Создать коннектора с обученной моделью (например, для распознавания объектов на видео или изображениях). Это позволит проверить обучение модели на новой порции данных. Подробнее описано в разделе [14.2.1 Проверка обученной модели на локальных данных](#).



## 9. Графическое представление информации на рабочей области

После сборки и успешного запуска блок-схемы на рабочей области есть возможность посмотреть результаты обучения созданной модели в виде графиков, таблиц и изображений, которые можно вывести прямо на рабочую область.

**Примечание:** под сборкой имеется в виду, что на рабочую область добавлены все элементы блок-схемы, и они последовательно соединены между собой. А успешным считается запуск блок-схемы, когда все ее элементы отработали с...



Вы можете нажать на те кнопки визуализации, которые подсвечиваются фиолетовым цветом. При этом голубым подсвечиваются только те иконки, которые актуальны для запущенной схемы, т.к. блоки могут содержать разные функции, имеющие разное графическое представление.

### - 9.1. Графики



Полный список доступных в Системе графиков с объяснением интерпретации результатов доступен в [Базе знаний](#). За каждым типом блока закреплен определенный набор графиков, например, для блока Анализ данных -> Анализ временных данных доступны следующие графики: Линейный график, ACF/PACF, Декомпозиция, Свечной график, Time profile, Extended, Bollinger Bands Stochastic Oscillator.

Для того чтобы добавить график на рабочую область из выпадающего списка выберите нужное название. Например, «Time profile временного ряда» в анализе временных рядов:

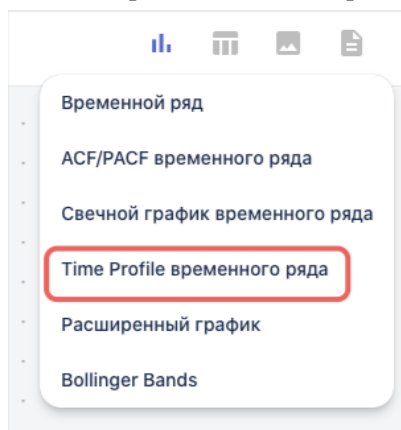


Рисунок 9.1.1 – Список доступных графиков для анализа временных рядов

В результате на рабочую область будет добавлен график:



Рисунок 9.1.2 – Time profile временного ряда на рабочей области

Для ряда графиков доступен выбор из выпадающего списка признака, для которого составляется визуализация. Например, можно два раза выбрать график Time profile и для одного указать признак Tq, а для другого Tw и сопоставлять их значения одновременно на рабочей области:

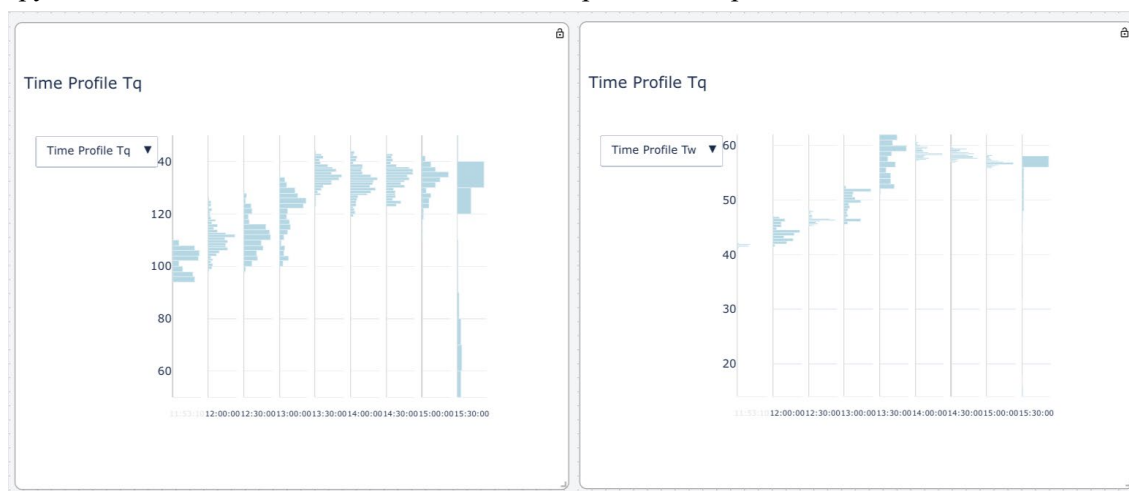


Рисунок 9.1.3 – Отображение одного типа графика для разных признаков

Также в правом углу рамки графической визуализации, при наведении на неё курсора, отображаются следующие кнопки:



Рисунок 9.1.4 – Кнопки для работы с визуализацией

где:

	- скачать график, как рисунок в формате png
	- при нажатии лупы, можно увеличить любой участок графика
	- при нажатии данной кнопки, можно передвигаться по плоскости графика в любую сторону
	- приблизить, увеличить
	- уменьшить
	- автоматически подогнать размер под границы
	- вернуться к исходному виду

Таблица 9.1 - Функционал кнопок работы с визуализацией

В программе также реализована возможность создания графиков с данными, получаемыми в режиме реального времени. Для таких графиков существует отдельный блок, находящийся в разделе «Анализ данных» -> «Визуализация Real Time». Главное отличие от стандартного блока «Визуализация» в том, что для каждого графика необходимо задавать число периодов в окне и период окна - параметры, которые определяют интервал, который будет отображаться на графике в рабочей области.

## - 9.2. Таблицы



Например:

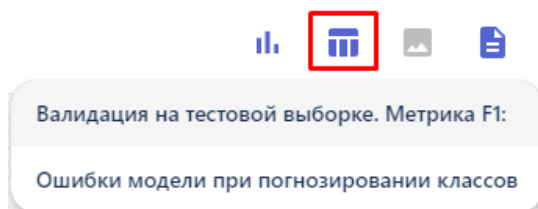
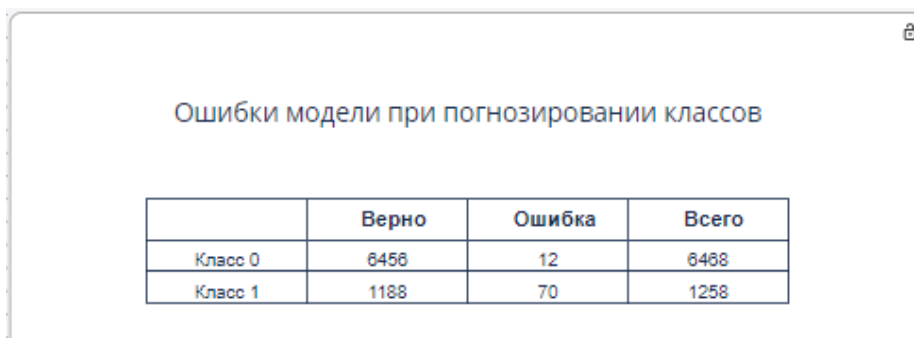


Рисунок 9.2.1 – Список таблиц для визуализации

Чтобы отобразить таблицу на рабочей области, нажмите на её название в выпадающем списке и она появится на экране, например:



Ошибки модели при прогнозировании классов

	Верно	Ошибка	Всего
Класс 0	6456	12	6468
Класс 1	1188	70	1258

Рисунок 9.2.2 – Блок визуализации «Ошибки модели при прогнозировании классов»

#### - 9.4. Описание модели



Описание варьируется в зависимости от функций, которые были применены в модели. Например, для блок схемы, где присутствовали элементы «Стандартизация», «Валидация» и «XGB классификации» описание будет выглядеть следующим образом:

Лучшие гиперпараметры при кросс-валидации:

max_depth	5
n_estimators	50

Лучшая метрика F1 при кросс-валидации:  
0.959

Время обучения полной обучающей выборки в сек:  
13.875

Модель

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1, enable_categorical=False, gamma=0, gpu_id=-1, importance_type=None, interaction_constraints="", learning_rate=0.300000012, max_delta_step=0, max_depth=5, min_child_weight=1, missing=nan, monotone_constraints='()'), n_estimators=50, n_jobs=40, num_parallel_tree=1, predictor='auto', random_state=42, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1, tree_method='exact', validate_parameters=1, verbosity=None)
```

Список преобразований целевых признаков:  
без преобразований.

Список преобразований признаков:  
Стандартизация.

Рисунок 9.4.1 – Вариант описания модели

Если на рабочей области размещены несколько блок схем, при нажатии на описание вы увидите информацию по каждой из них.

## 10. Работа с Дашбордами. Раздел «Визуализация».

*Дашборд* – это интерактивная рабочая область, которая наглядно представляет, визуализирует, объясняет и анализирует данные. На рабочую область пользователь может добавлять графики, таблицы, диаграммы, визуализацию пайплайнов для последующей работы с ними.

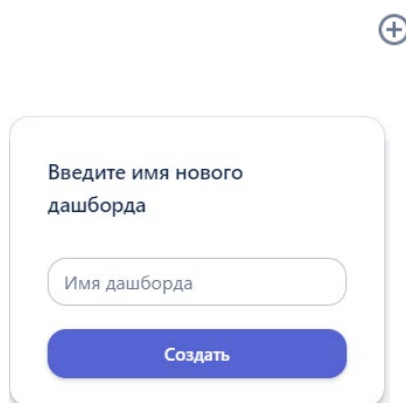


Рисунок 10.1 – Создание нового дашборда

В результате отобразится название дашборда (это текущий дашборд, с которым работает пользователь):

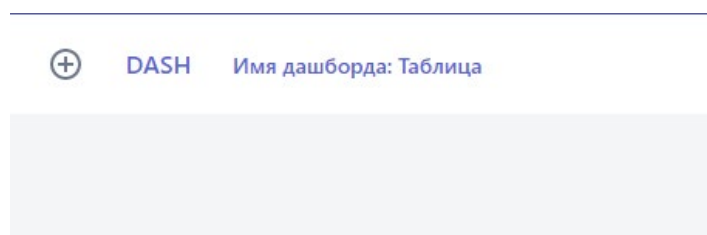


Рисунок 10.2 – Отображение наименования дашборда

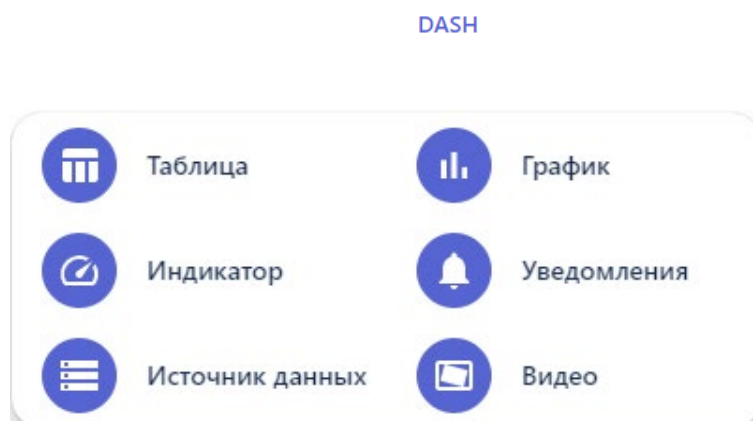


Рисунок 10.3. – Добавление интерактивного блока

В текущей версии Системы реализована работа с блоками: «Таблица», «График» и «Видео». Предварительным условием для добавления любого из типов блоков является создание коннектора в разделе «Соединения», в котором настроено получение данных. При этом это могут быть как данные,

получаемые из внешних источников, так и сгенерированные внутри Системы. Подробнее о создании коннекторов вы можете посмотреть в разделе [Работа со всеми типами коннекторов](#).

## - 10.1. Таблица

Интерактивный блок «Таблица» предназначен для:

- Отображения подключения к внешним базам данных, в виде табличных данных, обновляющихся в режиме реального времени. В таком случае настраивается подключение к коннектору с названиями типов баз данных (clickhouse, postgresql, mongo);
- Записи и сохранения в Системе информации, получаемой из внешних баз данных. Используется коннектор с типом «save\_table»;
- Прогнозирования целевых событий, когда данные для анализа поступают из внешних баз данных. Используется коннектор с типом «table\_app»;
- Отображение таблиц, полученных в результате отработки блок-схем, которые содержат блоки, имеющие в качестве выходной информации визуализации в виде таблиц. Используется коннектор с типом «constructor». Такой тип коннектора создается автоматически, после успешного запуска блок-схемы.

Например, чтобы настроить дашборд с подключением к коннектору с типом «clickhouse»:

1. Создайте новый дашборд.
2. На дашборд добавьте интерактивный блок с типом «Таблица».



### Connectors


Название	Создан	Статус	
<div style="display: flex; align-items: center;"> <div style="width: 15px; height: 15px; background-color: #ffc107; margin-right: 5px;"></div> <div>Коннектор ClickHouse</div> </div>	12.08.2022 14:02	Stopped	<div style="display: flex; gap: 5px;"> <span>▶</span> <span>✎</span> <span>🗑️</span> </div>
set_features_targets_(TRAFFIC)	03.08.2022 16:21	Started	<div style="display: flex; gap: 5px;"> <span>▶</span> <span>✎</span> <span>🗑️</span> </div>
load_data_(TRAFFIC)	03.08.2022 15:45	Started	<div style="display: flex; gap: 5px;"> <span>▶</span> <span>✎</span> <span>🗑️</span> </div>

Рисунок 10.4. – Список коннекторов

В открывшемся окне для коннекторов отображаются их состояния (Запущен (Started) / Остановлен(Stopped)). Если статус коннектора «остановлен», данные из внешних источников не поступают в Систему. Из этого окна можно запустить/остановить коннектор, выполнить его редактирование при необходимости, или удалить его.



5. Выбрать коннектор из списка нажатием на него левой кнопкой мыши.
6. На дашборде отобразится результат подключения к БД «ClickHouse» в виде таблицы с данными:



Data	Y5401	Y5402	Y5707	Y5708	Y5403	Y5404	Y5705	Y5706	P596
2021-09-01	2.701100111	1.108100056	2.03839993	5.59660005	6.69929981	0.02500000	0.66269999	4.93319988	0.885
2021-09-01	2.67519998	1.13929998	2.01699995	5.46570014	6.62709999	0.02390000	0.63190001	4.85069990	0.885
2021-09-01	2.73690009	1.15919995	2.01990008	5.43620014	6.67539978	0.02390000	0.79119998	5.02390003	0.885
2021-09-01	2.36220002	1.17149996	1.86570000	5.13219976	6.21269989	0.02380000	0.74739998	4.57070016	0.875
2021-09-01	2.05369997	1.26590001	1.77119994	4.87739992	5.67309999	0.02360000	0.63870000	4.03620004	0.870
2021-09-01	2.01929998	1.27520000	1.75940001	4.91300010	5.60090017	0.02380000	0.65689998	4.08449983	0.870
2021-09-01	2.00040006	1.31459999	1.75559997	4.91750001	5.59810018	0.02170000	0.67479997	4.09439992	0.860
2021-09-01	2.06110000	1.32720005	1.77139997	4.89720010	5.77670001	0.02170000	0.74140000	4.28620004	0.860
2021-09-01	2.07389998	1.28970003	1.76119995	4.87970018	5.73740005	0.02170000	0.71109998	4.22100019	0.860
2021-09-01	2.10010004	1.27939999	1.77030003	4.88380002	5.83239984	0.02170000	0.69459998	4.13929986	0.870

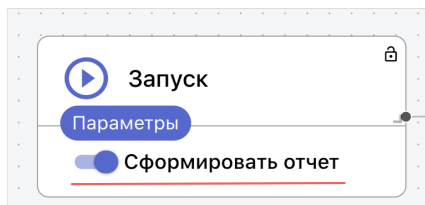
Рисунок 10.5 – Отображение табличных данных из коннектора clickhouse

Коннектор можно остановить или запустить прямо на дашборде, для этого нажмите кнопку «остановить»/«запустить». Для того чтобы удалить дашборд - нажмите на крестик в правом верхнем углу блока.



## 11. Создание отчета с результатами анализа данных

Более подробную информацию с результатами обучения модели можно посмотреть в отчете, который создается также после сборки блок-схемы. Для формирования отчета на первом блоке необходимо перевести бегунок вправо:



В названии отчета будет указано название блок-схемы, по которой формируется отчет и временная отметка его формирования. Для просмотра сформированного отчета после запуска блок-схемы нужно перейти в пункт меню «Отчеты» и выбрать из списка отчет, нажав на его название:

Название	Создан	
ПМИ тесты_2022-09-24_14:48:53.045715_UTC	24.09.2022 17:48	🗑️
ПМИ тесты_2022-09-24_14:48:27.691319_UTC	24.09.2022 17:48	🗑️
ПМИ тесты_2022-09-23_05:24:12.891941_UTC	23.09.2022 08:24	🗑️
NB_2022-09-22_15:12:52.825178_UTC	22.09.2022 18:12	🗑️
NB_2022-09-22_15:11:59.062263_UTC	22.09.2022 18:12	🗑️
NB_2022-09-22_14:10:34.206449_UTC	22.09.2022 17:10	🗑️
NB_2022-09-22_14:10:34.126945_UTC	22.09.2022 17:10	🗑️
NB_2022-09-22_14:10:34.049589_UTC	22.09.2022 17:10	🗑️
NB_2022-09-22_14:10:33.609967_UTC	22.09.2022 17:10	🗑️
NB_2022-09-22_14:10:33.851399_UTC	22.09.2022 17:10	🗑️

Рисунок 11.1 – Список автоматически сформированных отчетов

После нажатия на название отчета открывается отдельная вкладка с отчетом. В отчете кроме результатов обучения созданной модели (см. «Визуализацию») отображаются также: входные данные, отдельно выборки – обучающая и тестовая, датасет после стандартизации признаков, и т.д. Состав отчета отличается в зависимости от метода, который использовался в решении задачи ИИ.

## 12. Конвейер приложений

Конвейер приложений позволяет создать приложение на основе обученной модели. В таком приложении заложен шаблон, умеющий предсказывать наступление интересующих событий. Приложение можно развернуть отдельно за пределами системы, интегрировать с внешними системами, настроить получение входных данных и выполнять прогнозы.

Для того чтобы сформировать приложение необходимо перейти в раздел “Модели”:

Модели

Поиск

Название	Создан			
f1_v2_2023-04-05_12:57:48.321330.UTC	05.04.2023 15:57	+	↓	🗑
f1_v2_2023-04-05_12:57:08.711138.UTC	05.04.2023 15:57	+	↓	🗑
f1_v2_2023-04-05_12:56:22.382544.UTC	05.04.2023 15:56	+	↓	🗑
f1_v2_2023-04-05_12:51:16.073625.UTC	05.04.2023 15:51	+	↓	🗑
f1_v2_2023-04-05_12:34:47.057285.UTC	05.04.2023 15:34	+	↓	🗑
f1_v2_2023-04-05_12:29:17.609176.UTC	05.04.2023 15:29	+	↓	🗑
f1_v2_2023-04-05_10:46:27.303915.UTC	05.04.2023 13:46	+	↓	🗑
f1_v2_2023-04-05_10:19:04.289643.UTC	05.04.2023 13:19	+	↓	🗑
f1_v2_2023-04-05_10:15:47.392897.UTC	05.04.2023 13:15	+	↓	🗑

Рисунок 12.1 – Список сохраненных моделей



Приложения

Имя	Дата изменения	Размер
← Назад		
app_mei-01-auto_2023-04-05_10:59:07.431303.UTC_2023-04-05_11:54:47.320837.UTC.zip	05.04.2023 14:58	2.05 GB ...
app_dbcsan_anomalies_07_03_2023_2023-03-13_07:20:45.807045.UTC_2023-03-14_11:11:30.580853.UTC.zip	14.03.2023 14:23	2.05 GB ...
app_mei-01-auto_2023-03-31_17:53:08.729383.UTC_2023-04-02_13:19:14.541714.UTC.zip	02.04.2023 16:22	2.05 GB ...
app_mei-01-auto_2023-04-06_14:16:02.306016.UTC_2023-04-06_15:47:56.370469.UTC.zip	06.04.2023 18:51	2.05 GB ...
app_mei-01-auto_2023-03-27_09:09:55.595985.UTC_2023-03-27_10:55:39.559300.UTC.zip	27.03.2023 14:02	2.05 GB ...
app_mei-01-auto_2023-03-29_08:43:51.403738.UTC_2023-03-29_11:38:58.100680.UTC.zip	29.03.2023 14:45	2.05 GB ...
app_mei-01-auto_2023-03-29_08:11:37.255075.UTC_2023-03-29_17:13:22.407781.UTC.zip	29.03.2023 20:17	2.05 GB ...
app_mei-01-auto_2023-03-29_08:41:26.279374.UTC_2023-03-29_12:45:29.644548.UTC.zip	29.03.2023 15:48	2.05 GB ...

Рисунок 12.2 – Страница «Приложения»

Приложение упаковано в docker-контейнер и доступно для скачивания. Для того чтобы скачать приложение - нажмите на три точки в строке с названием приложения и скачайте его.

Комплектность приложения после создания и скачивания:

- **app.py** - файл приложения, которое принимает данные, вычисляет и возвращает результат прогнозирования на основе обученной модели
- **const.py** - переменные в приложении.
- **methods.py** - методы, используемые в модели

- **model.py** - загружает модель из файлов `model.pkl`, `model_vars_dict.pkl`, в которых хранится модель и параметры.
- **preprocess\_and\_predict.py** - подготовка и вычисление переменных для прогнозирования.
- **run.sh** - командный интерпретатор для Linux.
- **run.bat** - командный интерпретатор для Windows.
- **requirements.txt** - зависимости для сборки приложения.
- **dockerfile** - файл конфигурации, в котором расписано пошаговое создание среды для работы приложения.
- **docker-compose.yml** - файл с командами для запуска среды приложения.

Приложение позволяет решать задачи предиктивной аналитики для новых данных с использованием обученной модели. Предназначено для использования в сторонних системах.

## 13. Работа с проектом

Сущность «Проект» реализована с целью объединить *группу пользователей* для работы над одним проектом. При этом проект может объединять в себе такие сущности как: «модель», «рабочая область», «дашборд», «отчет», «файл», «коннектор». О том, как добавить каждую из этих сущностей в проект, написано в рамках данного раздела.

### - 13.1. Создание нового проекта

1. Перейдите в пункт меню «Проекты»:

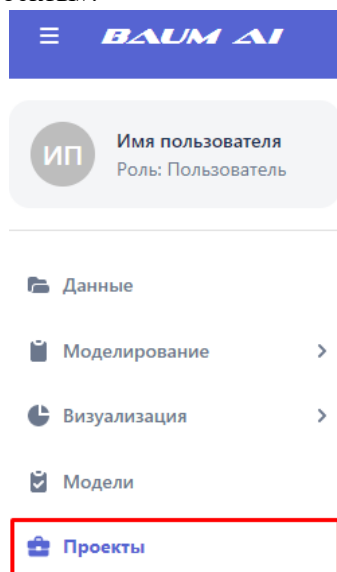


Рисунок 13.1 –Пункт меню «Проекты»

2. Откроется страница «Проекты», на которой отображаются все *проекты*, созданные *пользователем*:

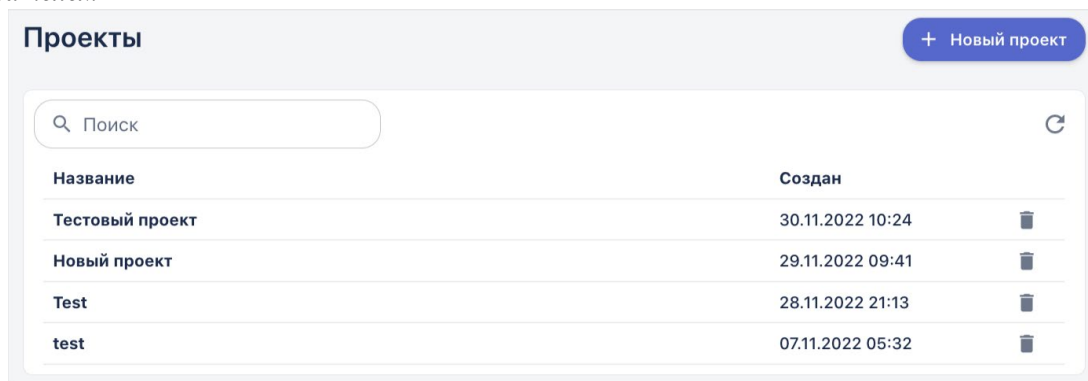


Рисунок 13.2 – Страница с проектами, доступными пользователю

3. Нажмите на кнопку «Новый проект», откроется окно создания нового проекта:

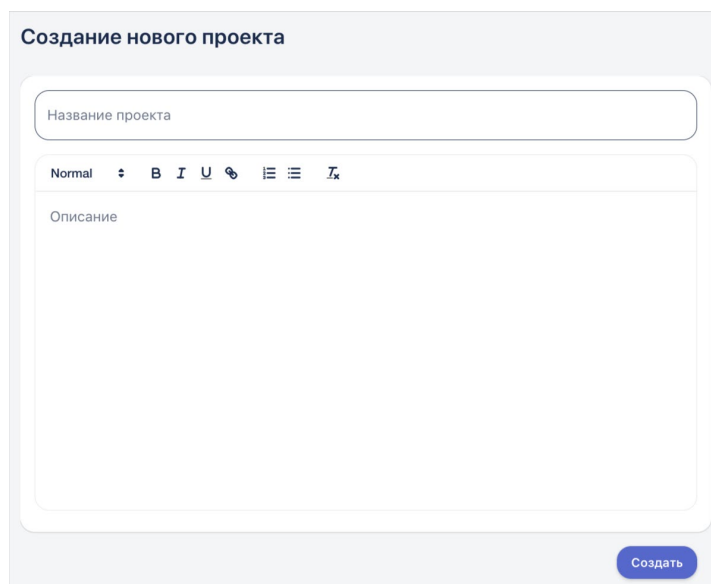


Рисунок 13.3 – Создание проекта

4. Задайте название проекта, например, «Проект тест» (обязательное поле)
5. Задайте описание проекта (необязательное поле), вы можете использовать средства форматирования текста:
  - Заголовки и подзаголовки
  - Жирный текст
  - Курсив
  - Подчеркивание
  - Вставка ссылки
  - Список
  - Нумерованный список
  - Кнопка очистки форматирования
6. Нажмите кнопку «Создать».
7. На страницу «Проекты» добавится новый проект.

## - 13.2. Редактирование проекта

Пользователь имеет возможность отредактировать название проекта и его описание:

1. На странице «Проекты» перейдите в проект, который требует редактирования
2. В открывшемся окне в верхнем правом углу нажмите кнопку «Редактировать» :

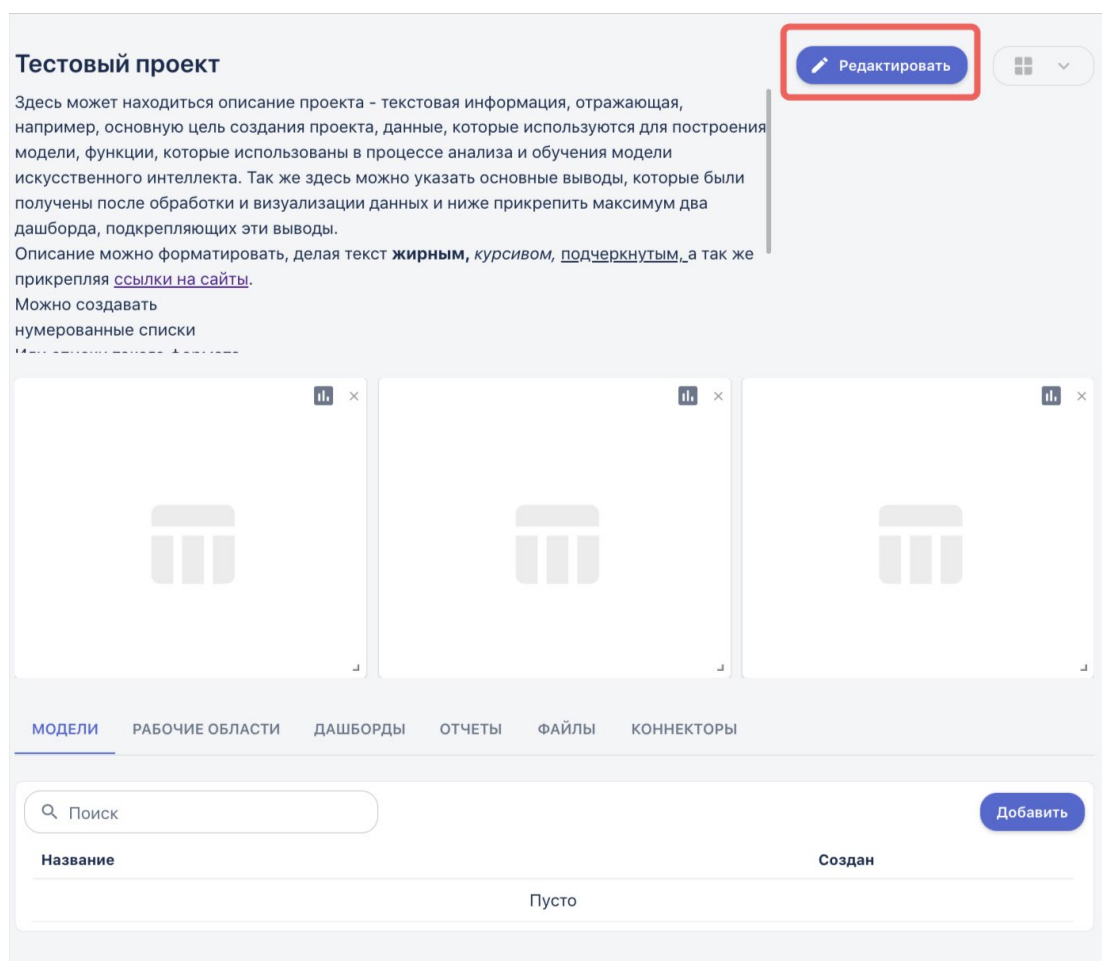


Рисунок 13.4 – Переход к редактированию существующего проекта

3. Задайте новое название проекта и его описание, затем нажмите кнопку «Сохранить»
4. Система вернется на страницу со списком проектов, для того чтобы посмотреть обновленное описание, перейдите в проект, нажав на его наименование.

### - 13.3. Наполнение проекта

После того, как проект создан, его можно наполнить *содержимым* – теми сущностями, над которыми предстоит совместно работать группе пользователей. Для этого:

1. На странице «Проекты» перейдите по ссылке с названием созданного проекта, кликнув на его название:

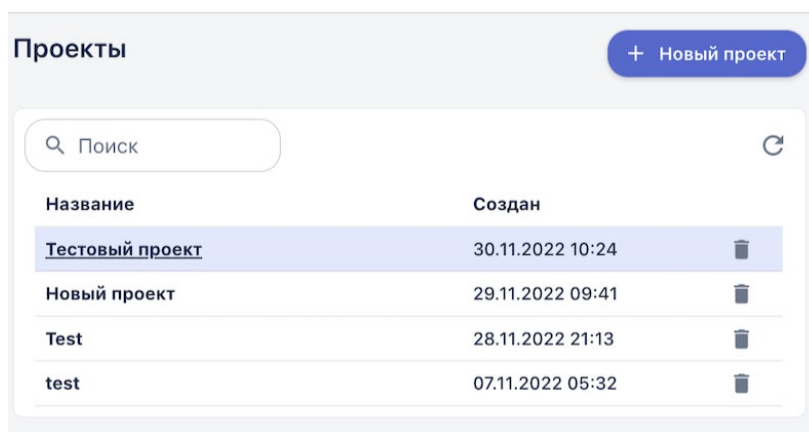


Рисунок 13.5 – Переход на страницу проекта для его просмотра и редактирования

## 2. Откроется *страница проекта* на первой вкладке «Модели»:

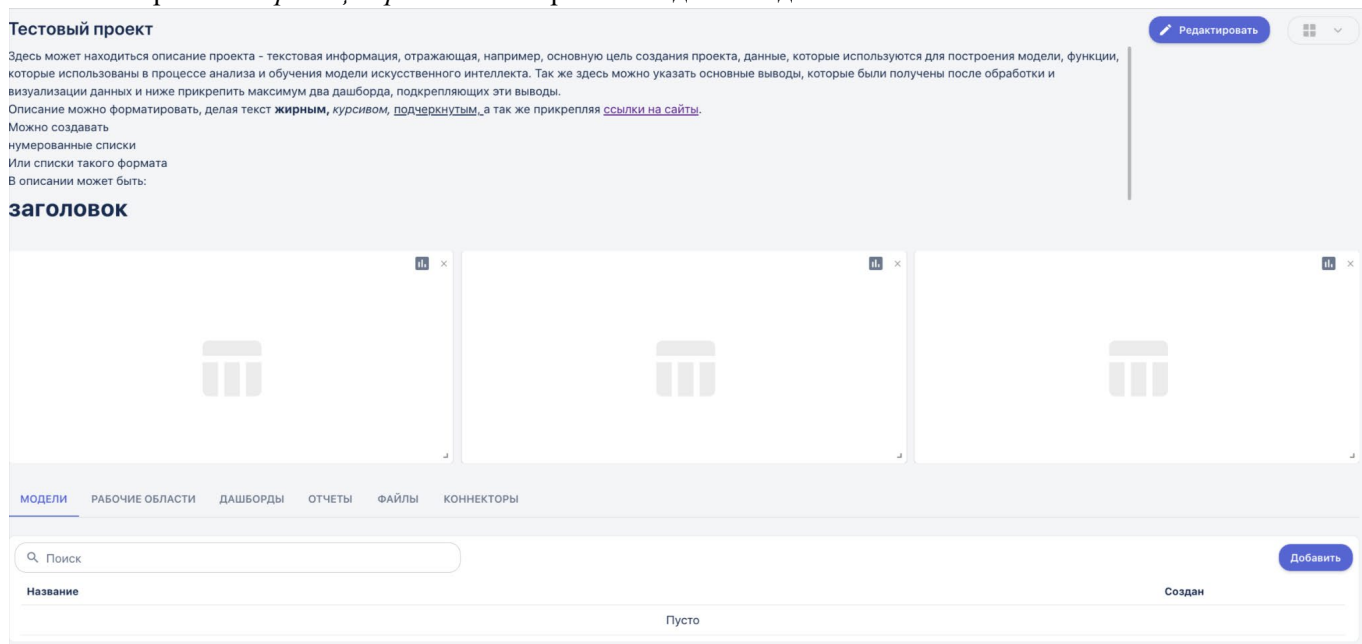
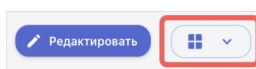


Рисунок 13.6 – Страница проекта

Целиком страница проекта состоит из вкладок с одноименными сущностями: модели, рабочие области, дашборды, отчеты, файлы, коннекторы.

## 3. Работа с дашбордами. Под описанием по умолчанию находятся три пустых дашборда. Вы можете сделать следующее в данном разделе:

- a. Удалить ненужные дашборды, для этого нажмите “x” в правом верхнем углу дашборда.



- b. Добавить новый дашборд. Для этого нажмите на кнопку «Добавить дашборд» и в выпадающем списке выберите нужный тип визуализации:

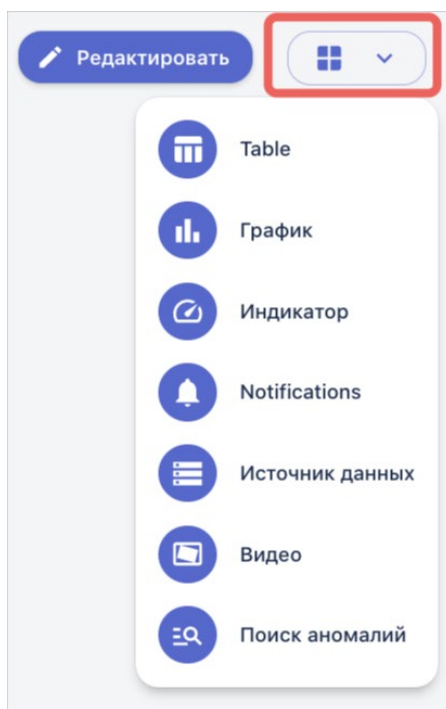


Рисунок 13.7 – Выбор типа визуализации

После этого под описание проекта добавится новый пустой дашборд выбранного типа

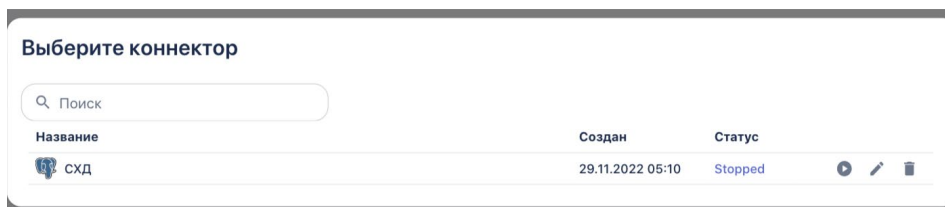


Рисунок 13.8 – Выбор коннектора

После этого визуализация отобразится на дашборде.

- Добавление сущности «Модель».** На вкладке «Модели» пользователь нажимает кнопку «Добавить», и открывается выпадающий список со всеми моделями ИИ, созданными в Системе. Выбирается модель и добавляется в проект.





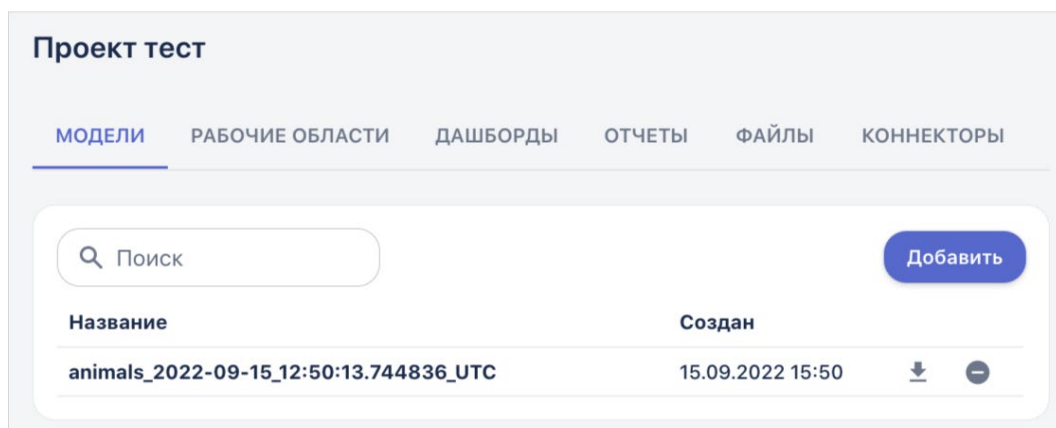


Рисунок 13.9 – Загруженная в проект модель

5. **Добавление сущности «Рабочая область».** Перейдите на вкладку «Рабочие области» и нажмите кнопку «Добавить». Из выпадающего списка выберите рабочую область для добавления в проект:

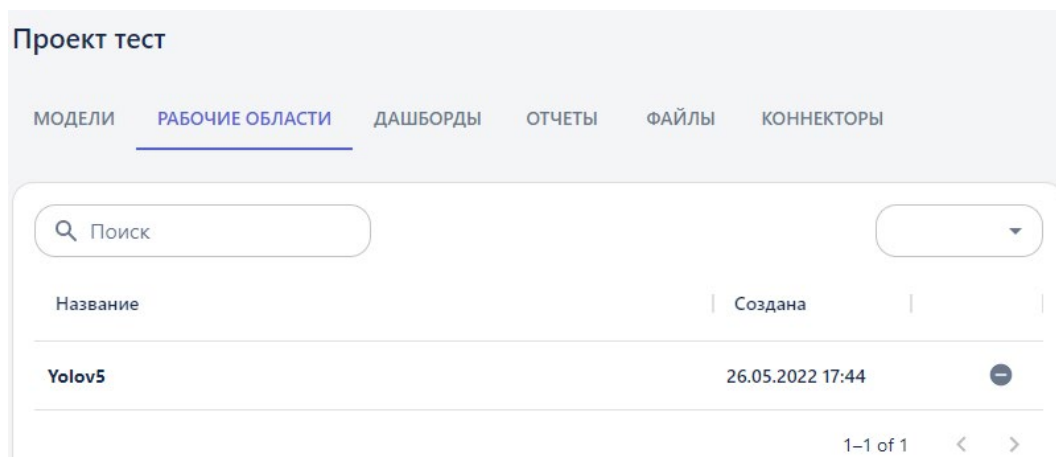


Рисунок 13.10 – Вкладка «Рабочие области»

6. **Добавление сущности «Дашборд».** Перейдите на вкладку «Дашборды» и нажмите кнопку «Добавить», после чего выберите дашборд из выпадающего списка:

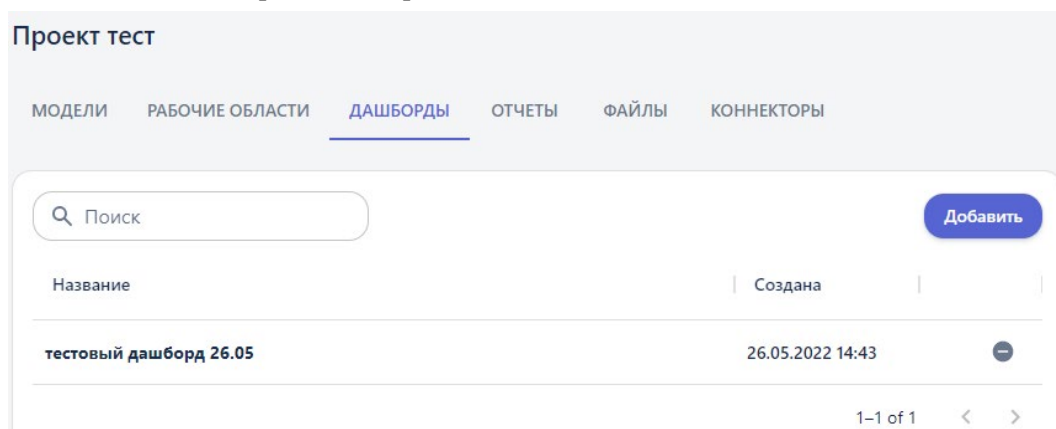


Рисунок 13.11 – Вкладка «Дашборды»

7. **Добавление сущности «Отчет».** Перейдите на вкладку «Отчеты» и нажмите кнопку «Добавить». Из выпадающего списка выберите отчет для добавления в проект.

8. **Добавление сущности «Файл».** Перейдите на вкладку «Файлы» и нажмите кнопку «Добавить». Из выпадающего списка выберите файл для добавления в проект:

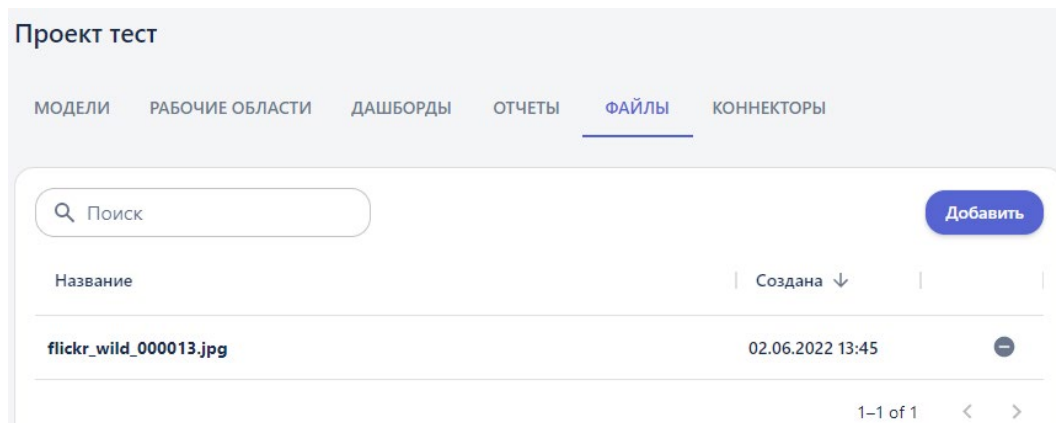


Рисунок 13.12 – Вкладка «Файлы»

9. **Добавление сущности «Коннектор».** Перейдите на вкладку «Коннекторы» и нажмите кнопку «Добавить». Из выпадающего списка выберите коннектор для добавления в проект:

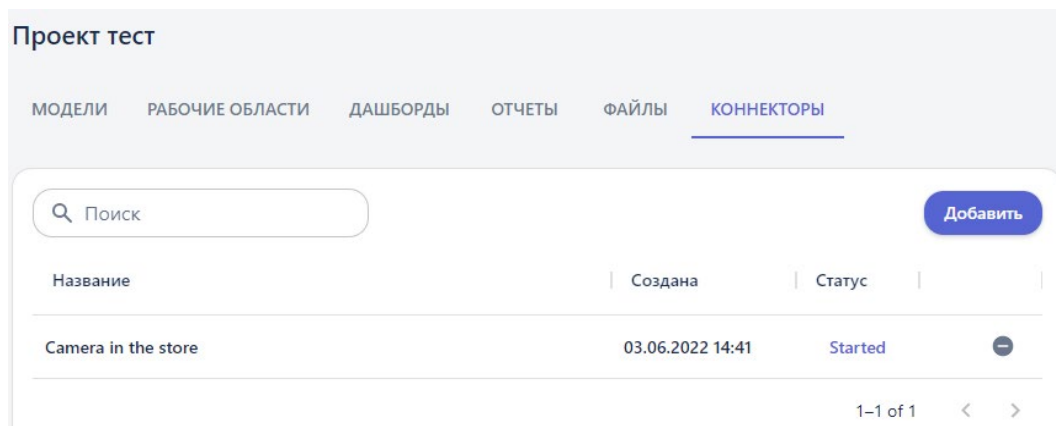


Рисунок 13.13 – Вкладка «Коннекторы»

## 14. Настройка подключения к источникам данных

В платформе «BAUM AI PREDICT» реализована возможность подключения к *внешним системам*, выступающим в качестве *источников данных* для Платформы. При этом данные в режиме реального времени могут поступать из следующих источников: БД (ClickHouse, PostgreSQL, MongoDB) и камеры видеонаблюдения. В данном разделе рассказывается, как настроить в Системе такие подключения. Подробно о визуализации информации из коннекторов на дашбордах написано в разделе «Работа с дашбордами»

Все подключения настраиваются в пункте меню «Соединения», где создаются сущности «Коннектор». *Коннекторы* объединяют в себе источник подключения и запрос на получение данных из него. В данном разделе описаны все типы коннекторов на платформе и сценарии работы с ними на примерах с тестовыми данными.

### - 14.1 Типы коннекторов

#### 1. «ClickHouse»

Данный коннектор предназначен для подключения к БД «ClickHouse». Настройка подключения описана выше в разделе [Настройка подключения на примере ClickHouse](#).

#### 2. «PostgreSQL»

Данный коннектор предназначен для подключения к БД «PostgreSQL». Настраивается по аналогии с коннектором «ClickHouse».

#### 3. «Mongo»

Данный коннектор предназначен для подключения к БД «MongoDB». Настраивается по аналогии с коннектором «ClickHouse».

#### 4. «Table\_app»

На вход коннектора поступают *табличные данные* в онлайн режиме, например из БД «PostgreSQL». Чтобы анализировать входные данные используется обученная в системе *модель* и используется тип коннектора `table_app`.

#### 5. «Save\_table»

Данный коннектор предназначен для сохранения в Системе в виде файлов (на данный момент реализовано сохранение файлов в формате csv) табличных данных, поступающих из сторонних систем. Директория для сохранения файлов в Системе – это раздел «Данные».

В поле «Коннектор» указывается коннектор для подключения к таблице внешнего источника. Устанавливается галочка в поле «Постоянное обновление».

Для такого типа коннектора обязательно указать на выбор:

- Количество строк данных - количество строк из таблицы, которые будут сохранены.
- Интервал - промежуток времени в секундах, в течение которого коннектор должен собирать информацию из таблицы, по истечению этого времени загрузка прекратится и файл будет сохранен.

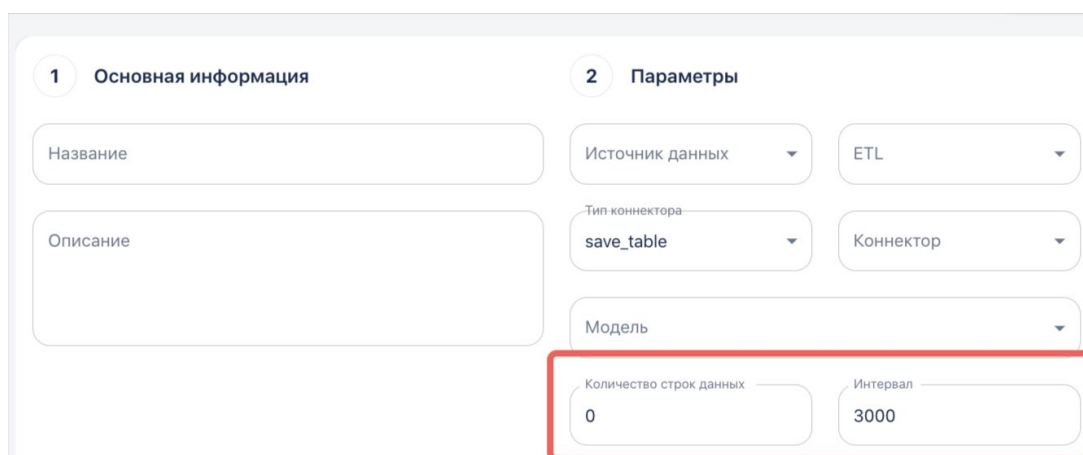


Рисунок 14.2 – Настройка коннектора «save\_table»

## 6. «Classification\_app»

Данный коннектор используется в задаче *классификации изображений*, где на вход коннектора для анализа подается серия из нескольких изображений. Пример настройки коннектора:

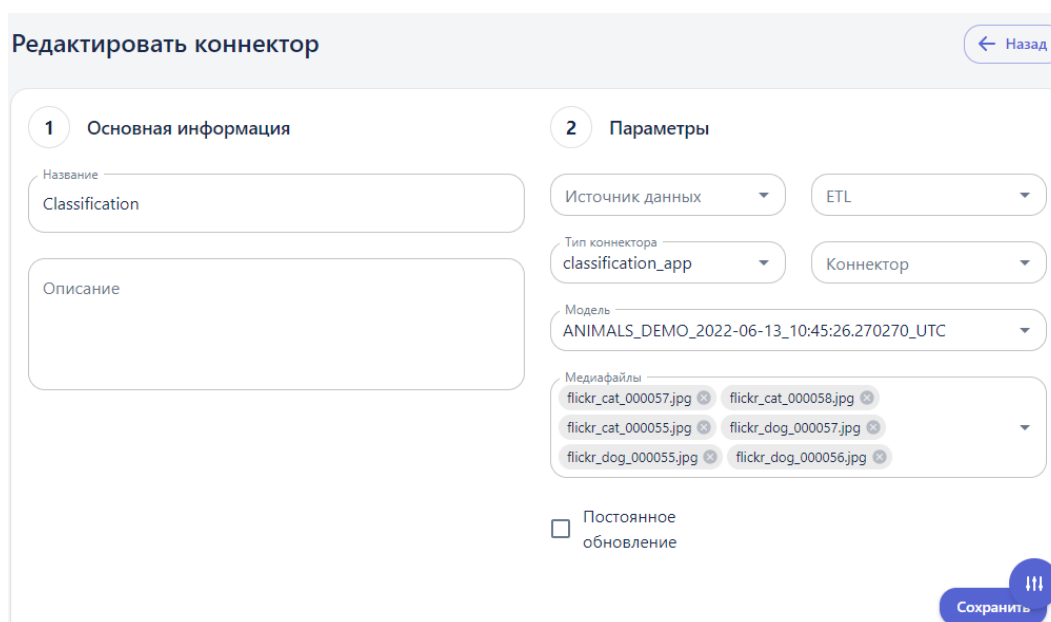


Рисунок 14.4 – Настройка коннектора «Classification\_app»

В поле «Модель» выбирается обученная модель классификации изображений. В поле «Медиафайлы» – серия изображений, которые необходимо классифицировать с использованием обученной модели.

\*Для данного типа коннектора не нужно устанавливать признак «Постоянное обновление», так как анализируются данные, загружаемые с локального устройства.

## 7. «Constructor» (автоматически создаваемый)

Условием создания коннектора является следующее: на Платформе создается блок-схема, где один из элементов имеет на выходе *визуализацию* – выходным параметром элемента является *таблица*, *график* и т.д. Пользователь запускает такую блок-схему, и после успешной отработки элемента с визуализацией в Системе создается коннектор. Число создаваемых коннекторов при запуске блок-схемы соответствует числу элементов с визуализацией на этой блок-схеме.

Название коннектора формируется из названия элемента и названия рабочей области, и должно являться уникальным в рамках Системы. Пример – *yolov5\_train\_(yolov5\_noses\_eyes)*, где название рабочей области указано в скобках.

Такой коннектор автоматически создается в статусе «Started» – пользователь не должен запускать коннектор, и может сразу же перейти к просмотру данных коннектора в окне дашборда.

## - 14.2 Порядок работы с коннекторами

Предварительно коннектор должен быть создан и запущен (за исключением служебных коннекторов, которые создаются и запускаются автоматически). Только после этого выполняется подключение к нему через окно дашборда.

### ■ 14.2.1 Создание коннектора

Создание коннекторов осуществляется в разделе меню «Соединения». Для коннекторов, в которых настраивается подключение к внешним источникам данных (к внешним базам данных, к камере видеонаблюдения), дополнительно создаются сущности – *источник данных*, и *ETL*. Для остальных типов, эти сущности не создаются, а сразу создается сущность «Коннектор».

Для того чтобы создать новый коннектор, перейдите на вкладку «Коннекторы» и нажмите на кнопку «Создать коннектор»:

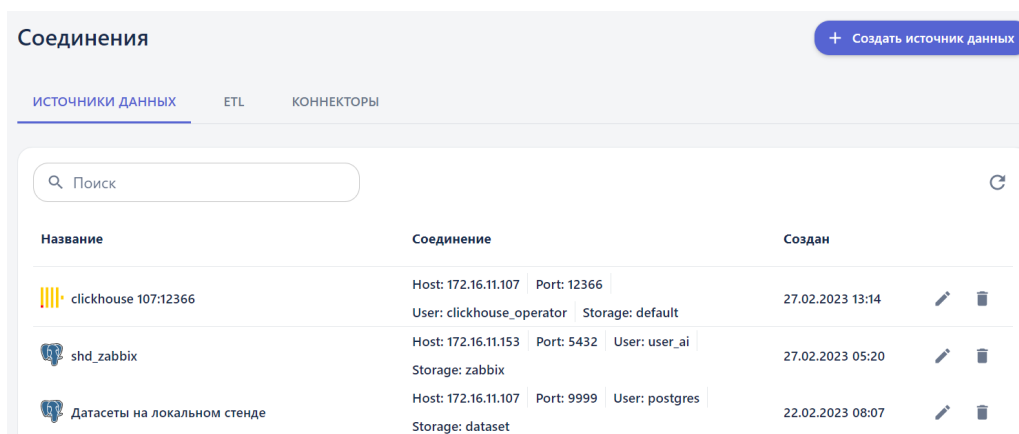


Рисунок 14.5 – Переход к созданию коннектора на вкладке «Коннекторы»

В открывшейся форме выберите тип коннектора в одноименном поле «Тип коннектора», и заполните поля:

Создать новый коннектор

← Назад

1 Основная информация

2 Параметры

Название

Описание

Источник данных: ETL

Тип коннектора: Коннектор

Модель

Медиафайлы

Постоянное обновление

Создать

Рисунок 14.6 – Форма создания нового коннектора

После заполнения формы нажмите кнопку «Создать».

#### ■ 14.2.2 Запуск коннектора



Соединения

+ Создать коннектор

ИСТОЧНИКИ ДАННЫХ ETL КОННЕКТОРЫ

Search

Название	Создан	Статус	
predict_overload_(Прогнозирование загрузки СХД)	14.04.2023 09:11	Started	⏸️ ✎️ 🗑️
linear_regression_(Прогнозирование загрузки СХД)	14.04.2023 09:11	Started	⏸️ ✎️ 🗑️
vis_overload_realtime_(Прогнозирование загрузки СХД) 1	12.04.2023 11:13	Started	⏸️ ✎️ 🗑️
shd_zabbix	27.02.2023 05:22	Stopped	▶️ ✎️ 🗑️

Рисунок 14.7 – Запуск коннектора

В результате коннектору присваивается статус «Started», и он готов к визуализации на дашборде.

#### ■ 14.2.3 Подключение к коннектору на дашборде



Откроется модальное окно со списком коннекторов для подключения.

## - 14.3 Настройка подключения на примере ClickHouse

Действия:

1. Перейдите в пункт меню «Соединения»:

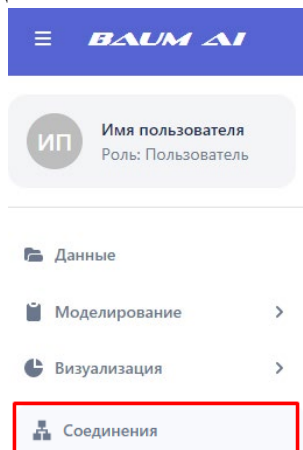


Рисунок 14.8 – Пункт меню Соединения

Откроется страница «Соединения» на первой вкладке «Источники данных», на которой отображаются все ранее созданные источники:

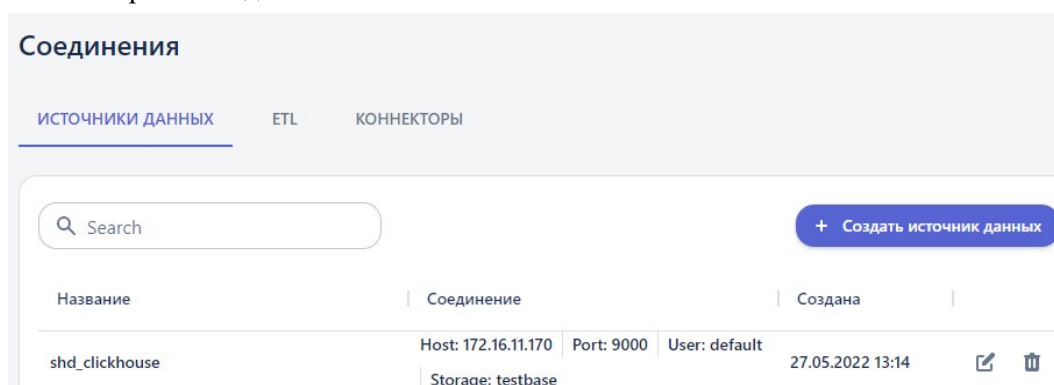


Рисунок 14.9 – Вкладка «Источники данных»

2. **Создание нового источника.** Нажмите на кнопку «Создать источник данных». Откроется окно «Создание нового источника данных»:

The screenshot shows the 'Создание нового источника данных' (Create new data source) form. It has a 'Назад' (Back) button at the top right. The form is divided into two sections: '1 Основная Информация' (Basic Information) and '2 Параметры' (Parameters). The 'Основная Информация' section contains fields for 'Название' (Name) and 'Описание' (Description). The 'Параметры' section contains fields for 'Хост' (Host), 'Порт' (Port), 'Имя хранилища' (Storage name), 'Тип хранилища' (Storage type), 'Имя пользователя' (Username), and 'Пароль' (Password). A blue 'Создать' (Create) button is located at the bottom right.

Рисунок 14.10 – Окно настройки источника данных

Заполните поля:

- *Название.* Пользователь задает название источника «TEST DATA SOURCE (clickhouse)», к которому будет настраиваться подключение.
- *Хост.* Указывается хост протокола TCP/IP, т.е. сетевой интерфейс устройства, предоставляющего сервис формата «клиент-сервер», где сервером выступает БД **ClickHouse**, а клиентом – платформа **BAUM AI PREDICT**. По сути это IP-адрес подключаемой БД. Необходимо указать «172.16.11.116».
- *Порт* – номер порта, по которому устанавливается соединение с сервером, на котором установлена БД **ClickHouse**. Указать «12366».
- *Имя хранилища* – название базы данных, которое указано на подключаемом сервере. Указать «default».
- *Тип хранилища.* Из выпадающего списка необходимо выбрать тип «clickhouse»:

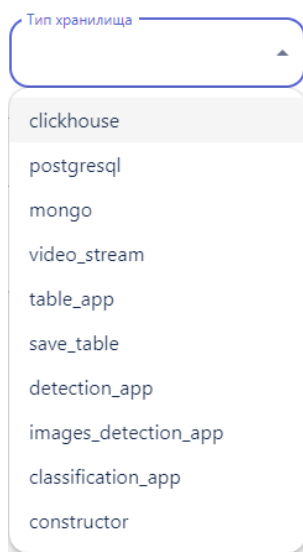


Рисунок 14.11 – Список возможных типов хранилища

Описание всех типов коннекторов представлено в разделе [Классификация коннекторов](#)

- *Имя пользователя, пароль* – параметры учетной записи администратора внешнего сервера для разрешения доступа к данным. Указать пользователя «clickhouse\_operator», и пароль для него «clickhouse\_operator\_password».
- *Описание.* Вводится дополнительная информация по источнику, необязательное поле.

Для регистрации в Системе источника нажмите кнопку «Создать».

3. **Создание нового ETL.** Сущность «ETL» (дословно Extract, Transform, Load – с англ. извлечение, преобразование загрузка) содержит в себе sql запрос для извлечения данных из источника. То есть в шаге 2 создается источник и прописывается запрос для извлечения данных из него.

- 3.1. Перейдите на вкладку «ETL»:



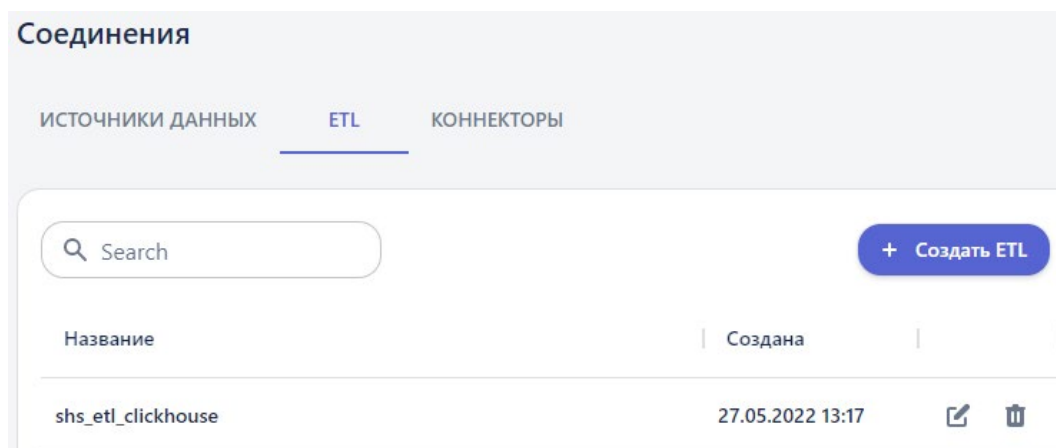


Рисунок 14.12 – Вкладка ETL

3.2. Нажмите на кнопку «Создать ETL». Откроется окно «Создать новый ETL»:

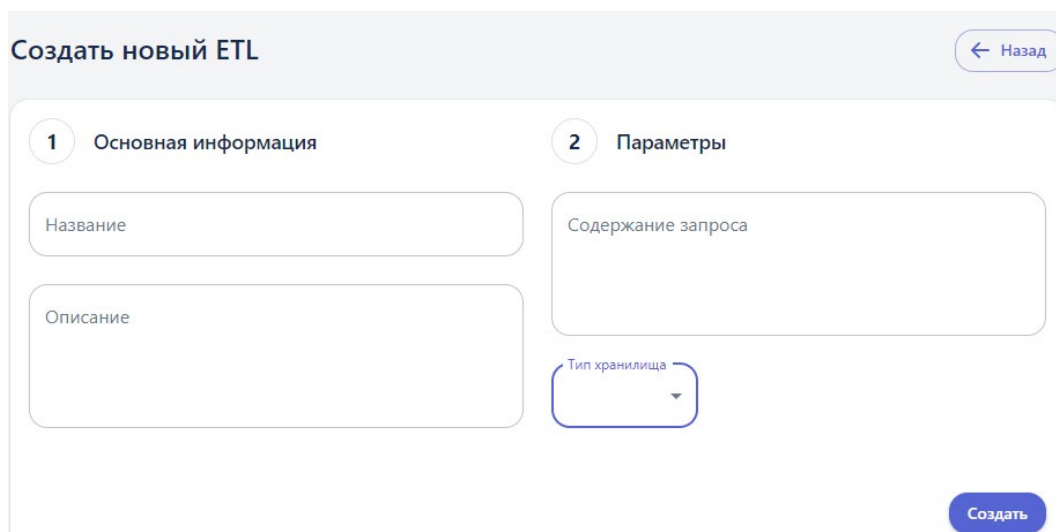


Рисунок 14.13 – Окно настройки ETL

3.3. Заполните поля:

- *Название.* Пользователь вручную задает название ETL «TEST ETL (clickhouse)» – запрос на извлечение данных.
- *Содержание запроса.* Прописывается непосредственно sql запрос для извлечения данных из внешнего сервера. При этом указывается название таблицы, из которой данные извлекаются, в запросе «SELECT \* FROM stock». Чтобы извлечь данные только из первых ста строк этой таблицы используется запрос «SELECT \* FROM stock LIMIT 100».
- *Тип хранилища.* Выбирается тип «clickhouse».
- *Описание* (необязательное поле).

3.4. Нажмите кнопку «Создать».

#### 4. Создание нового коннектора:

4.1. Перейдите на вкладку «Коннекторы»:

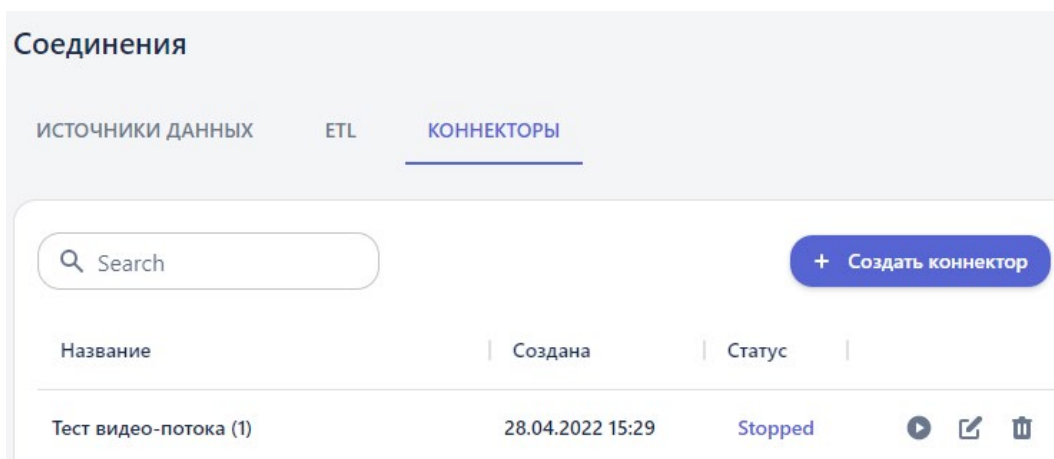


Рисунок 14.14 – Вкладка коннекторов

- 4.2. Нажмите на кнопку «Создать коннектор». Откроется окно «Создать новый коннектор»:

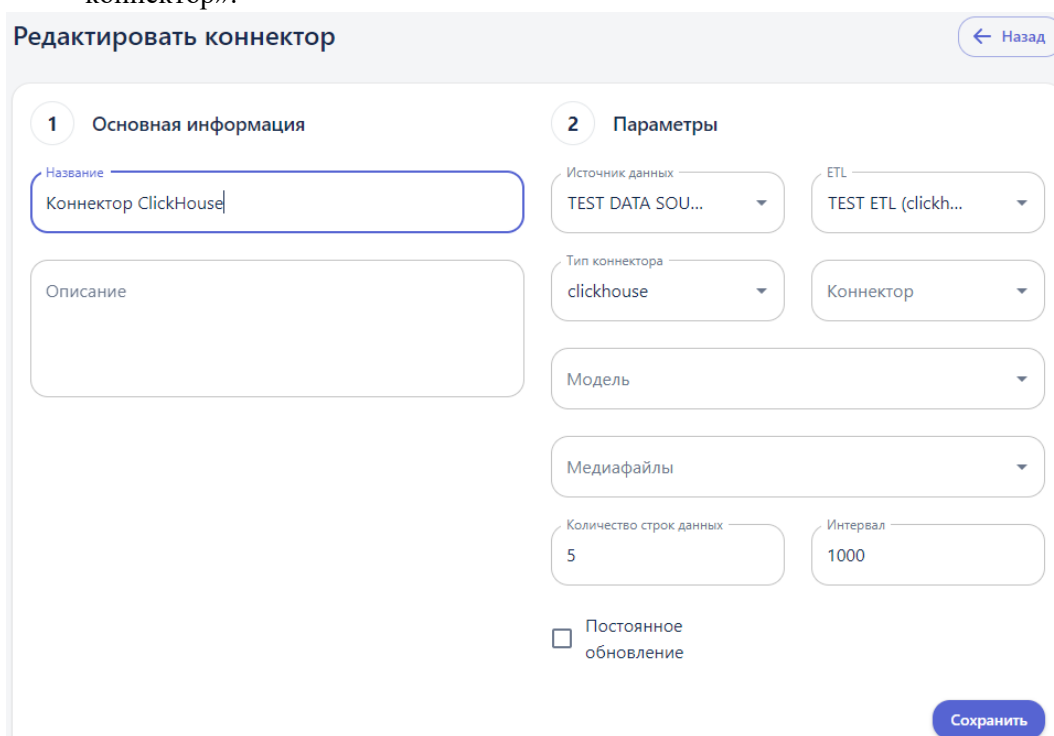


Рисунок 14.15 – Окно настройки коннектора

- 4.3. Заполните поля:

- *Название*. Пользователь вручную задает название создаваемого коннектора «Коннектор ClickHouse».
- *Источник данных*. Из списка выбирается источник «TEST DATA SOURCE (clickhouse)», созданный в шаге 2.
- *ETL*. Из списка выбирается ETL «TEST ETL (clickhouse)», созданный в шаге 3.
- *Тип коннектора*. По умолчанию выбирается первый тип из списка – clickhouse, оставить выбранное значение.
- *Количество строк данных*. Данные из внешней БД поступают порциями, по указанному или меньшему количеству строк за раз.

- *Интервал* – периодичность, с которой выполняются запросы во внешнюю БД. Заполняется числовым значением (в миллисекундах), или значением в формате «Дата» (*второй вариант не реализован в текущей версии*). Если задать интервал 1000 мс, и указать количество строк пять, каждую секунду будет запрашиваться пять записей.
- *Постоянное обновление*. Признак устанавливается, когда данные ожидаются бесконечно. Если признак не установить, запрос данных завершится, при получении их в полном объеме.
- *Описание*.

Остальные поля на форме создания коннектора для текущего сценария не заполняются, они используются при создании других типов коннекторов. Нажмите кнопку «Создать».

4.4. Сразу после создания коннектору присваивается статус «Stopped»:

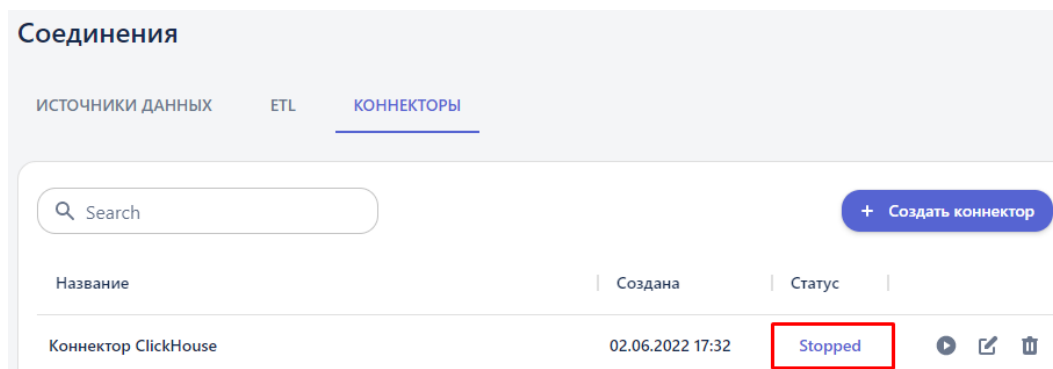


Рисунок 14.16 – Статус запуска коннектора



Вы так же можете отображать табличные данные ClickHouse в режиме реального времени, используя для визуализации сущность «Дашборд» и подключаясь к запущенному коннектору.

## 15. Примеры рабочих областей

### - 15.1 Вводная часть

-

#### 15.1.1. Описание проблемы

В работе оператора систем хранения данных возникает необходимость отслеживания состояния заполнения томов данных. При приближении заполнения тома к задаваемому заранее пороговому значению, оператор должен принимать решение об удалении ненужных данных или о вводе в строй новых томов для сохранения информации.

■

#### 15.1.2. Исходные данные

##### 15.1.2.1. Формат данных.

■ Формат входных данных алгоритма - поток данных с СХД по SNMP 123 .

Формат: Дата, Заполнение тома, Том, Время в сек (с 00:00 01.01.1970 года). Ключевые поля подчеркнуты

##### 15.1.2.2. Пример данных из СХД

```
2022-06-02 20:15:00,0.9748873805999756,Logicalused by volume: pool1/vol1,1654200900
2022-06-02 20:20:00,0.974887731552124,Logicalused by volume: pool1/vol1,1654201200
2022-06-02 20:25:00,0.9748881130218505,Logicalused by volume: pool1/vol1,1654201500
2022-06-02 20:30:00,0.9748884716033936,Logicalused by volume: pool1/vol1,1654201800
2022-06-02 20:35:00,0.9748888301849366,Logicalused by volume: pool1/vol1,1654202100
2022-06-02 20:40:00,0.9748892040252686,Logicalused by volume: pool1/vol1,1654202400
2022-06-02 20:45:00,0.9748895473480225,Logicalused by volume: pool1/vol1,1654202700
```

#### 15.1.3. Постановка задачи

Необходимо на базе платформы реализовать механизм для автоматизации процесса отслеживания состояния заполнения тома и уведомления оператора в случае заполнения тома выше определенного значения. При этом должны быть соблюдены следующие условия:

- 1) Для решения задачи отсутствует необходимость владения навыками программирования для обучения модели классификации. Специалист должен иметь возможность использования предобученной модели в режиме no-code.
- 2) Должен быть предусмотрен интерфейс для доступа к СХД через коннектор, который позволит получать данные в автоматическом режиме.

**Этапы решения проблемы:**

### Этап 1

1. Создать коннектор для подключения к Zabbix в режиме реального времени
2. Создать рабочую область
3. Создать дашборд для отображения результата

### Результат

1. Создан коннектор для подключения к Zabbix “shd\_zabbix”;
2. Создана рабочая область “СХД”;
3. Создан дашборд “СХД\_zabbix”.
4. Получен датасет логирования (буфер)

## - 15.2 Практическая часть

### 15.1.1. Функциональные требования

Сервер с доступом в Интернет.

### 15.1.2. Проект

Для удобства организации и работы с сохраняемыми объектами, на платформе можно создать проект, к которому прикрепляются все относящиеся к нему объекты. При нажатии на название проекта откроется окно, в котором можно работать над сущностями внутри проекта: модели, рабочие области, дашборды, отчеты, файлы и коннекторы (см. рис. 2-1). Для добавления сущности в проект необходимо сначала создать её в соответствующем разделе системы.

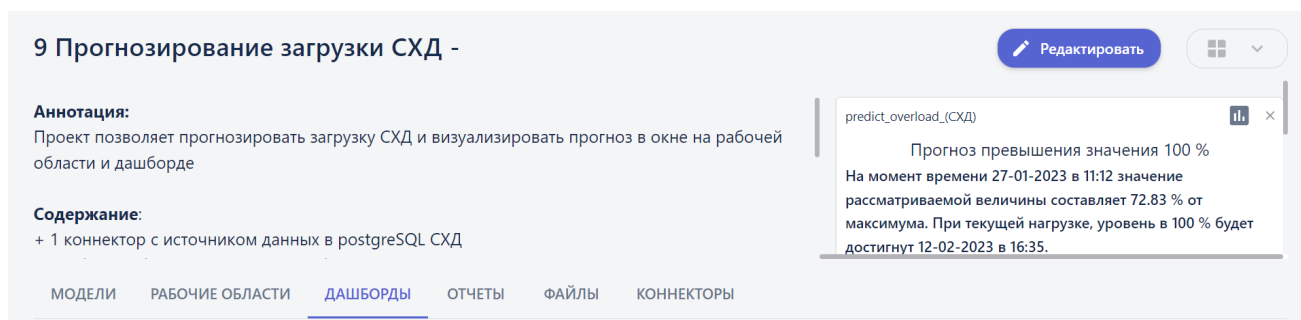


Рисунок 2-1 - Проект и его сущности

### 15.1.2.1. Рабочие области

#### 15.1.2.1.1. Пайплайн для прогнозирования загрузки

Пайплайн предназначен для обучения модели прогнозировать нагрузку на СХД на основе анализа данных, принятых из коннектора, применяя функцию линейной регрессии.

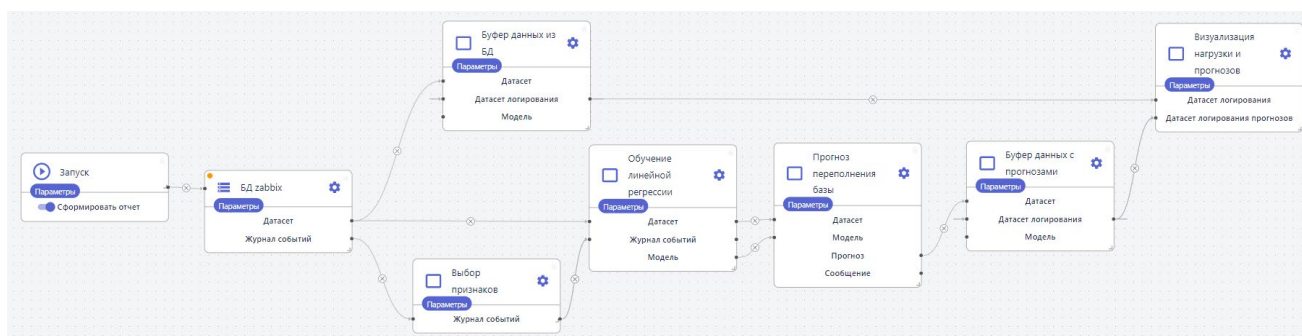


Рисунок 2.1.1-1 - Пайплайн для создания модели

Подробное описание блоков находится в Базе Знаний раздел “Моделирование”: [База знаний BAUM AI](#).

Параметры блоков:

- Блок “Запуск”. Блок является начальной точкой для пайплайна.

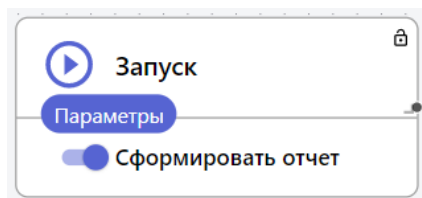


Рисунок 2.1.1-2 - Блок “Запуск”

**Комментарий:**

Кнопка “Сформировать отчет” появится после первого нажатия на запуск. Если активировать данную опцию, будет сформирован отчет об обучении модели предиктивной аналитики.

- Блок “Загрузка табличных данных их коннектора”. Блок реализует загрузку данных из коннектора **shz\_zabbix**. При выборе опции “Сохранить датасет” - обязательно выбрать ниже директорию для сохранения (в данном случае **shd\_test.csv**).

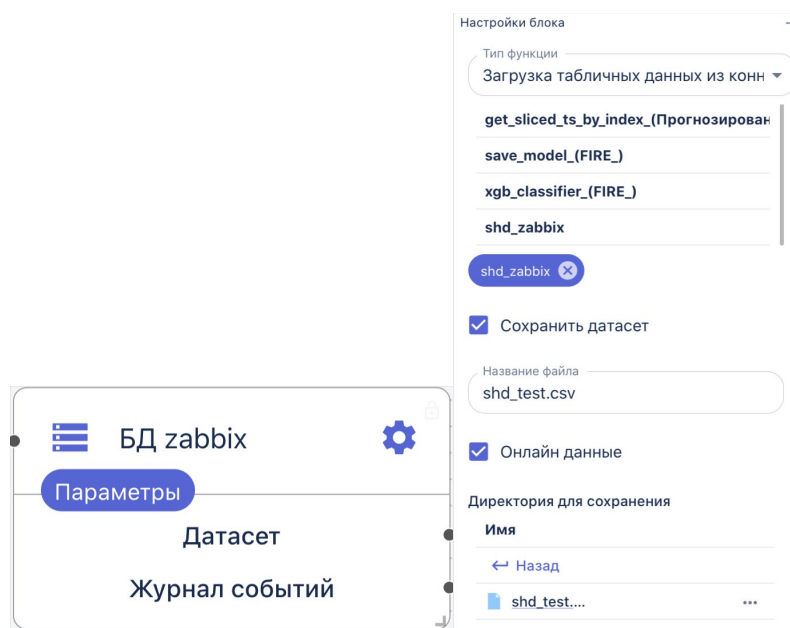


Рисунок 2.1.1-3 - Блок “Загрузка табличных данных их коннектора”

- Блок “**Выбор признаков**”. В блоке происходит выбор целевого признака (значение, которое предстоит научиться предсказывать модели) и признака (параметр, который будет исследоваться для выявления корреляция между ним и рассматриваемым целевым признаком).

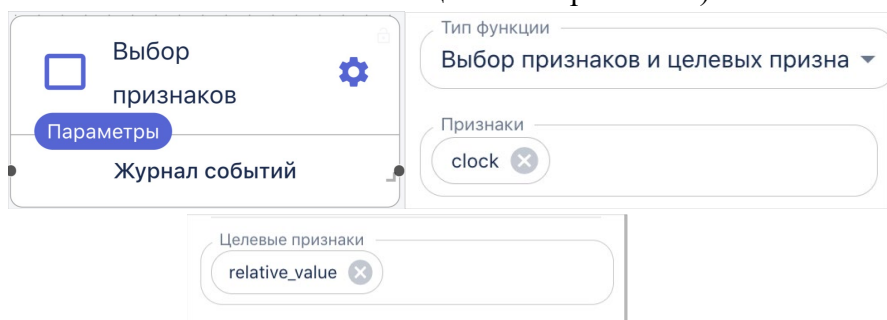


Рисунок 2.1.1-4 - Блок “Выбор признаков”

- Блок “**Обучение линейной регрессии**”. Блок прогнозирует целевую переменную (признак) на основании независимой переменной (признака).

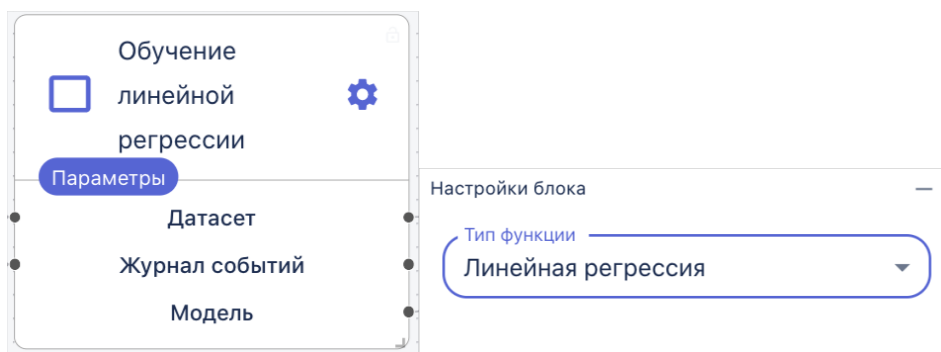


Рисунок 2.1.1-5 - Блок “Обучение линейной регрессии”

- Блок **“Буфер данных из БД”**. В данном блоке будет осуществляться логирование новой поступающей информации (Датасет) из коннектора в “растущий” буфер “Датасет логирования”.

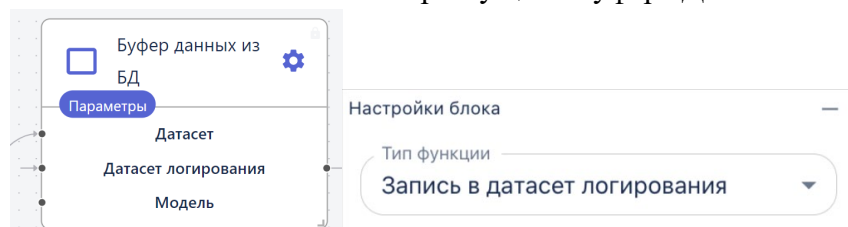


Рисунок 2.1.1-6 - Блок “Буфер данных из БД”

- Блок **“Прогноз переполнения базы”**. Блок будет осуществлять прогнозирование времени переполнения объема памяти с учетом заданного уровня, где 1 = 100%, то есть можно спрогнозировать, когда память заполнится полностью.

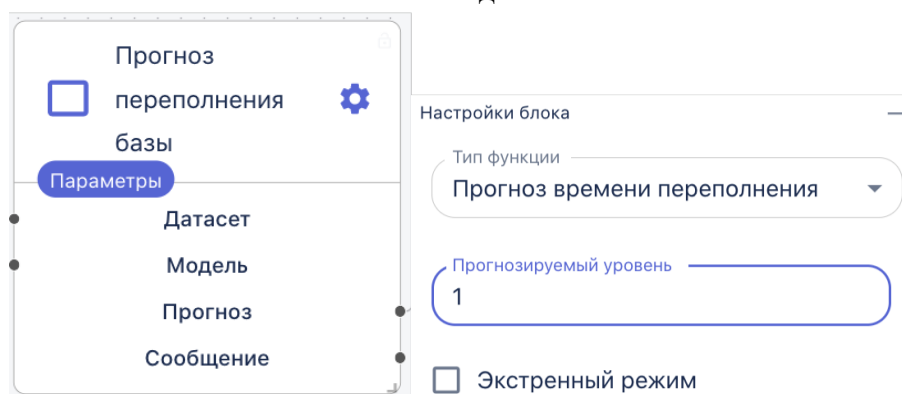


Рисунок 2.1.1-7 - Блок “Прогноз переполнения базы”

- Блок **“Буфер данных с прогнозами”**. В данном блоке будет



осуществляться логирование прогнозных значений в датасет.

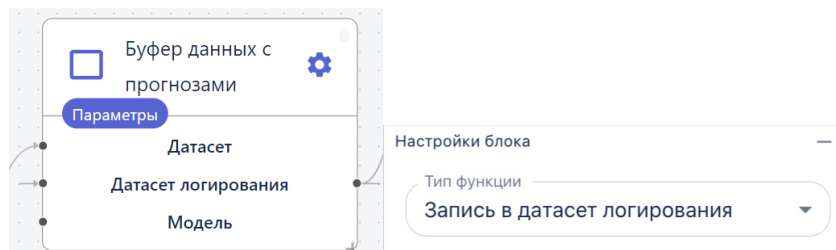


Рисунок 2.1.1-8 - Блок “Буфер данных с прогнозами”

- Блок “Визуализация нагрузки и прогнозов”. Блок позволяет создать визуализацию с помощью графика, на котором будут отображаться фактические и прогнозные значения заполнения памяти.

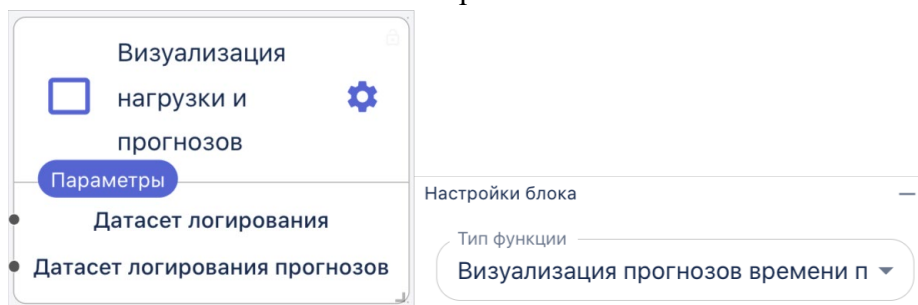


Рисунок 2.1.1-9 - Блок “Визуализация нагрузки и прогнозов”

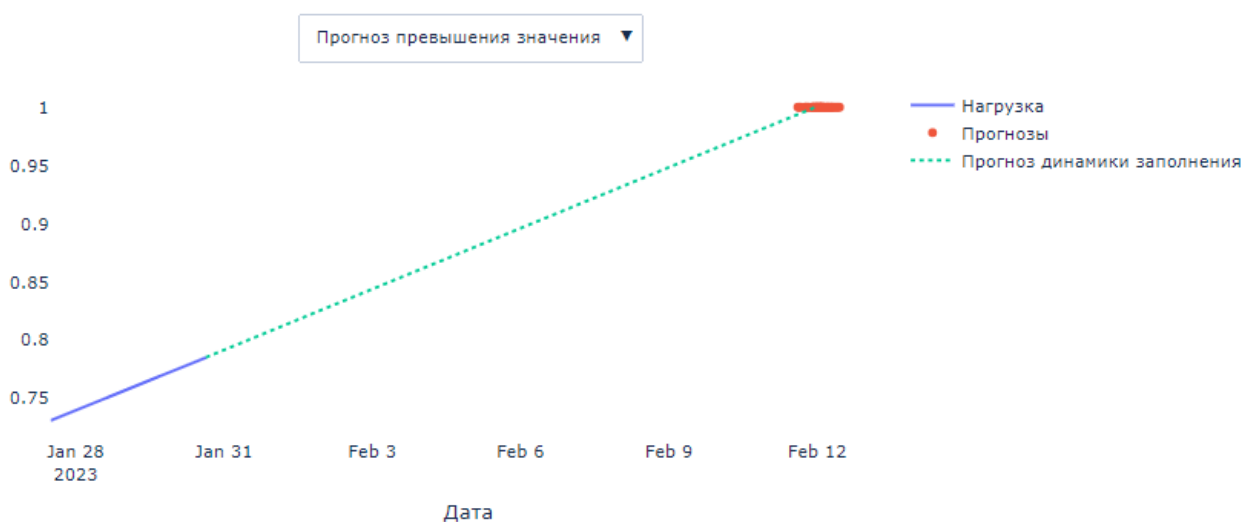


Рисунок 2.1.1-10 - Блок “Визуализация нагрузки и прогнозов”. График “Прогноз превышения значения”

15.1.2.1.2. Визуализация результатов

После того как все элементы схемы будут успешно обработаны на панели инструментов появляются кнопки (см. рис. 2.1.2-1):



Рисунок 2.1.2-1 - Меню результатов работы

Последовательно нажимая на кнопки визуализации, можно отобразить на рабочей области следующую информацию:

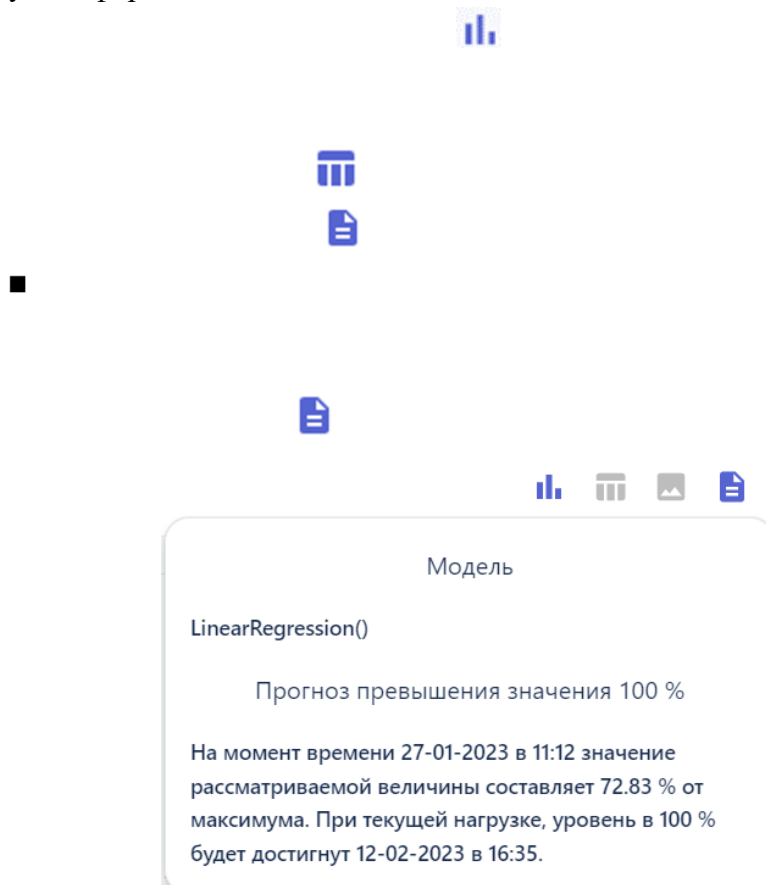


Рисунок 2.1.2-2 - Автогенерируемое описание

Более подробный отчет, включающий в себя дополнительно визуализацию результатов обучения модели, информацию по датасету и валидацию на тестовой выборке можно посмотреть в разделе “**Отчеты**”.

### 15.1.2.2. Дашборды

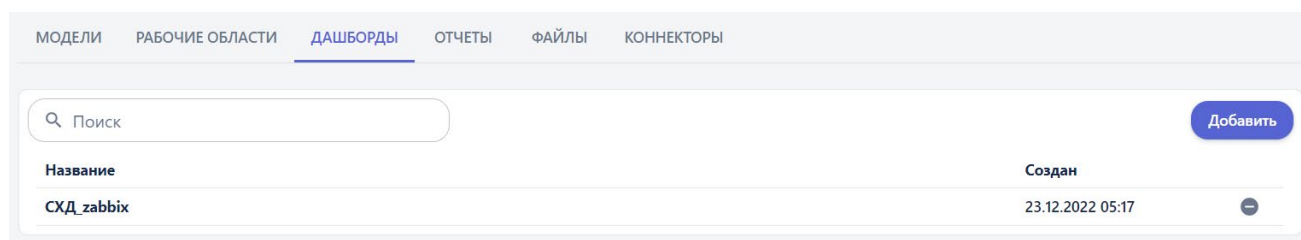


Рисунок 2.2-1 - Дашборд в проекте

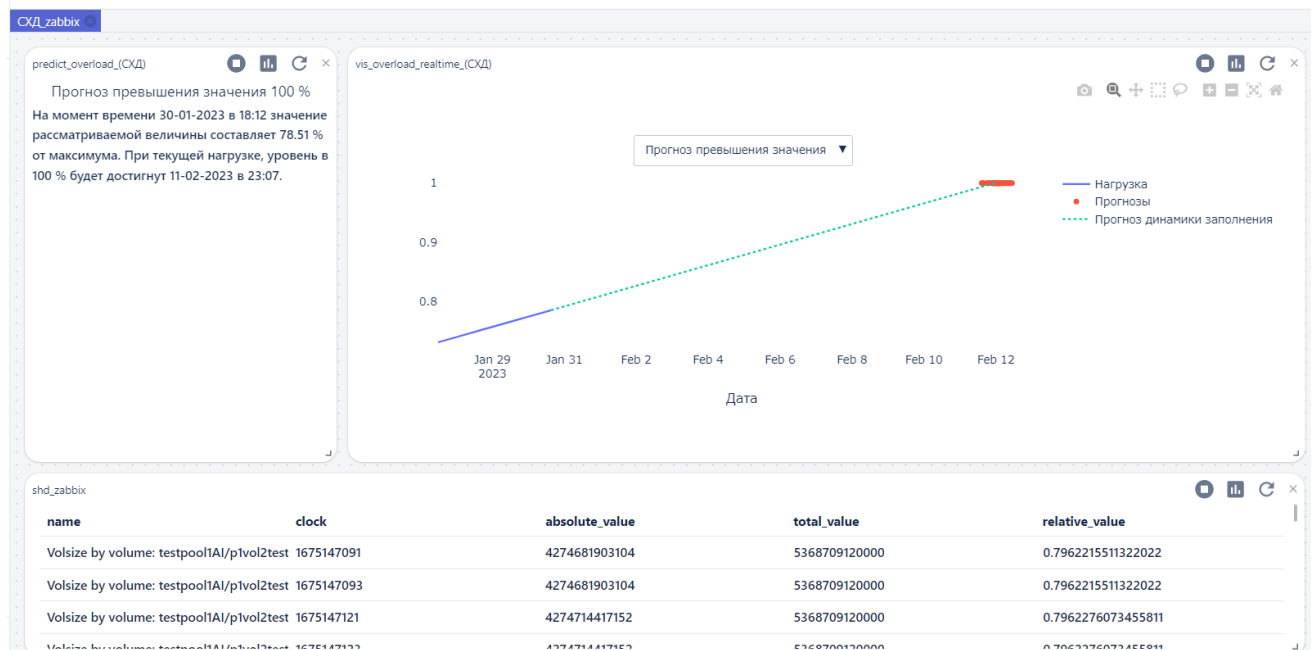


Рисунок 2.2-1 - Отображение информации на дашборде

### 15.1.2.3. Файлы

МОДЕЛИ	РАБОЧИЕ ОБЛАСТИ	ДАШБОРДЫ	ОТЧЕТЫ	ФАЙЛЫ	КОННЕКТОРЫ
Поиск					
<a href="#">Добавить</a>					
Название					Создан
shd_test.csv					26.01.2023 19:27

Рисунок 2.4-1 - Сохраняемый датасет в проекте

### 15.1.2.4. Коннекторы

Коннектор предназначен для создания соединения с внешней базой данных и позволяет получать информацию напрямую из системы по заданному интервалу. Соединение состоит из трех последовательно настраиваемых элементов: источник данных, ETL и коннектор.

#### 15.1.2.5.1. Источник данных

Здесь настраивается соединение с сервером с указанием его параметров, а именно IP-адрес хоста, порт, имя хранилища, тип хранилища, имя пользователя и пароль (скрыт на рис. 2.5.1-1).

Редактировать источник данных

← Назад

**1** Основная Информация

Название  
shd\_zabbix

Описание

**2** Параметры

Хост  
172.16.11.153

Порт  
5432

Имя хранилища  
zabbix

Тип хранилища  
postgresql

Имя пользователя  
login

Пароль  
pass

Сохранить

Рисунок 2.5.1-1 - Настройки “Источника данных”

#### 15.1.2.5.2. ETL

Здесь прописывается SQL запрос для извлечения данных из источника, и указывается тип хранилища.

Редактировать источник данных

← Назад

**1** Основная Информация

Название  
shd\_zabbix\_etL

Описание

**2** Параметры

Содержание запроса  
WITH total\_size AS  
(  
SELECT  
history\_uint.clock,  
history\_uint.value AS total\_value,  
items.name  
FROM history\_uint  
JOIN items ON (history\_uint.itemid = items.itemid)

Тип хранилища  
postgre...

Сохранить

Рисунок 2.5.2-1 - Настройки “ETL”

Содержание запроса (название таблиц в запросе может изменяться в зависимости от названия таблиц у заказчика):

```
WITH total_size AS  
(  
SELECT  
    history_uint.clock,  
    history_uint.value AS total_value,  
    items.name  
FROM history_uint  
JOIN items ON (history_uint.itemid = items.itemid)
```

```
AND items.snmp_oid LIKE '.1.3.6.1.4.1.91919191.1.3.6.4.2.1.8.%'  
WHERE name LIKE '%testpool1AI/p1vol2testAI'  
LIMIT 1  
)  
  
SELECT  
    total_size.name,  
    used_size.clock,  
    used_size.absolute_value,  
    total_size.total_value,  
    absolute_value / total_size.total_value AS relative_value  
FROM total_size, (  
    SELECT  
        history_uint.clock,  
        history_uint.value AS absolute_value,  
        items.name  
    FROM history_uint  
    JOIN items ON (history_uint.itemid = items.itemid)  
    AND items.snmp_oid LIKE '.1.3.6.1.4.1.91919191.1.3.6.4.2.1.9.%'  
    WHERE name LIKE '%testpool1AI/p1vol2testAI' AND clock > 1668508000) AS  
used_size  
ORDER BY clock ASC
```

#### 15.1.2.5.3. Коннектор

В коннекторе указывается уже настроенный ранее источник соединения данных и ETL и указывается тип хранилища.

**1 Основная информация**

Название: shd\_zabbix

Описание:

**2 Параметры**

Источник данных: shd\_zabbix

ETL: shd\_zabbix\_etl

Тип коннектора: postgresql

Коннектор:

Модель:

Количество строк данных: 0

Интервал: 3600000

Постоянное обновление

Рисунок 2.5.3-1 - Настройки “Коннектора”

ИСТОЧНИКИ ДАННЫХ    ETL    **КОННЕКТОРЫ**

Search

Название	Создан	Статус
shd_zabbix	14.12.2022 14:05	Started

Рисунок 2.5.3-2 - Коннектора shd\_zabbix

### 15.1.3. Результаты

На базе платформы создана рабочая область с пайплайном и дашбордом для отслеживания состояния загрузки СХД и превентивного принятия решения оператором.

### 15.1.4. Сохраняемые объекты

1. Коннектор;
2. Рабочая область;
3. Дашборд;
4. Датасет логирования.



## 16. Администрирование Платформы

### - 16.1 Пользователи и группы

Платформа позволяет разделять уровни доступа к разделам меню для разных пользователей в зависимости от требования проекта. Для этого создаются Группы, которые впоследствии назначаются отдельным пользователям.

Создание новой Группы или редактирование уже созданной группы осуществляется в разделе Администрирование -> Группы:

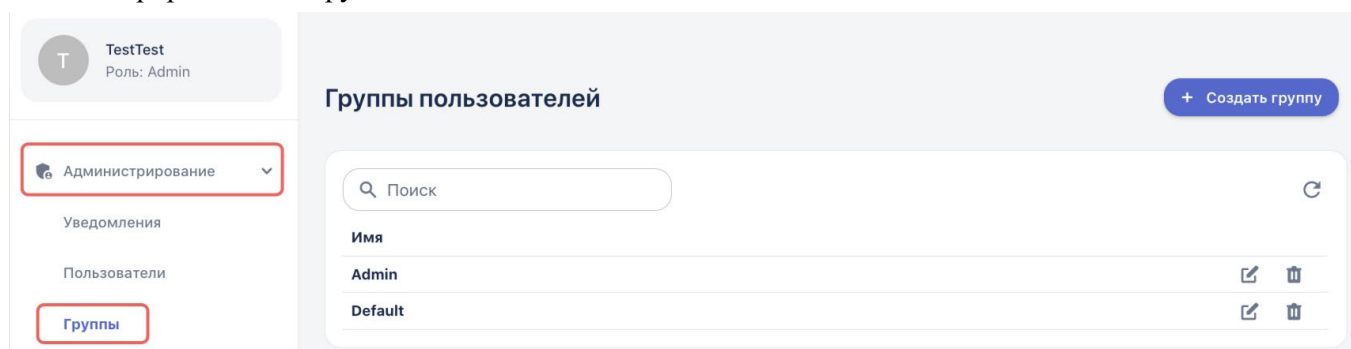


Рисунок 16.1 – Раздел Группы меню Администрирование

В списке отображаются уже созданные группы. Для создания новой группы, нажимается кнопка «Создать группу» в правом верхнем угле, в открывшемся окне отобразятся все доступные для настройки параметров:

Рисунок 16.2 – Настройка новой группы пользователей

Сначала заполняется основная информация:

- *Название* - задать уникальное название для Группы (обязательное поле)



- *Описание* - внести краткую информацию о применимости данной роли (необязательное поле)

Далее из списка доступов выбираются разделы данных, к которым у пользователей, принадлежащих к этой Группе, должны быть доступны:

- admin - Администрирования
- application - Приложения
- api - API
- connection - Соединения
- data - Данные
- modeling - Моделирование
- model - Модели
- project - Проекты
- report - Отчеты
- visual - Визуализация

После того, как галочки для соответствующих разделов проставлены, нажимается кнопка «Создать». Новая группа появится в списке.

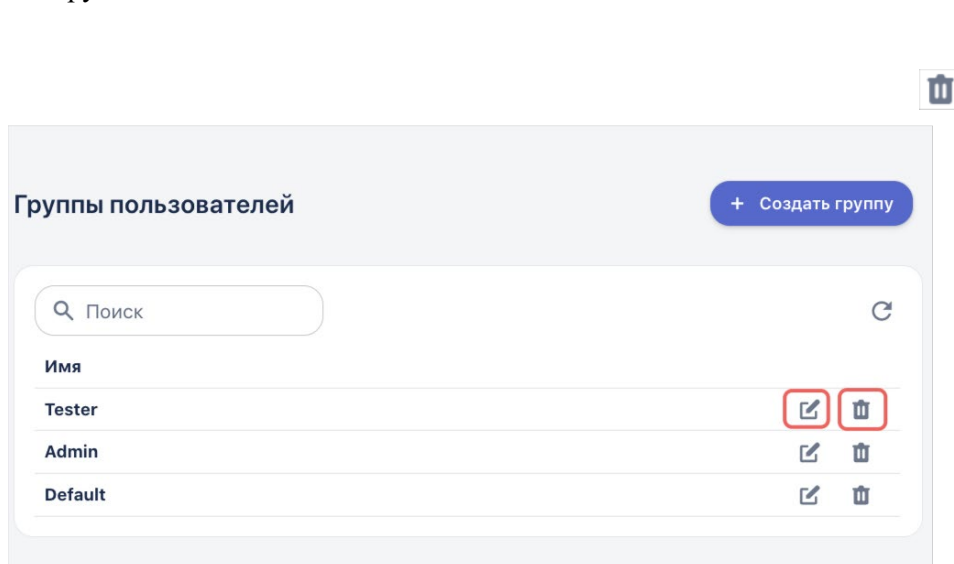


Рисунок 16.3 – Кнопки удаления и редактирования группы пользователей

Список всех пользователей представлен в разделе Администрирование -> Пользователи:

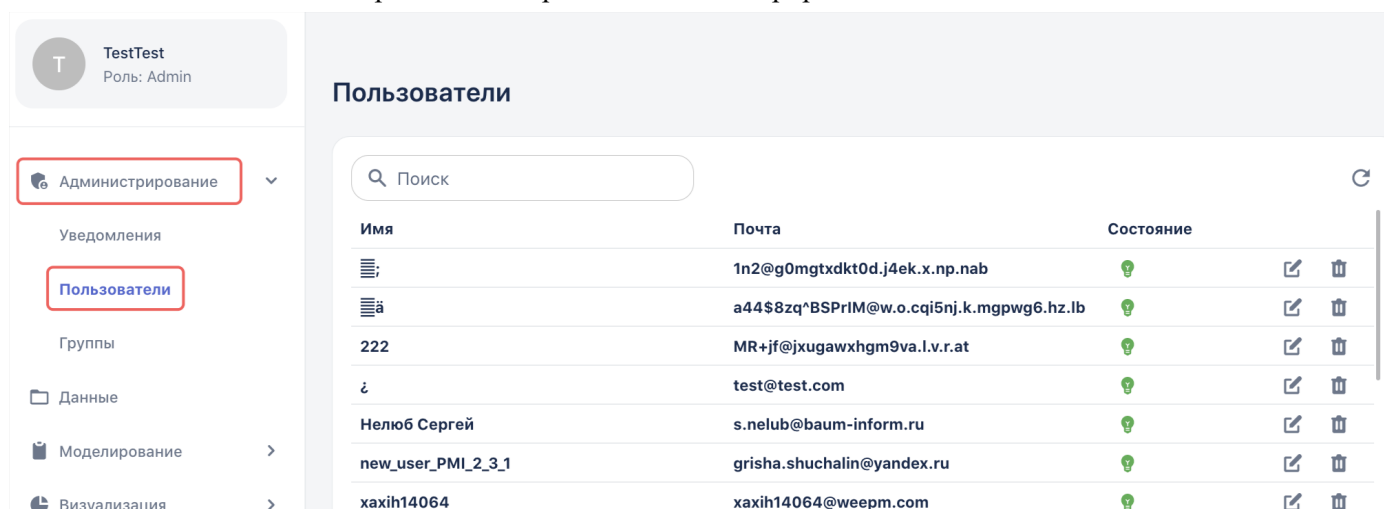


Рисунок 16.4 – Список пользователей в меню Администрирование

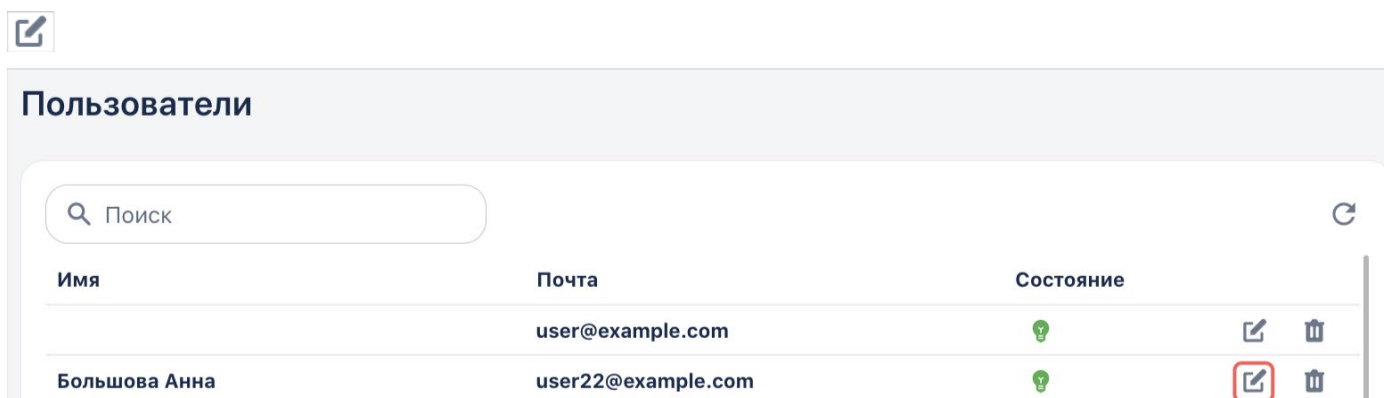


Рисунок 16.5 – Редактирование данных пользователя

В открывшемся окне есть возможность изменить имя пользователя, адрес электронной почты, а также назначить Группу:

**Редактирование пользователя**

Имя  
Olga Shpak

Эл. почта \*  
o.guryanova@baumlab.pro

Группа пользователей  
Default

Сохранить

Рисунок 16.6 – Параметры для редактирования

Обратите внимание, что всем новым пользователям по умолчанию присваивается Группа Default (группа по умолчанию).



Рекомендуется оставлять одного-двух пользователей Администраторов, которые смогут управлять доступами и отвечать за настройки.

## - 16.2 Настройка отправки уведомлений

На платформе реализована возможность отправки уведомлений - сообщений в телеграм или на почту. Этот функционал позволяет пользователям получать автоматические оповещения о результатах работы пайплайнов, при выполнении заданных условий. Создание и настройка каналов осуществляется следующим образом:

1. Перейти в раздел Администрирование -> Уведомления:

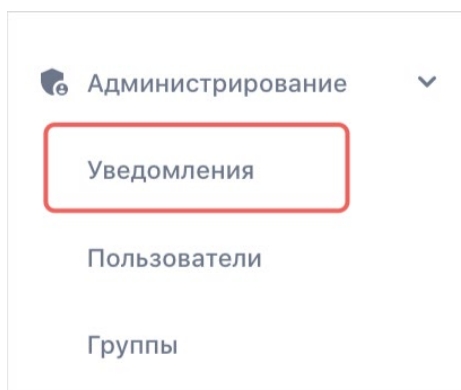


Рисунок 16.6 - Раздел Уведомления

2. В открывшемся окне отобразится список всех существующих каналов уведомлений:

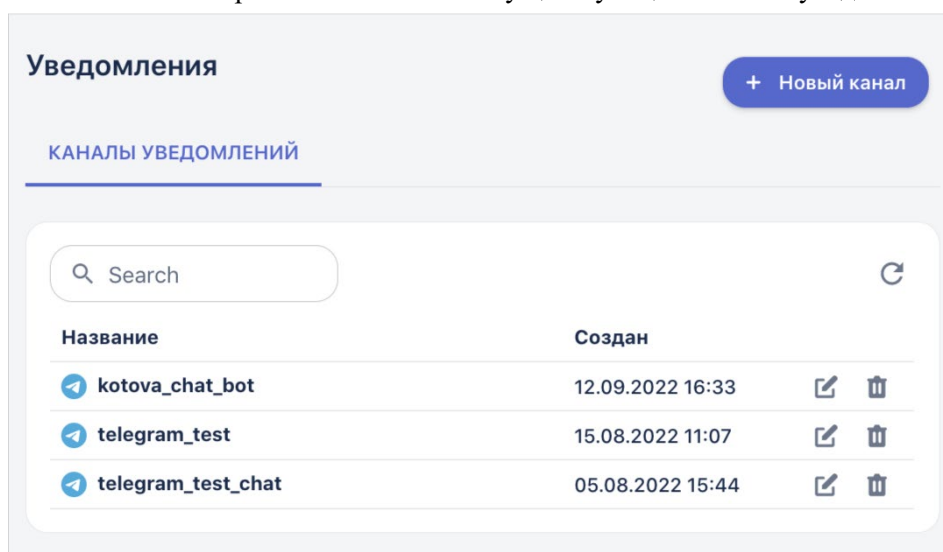


Рисунок 16.6 - Список созданных каналов

3. Для создания нового канала уведомлений нажимается кнопка «Новый канал» в правом верхнем углу
4. Открывается окно создания нового канала уведомлений:

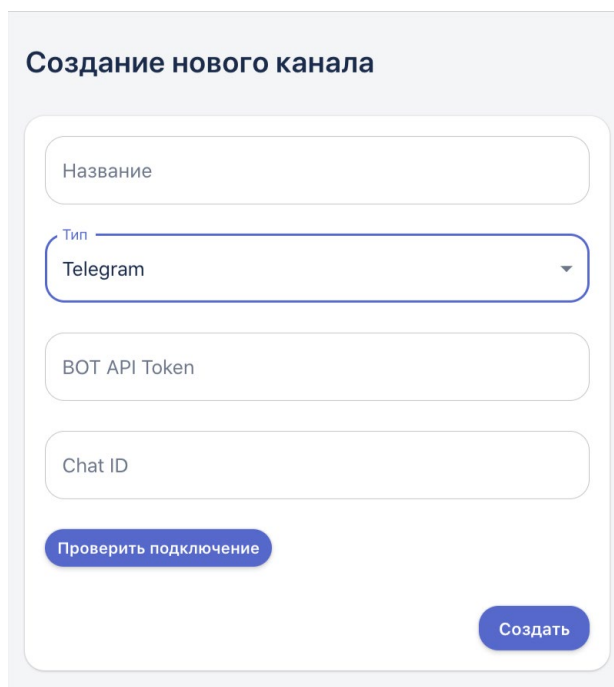


Рисунок 16.7 - Окно создания нового канала

5. Уведомления можно отправлять в телеграм бот или на почту.
  - 5.1. Для создания уведомления с использованием **телеграм бота** поля заполняются следующим образом:
    - *Название* - пользователь задаёт название бота из телеграмм
    - *Тип* - Telegram
    - *BOT API Token* - токен бота, полученный в телеграм при создании бота
    - *Chat ID* - уникальный численный идентификатор чат бота

**Примечание:** сначала пользователь должен создать чат бота в телеграм или получить его

token и id для настройки канала уведомлений. Создание бота описано

- 5.2. Для создания уведомления **на почту** поля заполняются следующим образом:
  - *Название* - пользователь задает название канала
  - *Тип* - Почта
  - *Почты* - указываются адреса, на которые должны быть отправлены уведомления. При этом после ввода первого адреса, нужно нажать Enter, потом перейти к указанию следующего адреса и т.д.

**Примечание:** уведомления на почту не реализованы в текущей версии системы

6. После того, как все поля заполнены, нажимается кнопка «Проверить подключение», если все настройки были указаны верно, система выдаст сообщение об успешности подключения в верхнем правом углу окна. В противном случае, отобразится сообщение «Не удалось подключиться»
7. Завершающим этапом нажимается кнопка «Создать»
8. Новый канал отобразится в списке.

Для того чтобы начать получать уведомления в телеграм, пользователь должен создать пайплайн, который будет содержать блок «Процесс» с функцией «Отправка уведомлений», где указывается необходимый канал:

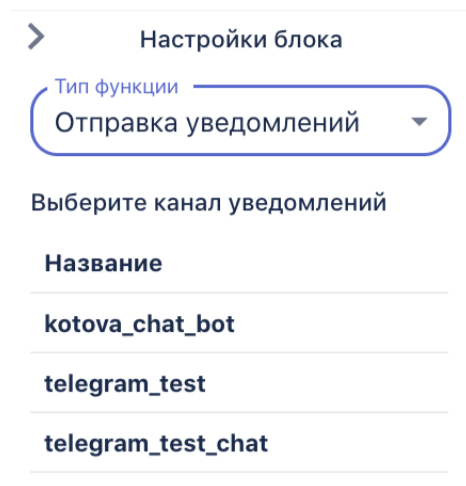


Рисунок 16.8 - Настройка блока «Отправка уведомлений»

Условием для отправки будет служить блок, идущий перед блоком уведомлений. Например, это может быть шлюз, где задаётся параметр, при котором уведомление должно быть отправлено:

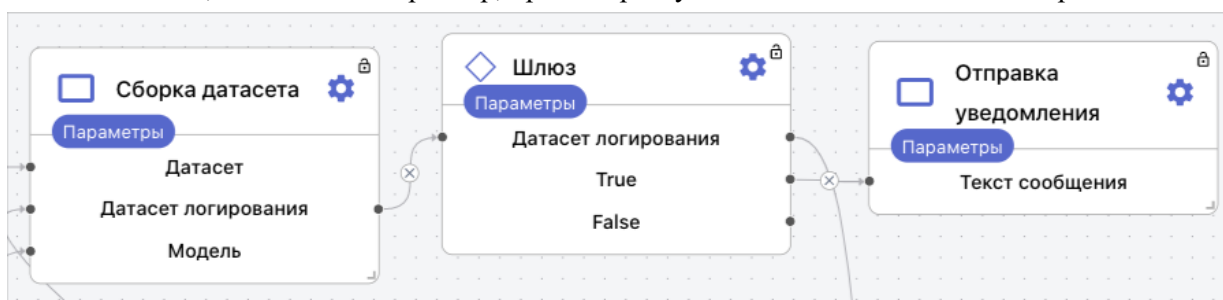


Рисунок 16.9 - Пример построения пайплайна с блоком уведомления

Например, в результате отработки такого пайплайна, пользователи будут получать следующую информацию:

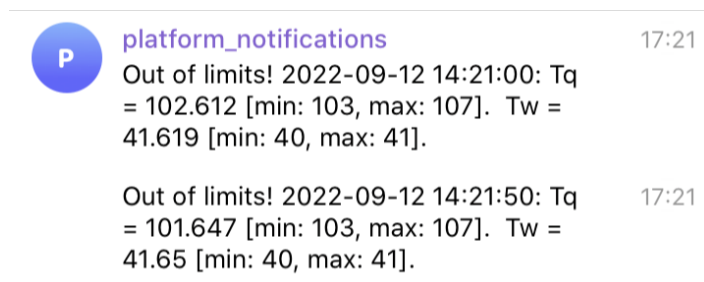


Рисунок 16.10 - Пример уведомлений в телеграм канале

## 17. Дополнительные возможности Платформы

### - 17.1. Обращение в службу поддержки

В верхнем правом углу окна системы находится кнопка «Сообщить об ошибке»:

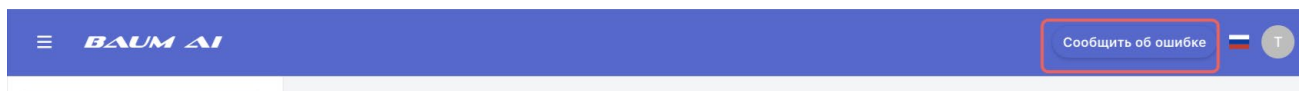


Рисунок 17.1 – Кнопка «Сообщить об ошибке»

В случае, если пользователь сталкивается с какой-то проблемой в работе платформы, поведением, не описанным в данном руководстве, есть возможность связаться с командой поддержки системы и сообщить о проблеме.

Пользователю необходимо максимально подробно описать проблему и шаги, которые привели к её появлению в текстовом окне сообщения об ошибке, а после нажать кнопку «Отправить»:

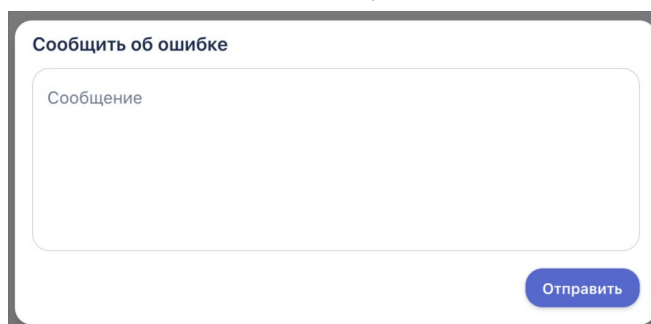
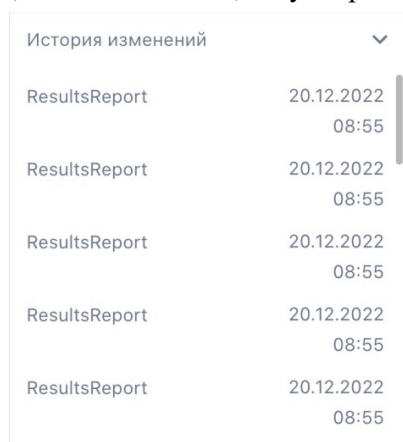
The image shows a form titled 'Сообщить об ошибке' (Report error). It contains a large text input field with the placeholder text 'Сообщение' (Message). At the bottom right of the form is a blue button labeled 'Отправить' (Send).

Рисунок 17.2 – Описание найденной в Системе ошибки, и отправка на обработку Разработчику

Ограничение размера сообщения составляет 256 символов.

### - 17.2. История изменений

История изменений - это раздел системы, отображающий информацию о том в каком блоке и в какое время было совершено какое либо изменение. Раздел содержит информацию о всех запусках пайплайнов, сформированных отчетах, названия блоков, дату и время каждого изменения.

The image shows a scrollable list titled 'История изменений' (Change History). The list contains five entries, each with the text 'ResultsReport' and a timestamp '20.12.2022 08:55'. A vertical scrollbar is visible on the right side of the list.

История изменений	
ResultsReport	20.12.2022 08:55
ResultsReport	20.12.2022 08:55
ResultsReport	20.12.2022 08:55
ResultsReport	20.12.2022 08:55
ResultsReport	20.12.2022 08:55

Рисунок 17.3 – История изменений

## 18. Приложения

### - Приложение 1. Автоматизированные функции

\*Журнал преобразований – технический объект, который позволяет пробрасывать данные между элементами блок-схемы.

\*\*Выходным параметром для всех функций является *словарь с данными*, который может включать в себя: таблицы, графики, текстовое описание.

Таблица 18.1 – Перечень автоматизированных функций элемента «Источник данных»

Функция	Назначение	Параметры	Выходная информация		
<b>Группа «Загрузка данных»</b>					
<b>Загрузка табличных данных</b>	Преобразует загруженные данные во временной ряд. При этом автоматически определяет: формат файла (csv, txt, xls, xlsx), разделитель, размерность, где <i>разделитель</i> – это символ в датасете с временным рядом, обозначающий деление на ячейки. Для <i>ресемплирования</i> использует агрегирующую функцию	<table border="1"> <tr> <td><b>Выберите файл</b></td> <td>Название файла, содержащего временной ряд для анализа</td> </tr> </table>	<b>Выберите файл</b>	Название файла, содержащего временной ряд для анализа	**
<b>Выберите файл</b>	Название файла, содержащего временной ряд для анализа				
<b>Загрузка изображений для классификации</b>	Загружаются изображения с объектами для дальнейшей классификации этих объектов на новых данных. При этом предусмотрено, что изображения с объектами	<table border="1"> <tr> <td><b>Группа обучающих изображений</b></td> <td>Выбрать папку для <i>обучения</i> нейронной сети, который представляет собой каталог с папками, где каждая папка содержит изображения с объектами одного класса.</td> </tr> </table>	<b>Группа обучающих изображений</b>	Выбрать папку для <i>обучения</i> нейронной сети, который представляет собой каталог с папками, где каждая папка содержит изображения с объектами одного класса.	В БД сохраняются каталоги датасета обучения и валидации, прошедшие предварительную обработку перед загрузкой.
<b>Группа обучающих изображений</b>	Выбрать папку для <i>обучения</i> нейронной сети, который представляет собой каталог с папками, где каждая папка содержит изображения с объектами одного класса.				

	<p>одного класса распределены по отдельным папкам. При загрузке выполняется <i>ресайз</i> изображений – изменение (чаще уменьшение) размера изображений до заданного формата. Данные загружаются маленькими порциями, так называемыми <i>мини-батчами</i> (например, за один раз подается два изображения).</p>	<table border="1"> <tr> <td data-bbox="835 244 1133 347"></td> <td data-bbox="1133 244 1722 347"> <p>Например – папки с изображениями собак, кошек и диких животных (три класса).</p> </td> </tr> <tr> <td data-bbox="835 347 1133 523"> <p><b>Группа тестовых изображений</b></p> </td> <td data-bbox="1133 347 1722 523"> <p>Выбрать папку для <i>валидации</i> обученной нейронной сети, который также представлен в виде каталога с папками отдельных классов.</p> </td> </tr> <tr> <td data-bbox="835 523 1133 667"> <p><b>Размер мини-батча</b></p> </td> <td data-bbox="1133 523 1722 667"> <p>Указывается количество изображений, которое за один раз подается на вход нейронной сети для её обучения.</p> </td> </tr> <tr> <td data-bbox="835 667 1133 730"> <p><b>Новая высота</b></p> </td> <td data-bbox="1133 667 1722 730"> <p>Новая высота изображения</p> </td> </tr> <tr> <td data-bbox="835 730 1133 799"> <p><b>Новая ширина</b></p> </td> <td data-bbox="1133 730 1722 799"> <p>Новая ширина изображения</p> </td> </tr> </table>		<p>Например – папки с изображениями собак, кошек и диких животных (три класса).</p>	<p><b>Группа тестовых изображений</b></p>	<p>Выбрать папку для <i>валидации</i> обученной нейронной сети, который также представлен в виде каталога с папками отдельных классов.</p>	<p><b>Размер мини-батча</b></p>	<p>Указывается количество изображений, которое за один раз подается на вход нейронной сети для её обучения.</p>	<p><b>Новая высота</b></p>	<p>Новая высота изображения</p>	<p><b>Новая ширина</b></p>	<p>Новая ширина изображения</p>	
	<p>Например – папки с изображениями собак, кошек и диких животных (три класса).</p>												
<p><b>Группа тестовых изображений</b></p>	<p>Выбрать папку для <i>валидации</i> обученной нейронной сети, который также представлен в виде каталога с папками отдельных классов.</p>												
<p><b>Размер мини-батча</b></p>	<p>Указывается количество изображений, которое за один раз подается на вход нейронной сети для её обучения.</p>												
<p><b>Новая высота</b></p>	<p>Новая высота изображения</p>												
<p><b>Новая ширина</b></p>	<p>Новая ширина изображения</p>												
<p><b>Загрузка табличных данных из коннектора</b></p>	<p>Функция предназначена для подключения к источникам данным в виде баз данных «ClickHouse» или «Postgresql». При этом используется сущность «Коннектор», в которой прописываются настройки для подключения к этим базам данных. Реализована возможность формирования датасета на основании полученных данных. Для этого в настройках функции активируется галочка в поле «Сохранить датасет», и в поле «Название файла»</p>	<table border="1"> <tr> <td data-bbox="835 799 1133 1102"> <p><b>Выберите файл</b></p> </td> <td data-bbox="1133 799 1722 1102"> <p>Из списка выбирается <i>коннектор</i> – источник данных, подключение к которому необходимо выполнить. Источником выступает сторонняя база данных – ClickHouse или Postgresql (соответственно из списка выбирается коннектор с таким типом).</p> </td> </tr> <tr> <td data-bbox="835 1102 1133 1410"> <p><b>Сохранить датасет</b></p> </td> <td data-bbox="1133 1102 1722 1410"> <p>Чтобы сформировать бэкап таблицы внешней БД в поле «Сохранить датасет» устанавливается галочка. Иначе, если не установить галочку в поле «Сохранить датасет», выполняется подключение к внешней БД в её состоянии на текущий момент времени, без</p> </td> </tr> </table>	<p><b>Выберите файл</b></p>	<p>Из списка выбирается <i>коннектор</i> – источник данных, подключение к которому необходимо выполнить. Источником выступает сторонняя база данных – ClickHouse или Postgresql (соответственно из списка выбирается коннектор с таким типом).</p>	<p><b>Сохранить датасет</b></p>	<p>Чтобы сформировать бэкап таблицы внешней БД в поле «Сохранить датасет» устанавливается галочка. Иначе, если не установить галочку в поле «Сохранить датасет», выполняется подключение к внешней БД в её состоянии на текущий момент времени, без</p>	<p>**</p>						
<p><b>Выберите файл</b></p>	<p>Из списка выбирается <i>коннектор</i> – источник данных, подключение к которому необходимо выполнить. Источником выступает сторонняя база данных – ClickHouse или Postgresql (соответственно из списка выбирается коннектор с таким типом).</p>												
<p><b>Сохранить датасет</b></p>	<p>Чтобы сформировать бэкап таблицы внешней БД в поле «Сохранить датасет» устанавливается галочка. Иначе, если не установить галочку в поле «Сохранить датасет», выполняется подключение к внешней БД в её состоянии на текущий момент времени, без</p>												



	<p>указывается имя для бэкапа таблицы базы данных в настоящий момент времени. В результате в разделе «Данные» сохраняется файл в формате .csv с данными из коннектора. <b>Важно</b> – Подключение выполняется к БД в её состоянии на текущий момент времени.</p>	<table border="1"> <tr> <td data-bbox="842 244 1133 347"></td> <td data-bbox="1133 244 1715 347">дополнительного формирования датасета во внутренней БД.</td> </tr> <tr> <td data-bbox="842 347 1133 443"><b>Название файла</b></td> <td data-bbox="1133 347 1715 443">Указывается название файла бэкапа таблицы для сохранения во внутренней БД.</td> </tr> <tr> <td data-bbox="842 443 1133 539"><b>Онлайн данные</b></td> <td data-bbox="1133 443 1715 539">Указывается, осуществляется ли получение данных в режиме он лайн.</td> </tr> </table>		дополнительного формирования датасета во внутренней БД.	<b>Название файла</b>	Указывается название файла бэкапа таблицы для сохранения во внутренней БД.	<b>Онлайн данные</b>	Указывается, осуществляется ли получение данных в режиме он лайн.	
	дополнительного формирования датасета во внутренней БД.								
<b>Название файла</b>	Указывается название файла бэкапа таблицы для сохранения во внутренней БД.								
<b>Онлайн данные</b>	Указывается, осуществляется ли получение данных в режиме он лайн.								
<b>Загрузка модели</b>	<p>Функция предназначена для использования в качестве источника данных ранее обученной модели. При этом система при обработке данных пайплайна применяет ранее полученные знания для построения прогнозов.</p>	<table border="1"> <tr> <td data-bbox="842 619 1133 751"><b>Модель</b></td> <td data-bbox="1133 619 1715 751">Выбор из списка ранее сохраненных моделей</td> </tr> </table>	<b>Модель</b>	Выбор из списка ранее сохраненных моделей					
<b>Модель</b>	Выбор из списка ранее сохраненных моделей								
<b>Загрузка текстовых файлов для классификации</b>	<p>Данная функция предназначена для загрузки текстов, принадлежащих к тем или иным классам, для обучения нейронной сети определять эти классы на новых данных. Функция обязательно используется при решении задач классификации текстов.</p>	<table border="1"> <tr> <td data-bbox="842 895 1133 1198"><b>Группа обучающих текстов</b></td> <td data-bbox="1133 895 1715 1198">Выбор папки для обучения нейронной сети, которая должна содержать в себе подпапки с названиями классов объектов. Данные подпапки содержат тексты, принадлежащие определенному классу. Например, это могут быть: «Пушкин», «Лермонтов», «Голстой».</td> </tr> <tr> <td data-bbox="842 1198 1133 1326"><b>Группа тестовых текстов</b></td> <td data-bbox="1133 1198 1715 1326">Выбор папки для валидации обученной нейронной сети. Папка должна иметь такую же структуру, как и обучающая.</td> </tr> </table>	<b>Группа обучающих текстов</b>	Выбор папки для обучения нейронной сети, которая должна содержать в себе подпапки с названиями классов объектов. Данные подпапки содержат тексты, принадлежащие определенному классу. Например, это могут быть: «Пушкин», «Лермонтов», «Голстой».	<b>Группа тестовых текстов</b>	Выбор папки для валидации обученной нейронной сети. Папка должна иметь такую же структуру, как и обучающая.			
<b>Группа обучающих текстов</b>	Выбор папки для обучения нейронной сети, которая должна содержать в себе подпапки с названиями классов объектов. Данные подпапки содержат тексты, принадлежащие определенному классу. Например, это могут быть: «Пушкин», «Лермонтов», «Голстой».								
<b>Группа тестовых текстов</b>	Выбор папки для валидации обученной нейронной сети. Папка должна иметь такую же структуру, как и обучающая.								

		<table border="1"> <tr> <td data-bbox="846 252 1133 411"><b>Группа текстов для классификации</b></td> <td data-bbox="1133 252 1715 411">Здесь можно сразу выбрать файл или папку с файлами, которые необходимо классифицировать с применением обученной модели.</td> </tr> </table>	<b>Группа текстов для классификации</b>	Здесь можно сразу выбрать файл или папку с файлами, которые необходимо классифицировать с применением обученной модели.	
<b>Группа текстов для классификации</b>	Здесь можно сразу выбрать файл или папку с файлами, которые необходимо классифицировать с применением обученной модели.				
<b>Загрузка текстовых файлов для кластеризации</b>	Функция обязательно используется при решении задач кластеризации текстов, когда необходимо определить кластеры к которым принадлежат тексты.	<table border="1"> <tr> <td data-bbox="846 515 1133 651"><b>Группа текстов для кластеризации</b></td> <td data-bbox="1133 515 1715 651">Выбор файла, содержащего однотипные данные, подлежащие разделению на кластеры.</td> </tr> </table>	<b>Группа текстов для кластеризации</b>	Выбор файла, содержащего однотипные данные, подлежащие разделению на кластеры.	
<b>Группа текстов для кластеризации</b>	Выбор файла, содержащего однотипные данные, подлежащие разделению на кластеры.				
<b>Загрузка графа</b>	Функция предназначена для загрузки и дальнейшего преобразования файлов с форматом .graphml в переменную graph_out с типом данных networkx.MultiDiGraph, предназначенных для решения задач с применением теории графов. Граф — это геометрическая фигура, которая состоит из точек и линий, которые их соединяют. Точки называют вершинами графа, а линии — ребрами. Графы имеют очень широкое применение: с их помощью выбирают наиболее выгодное расположение зданий, графами представлены схемы	<table border="1"> <tr> <td data-bbox="846 754 1133 858"><b>Выберите файл с графом</b></td> <td data-bbox="1133 754 1715 858">Выбирается ранее загруженный в систему файл в формате .graphml.</td> </tr> </table>	<b>Выберите файл с графом</b>	Выбирается ранее загруженный в систему файл в формате .graphml.	
<b>Выберите файл с графом</b>	Выбирается ранее загруженный в систему файл в формате .graphml.				

	метро, маршруты, схемы игр, блок схемы процессов и т.д.				
<b>Группа «Spark»</b>					
<b>Загрузка табличных данных из файла CSV (Spark)</b>	<p>При помощи данной функции осуществляется загрузка в систему табличных данных с помощью фреймворка для распределенных вычислений Apache Spark, конкретно, с помощью библиотеки PySpark для Python. Датафрейм в PySpark — это таблица, строки которой хранятся в RDD (Отказоустойчивый распределенный набор данных (англ. Resilient Distributed Dataset, RDD) — тип структуры данных, который можно распределить между несколькими узлами в кластере). Работа с датафреймами ведётся по принципу «ленивых вычислений» (англ. lazy evaluations). Это вычисления, которые откладываются до тех пор, пока пользователь не запросит их результат. Данная функция работает только для файлов в формате csv, содержащих big data.</p>	<table border="1"> <tr> <td><b>Выберите файл для загрузки</b></td> <td>Выбор из списка <i>файла</i> для дальнейшего анализа.</td> </tr> </table>		<b>Выберите файл для загрузки</b>	Выбор из списка <i>файла</i> для дальнейшего анализа.
		<b>Выберите файл для загрузки</b>	Выбор из списка <i>файла</i> для дальнейшего анализа.		

<b>Загрузка табличных данных из папки CSV (Spark)</b>	Фреймворк «Apache Spark» распределяет хранимые данные по серверам и директориям. Чтобы обратиться к файлам на уровне папки, в которой они хранятся, используется данный метод.	<table border="1"> <tr> <td data-bbox="846 284 1126 443"> <b>Выберите директорию с датасетом для загрузки</b> </td> <td data-bbox="1126 284 1718 443">                     Выбор из списка <i>папки</i> для дальнейшего анализа.                 </td> </tr> </table>	<b>Выберите директорию с датасетом для загрузки</b>	Выбор из списка <i>папки</i> для дальнейшего анализа.	
<b>Выберите директорию с датасетом для загрузки</b>	Выбор из списка <i>папки</i> для дальнейшего анализа.				
<b>Загрузка модели</b>	Функция предназначена для использования в качестве источника данных ранее обученной модели ИИ «Spark».	<table border="1"> <tr> <td data-bbox="846 555 1126 754"> <b>Модель</b> </td> <td data-bbox="1126 555 1718 754">                     Выбор из списка ранее сохраненных моделей Spark (в разработке отдельное <i>API</i> для моделей Spark. Модели будут объединяться в одну группу по семантическому типу).                 </td> </tr> </table>	<b>Модель</b>	Выбор из списка ранее сохраненных моделей Spark (в разработке отдельное <i>API</i> для моделей Spark. Модели будут объединяться в одну группу по семантическому типу).	
<b>Модель</b>	Выбор из списка ранее сохраненных моделей Spark (в разработке отдельное <i>API</i> для моделей Spark. Модели будут объединяться в одну группу по семантическому типу).				
<b>Загрузка табличных данных из коннектора (Spark)</b>	Функция предназначена для получения табличных данных через коннектор с типом «ClickHouse», с использованием библиотеки «Spark».	<i>(в разработке)</i>			

**Таблица 18.2 – Перечень автоматизированных функций элемента «Процесс»**

Функция	Назначение	Параметры	Выходная информация
<b>1.Группа «Анализ данных»</b>			

<p><b>Выбор признаков и целевых признаков</b></p>	<p>В датасете выбираются: признаки – измеримые характеристики исследуемого объекта/процесса, и целевые (зависимые) переменные, значения которых предстоит предсказывать модели.</p>	<table border="1"> <tr> <td data-bbox="837 284 1133 411"> <p><b>Признаки</b></p> </td> <td data-bbox="1133 284 1724 411"> <p>Характеристики, которые исследуются и выявляется корреляция между ними и рассматриваемым целевым признаком</p> </td> </tr> <tr> <td data-bbox="837 411 1133 512"> <p><b>Целевые признаки</b></p> </td> <td data-bbox="1133 411 1724 512"> <p>Предсказываемые переменные</p> </td> </tr> </table>	<p><b>Признаки</b></p>	<p>Характеристики, которые исследуются и выявляется корреляция между ними и рассматриваемым целевым признаком</p>	<p><b>Целевые признаки</b></p>	<p>Предсказываемые переменные</p>	<p>Датасет с размеченными признаками и целевыми признаками</p>
<p><b>Признаки</b></p>	<p>Характеристики, которые исследуются и выявляется корреляция между ними и рассматриваемым целевым признаком</p>						
<p><b>Целевые признаки</b></p>	<p>Предсказываемые переменные</p>						
<p><b>Матрица корреляции</b></p>	<p>1.Алгоритм сначала рассчитывает коэффициенты корреляции по всем признакам (общая матрица корреляции). 2.Затем в этой матрице отбираются топ-k максимальных (ближе к 1) значений коэффициентов корреляции. 3.Строится новая матрица корреляции, состоящая из признаков, для которых найдены максимальные значения коэффициентов.</p>	<table border="1"> <tr> <td data-bbox="837 651 1133 762"> <p><b>Топ k-значений для корреляции</b></p> </td> <td data-bbox="1133 651 1724 762"> <p>Количество максимальных значений корреляции (int)</p> </td> </tr> </table>	<p><b>Топ k-значений для корреляции</b></p>	<p>Количество максимальных значений корреляции (int)</p>	<p>1. Матрица корреляции с топ-k признаков, имеющих максимальные значения корреляции. 2.Матрица корреляции по всем признакам.</p>		
<p><b>Топ k-значений для корреляции</b></p>	<p>Количество максимальных значений корреляции (int)</p>						

<p><b>Косинусное расстояние</b></p>	<p>Вычисляется <i>косинусное расстояние</i> между значениями во входном векторе и значениями выбранных столбцов в наблюдениях.</p>	<table border="1"> <tr> <td data-bbox="844 284 1124 355"><b>Датасет</b></td> <td data-bbox="1124 284 1718 355">Датасет с исходными данными.</td> </tr> <tr> <td data-bbox="844 355 1124 427"><b>Признаки</b></td> <td data-bbox="1124 355 1718 427">Признаки для анализа.</td> </tr> <tr> <td data-bbox="844 427 1124 547"><b>Вектор</b></td> <td data-bbox="1124 427 1718 547">Вектор такой же длины, что и количество выбранных в датасете признаков (1D-array).</td> </tr> </table>	<b>Датасет</b>	Датасет с исходными данными.	<b>Признаки</b>	Признаки для анализа.	<b>Вектор</b>	Вектор такой же длины, что и количество выбранных в датасете признаков (1D-array).	<p>Таблица с наблюдениями, наиболее схожими с входным вектором, где первые топ-5 наблюдений выделены жирным шрифтом (по ним рассчитанная косинусная мера имеет значение, наиболее близкое к 0).</p>
<b>Датасет</b>	Датасет с исходными данными.								
<b>Признаки</b>	Признаки для анализа.								
<b>Вектор</b>	Вектор такой же длины, что и количество выбранных в датасете признаков (1D-array).								
<p><b>Поиск пропущенных значений</b></p>	<p>Для каждого выбранного признака метод находит пропущенные значения в наблюдениях.</p>	<table border="1"> <tr> <td data-bbox="844 627 1124 722"><b>Признаки</b></td> <td data-bbox="1124 627 1718 722">Выбираются признаки, в которых необходимо найти пропущенные значения.</td> </tr> </table>	<b>Признаки</b>	Выбираются признаки, в которых необходимо найти пропущенные значения.	<p>Словарь, в котором по каждому выбранному признаку отображается количество пропущенных значений в наблюдениях, и процент пропусков.</p>				
<b>Признаки</b>	Выбираются признаки, в которых необходимо найти пропущенные значения.								
<p><b>Анализ временных рядов</b></p>	<p>Совокупность математико-статистических методов анализа, предназначенных для выявления структуры временных рядов и для их прогнозирования</p>	<table border="1"> <tr> <td data-bbox="844 834 1124 994"><b>Список графиков</b></td> <td data-bbox="1124 834 1718 994">Выбираются графики, которые будут использованы для визуализации анализа временных рядов и задаются параметры для них</td> </tr> </table>	<b>Список графиков</b>	Выбираются графики, которые будут использованы для визуализации анализа временных рядов и задаются параметры для них	<p>Линейный график, ACF/PACF, Декомпозиция, Свечной график, Time profile, Extended, Bollinger Bands, Stochastic Oscillator</p>				
<b>Список графиков</b>	Выбираются графики, которые будут использованы для визуализации анализа временных рядов и задаются параметры для них								
<p><b>Запись в датасет логирования</b></p>	<p>Данная функция применяется при построении пайплайнов в режиме реального времени. В процессе логирования записываются все новые поступающие значения в датасет для дальнейшего использования при валидации и работе с моделью. При этом в</p>	<p>-</p>	<p>Датасет с данными</p>						

	<p>процессе логирования записываются:</p> <ul style="list-style-type: none"> <li>- фактические значения</li> <li>- промежуточные трансформации при препроцессинге до подачи в модель (если в качестве источника данных выбрана модель)</li> <li>- прогнозные значения (если в пайплайне настроен прогноз)</li> </ul>								
<p><b>Визуализация Real-time</b></p>	<p>Представление в виде графиков и диаграмм результатов работы пайплайна, содержащего данные, получаемые в режиме реального времени</p>	<table border="1"> <tr> <td data-bbox="846 722 1126 887"> <p><b>Список графиков</b></p> </td> <td data-bbox="1126 722 1720 887"> <p>Выбираются графики, которые будут использованы для визуализации данных в режиме он лайн и задаются параметры для них</p> </td> </tr> </table>	<p><b>Список графиков</b></p>	<p>Выбираются графики, которые будут использованы для визуализации данных в режиме он лайн и задаются параметры для них</p>	<p>Линейный график, Свечной график, Time Profile, Extended, Bollinger Bands, Stochastic Oscillator</p>				
<p><b>Список графиков</b></p>	<p>Выбираются графики, которые будут использованы для визуализации данных в режиме он лайн и задаются параметры для них</p>								
<p><b>Загрузка данных</b></p>									
<p><b>Преобразование данных во временной ряд</b></p>	<p>Метод редактирует исходные данные, исключая в них аномалии и искаженные наблюдения, которые могли быть зафиксированы в результате помех. Далее выполняется <i>дискретизация</i> – определяются точки (моменты времени), в которых должны быть произведены выборки значений. Дискретизация</p>	<table border="1"> <tr> <td data-bbox="846 986 1126 1185"> <p><b>Шаг ресемплирования</b></p> </td> <td data-bbox="1126 986 1720 1185"> <p><i>Дискретность</i> для временного ряда – частота фиксирования наблюдений, значения начиная с нано-, микро-, милли-, секунд и заканчивая годами. Указывается оптимальный интервал дискретности.</p> </td> </tr> <tr> <td data-bbox="846 1185 1126 1281"> <p><b>Частота ресемплирования</b></p> </td> <td data-bbox="1126 1185 1720 1281"> <p>Единица измерения, в которой фиксируются наблюдения.</p> </td> </tr> <tr> <td data-bbox="846 1281 1126 1412"> <p><b>Агрегирующая функция</b></p> </td> <td data-bbox="1126 1281 1720 1412"> <p>Функция, вычисляющая результат по набору значений группы, где группа – наблюдения в пределах шага</p> </td> </tr> </table>	<p><b>Шаг ресемплирования</b></p>	<p><i>Дискретность</i> для временного ряда – частота фиксирования наблюдений, значения начиная с нано-, микро-, милли-, секунд и заканчивая годами. Указывается оптимальный интервал дискретности.</p>	<p><b>Частота ресемплирования</b></p>	<p>Единица измерения, в которой фиксируются наблюдения.</p>	<p><b>Агрегирующая функция</b></p>	<p>Функция, вычисляющая результат по набору значений группы, где группа – наблюдения в пределах шага</p>	<p>Создается временной ряд, с заданным шагом ресемплирования</p>
<p><b>Шаг ресемплирования</b></p>	<p><i>Дискретность</i> для временного ряда – частота фиксирования наблюдений, значения начиная с нано-, микро-, милли-, секунд и заканчивая годами. Указывается оптимальный интервал дискретности.</p>								
<p><b>Частота ресемплирования</b></p>	<p>Единица измерения, в которой фиксируются наблюдения.</p>								
<p><b>Агрегирующая функция</b></p>	<p>Функция, вычисляющая результат по набору значений группы, где группа – наблюдения в пределах шага</p>								

	производится через равные промежутки времени.		ресемплирования. По умолчанию, значение вычисляется функцией медианы.	
		<b>Столбец с временной меткой</b>	Время фиксирования наблюдения. По умолчанию, нулевой столбец в датасете.	
<b>Препроцессинг</b>				
<b>Стабилизация дисперсии</b>	Уменьшает разброс исследуемых данных, чтобы сделать их более компактными и пригодными для работы.	<b>Замена значений столбцов</b>	Преобразование оригинального временного ряда, загруженного в систему. Позволяет заменять трансформируемые столбцы или добавлять новые	Преобразованный датасет. При этом преобразования над целевыми признаками проводятся отдельно.
		<b>Стандартизация</b>	Преобразование значений признака, адаптирующая признаки с разными диапазонами значений к моделям машинного обучения	
		<b>Метод</b>	Выбирается метод, с помощью которого проводится стабилизация дисперсии – приведение данных к нормальному распределению. На выбор два метода – уео-johnson и box-cox. Метод уео-johnson работает как с отрицательными, так и с положительными значениями, а метод box-cox только с положительными	
		<b>Флаг признака</b>	Показатели датасета, значения которых предстоит предсказывать модели машинного обучения	



<p><b>Стандартизация</b></p>	<p>Чтобы сгладить большие различия между диапазонами признаков датасета и предотвратить искаженное восприятие данных моделью машинного обучения выполняется <i>стандартизация</i> – преобразование и приведение признаков датасета к единому формату</p>	<table border="1"> <tr> <td data-bbox="840 296 1160 475"> <p><b>Замена значений столбцов</b></p> </td> <td data-bbox="1160 296 1711 475"> <p>Подтверждение преобразования оригинального временного ряда (заменой трансформируемых столбцов или добавлением новых)</p> </td> </tr> <tr> <td data-bbox="840 475 1160 603"> <p><b>Флаг признака</b></p> </td> <td data-bbox="1160 475 1711 603"> <p>Выбрать столбцы для преобразования – все, кроме столбцов с датой и целевым признаком</p> </td> </tr> </table>	<p><b>Замена значений столбцов</b></p>	<p>Подтверждение преобразования оригинального временного ряда (заменой трансформируемых столбцов или добавлением новых)</p>	<p><b>Флаг признака</b></p>	<p>Выбрать столбцы для преобразования – все, кроме столбцов с датой и целевым признаком</p>	<p>Преобразованные значения показателей временного ряда, кроме целевого признака</p>
<p><b>Замена значений столбцов</b></p>	<p>Подтверждение преобразования оригинального временного ряда (заменой трансформируемых столбцов или добавлением новых)</p>						
<p><b>Флаг признака</b></p>	<p>Выбрать столбцы для преобразования – все, кроме столбцов с датой и целевым признаком</p>						
<p><b>Дифференцирование временного ряда</b></p>	<p>Выполняется дифференцирование целевых признаков (таргетов) временного ряда. При этом временной ряд сдвигается на указанное число шагов в разрезе каждого целевого признака. Если есть сезонность, сначала проводится сезонное дифференцирование. Желательно дифференцировать ряд как можно меньше раз, потому что с увеличением количества дифференцирований растет дисперсия ошибки прогноза</p>	<table border="1"> <tr> <td data-bbox="840 691 1093 1094"> <p><b>Шаг дифференцирования для каждого целевого признака</b></p> </td> <td data-bbox="1093 691 1711 1094"> <p>Есть возможность задать шаг дифференцирования для каждого таргета, в формате [сдвиг для таргета 1, сдвиг для таргета 2, ...], где сдвиг на один шаг применяется для обычного (для избавления от тренда) дифференцирования, сдвиг на несколько шагов – для сезонного, а сдвиг, равный нулю означает, что дифференцирование для данного таргета не проводится. Например, [1, 0, 3].</p> </td> </tr> </table>	<p><b>Шаг дифференцирования для каждого целевого признака</b></p>	<p>Есть возможность задать шаг дифференцирования для каждого таргета, в формате [сдвиг для таргета 1, сдвиг для таргета 2, ...], где сдвиг на один шаг применяется для обычного (для избавления от тренда) дифференцирования, сдвиг на несколько шагов – для сезонного, а сдвиг, равный нулю означает, что дифференцирование для данного таргета не проводится. Например, [1, 0, 3].</p>	<p>К датасету временного ряда добавляются новые столбцы с окончанием '_diff' для каждого указанного таргета. При этом замена колонок не предусмотрена – оригинальные колонки сохраняются для задачи обратного дифференцирования. Отображаются графики настоящего и сдвинутого временных рядов.</p>		
<p><b>Шаг дифференцирования для каждого целевого признака</b></p>	<p>Есть возможность задать шаг дифференцирования для каждого таргета, в формате [сдвиг для таргета 1, сдвиг для таргета 2, ...], где сдвиг на один шаг применяется для обычного (для избавления от тренда) дифференцирования, сдвиг на несколько шагов – для сезонного, а сдвиг, равный нулю означает, что дифференцирование для данного таргета не проводится. Например, [1, 0, 3].</p>						
<p><b>One-Hot Encoding</b></p>	<p>Метод One Hot Encoding (ONE) применяется, когда в датасете необходимо</p>	<table border="1"> <tr> <td data-bbox="840 1337 1077 1433"> <p><b>Флаг удаления</b></p> </td> <td data-bbox="1077 1337 1688 1433"> <p>Устанавливается, чтобы удалить из итоговой таблицы столбец с признаком, над</p> </td> </tr> </table>	<p><b>Флаг удаления</b></p>	<p>Устанавливается, чтобы удалить из итоговой таблицы столбец с признаком, над</p>	<p>Таблица с новыми столбцами, в которых отражается принадлежность наблюдений к</p>		
<p><b>Флаг удаления</b></p>	<p>Устанавливается, чтобы удалить из итоговой таблицы столбец с признаком, над</p>						

	<p>закодировать категориальные признаки (текстовые), перед подачей в модель. Для кодируемого категориального признака создаются N новых столбцов в датасете, где N – количество уникальных категорий. Значения в новых столбцах – 0 или 1, в зависимости от принадлежности к категории. Так каждый новый признак – бинарный характеристический признак категории.</p>	<p><b>первого признака</b></p>	<p>которым были выполнены преобразования. Так как новые столбцы отражают принадлежность наблюдения к той или иной категории признака, удаление первого признака не повлияет на результат.</p>	<p>тем или иным категориям преобразованных признаков.</p>
<p><b>Создание признаков для временного ряда</b></p>	<p>Для временного ряда создаются новые признаки, в которых значения таргетов сдвигаются на указанное число шагов. Например, если для одномерного (с одним таргетом) временного ряда задать сдвиг в один шаг, создается новая колонка y, в которой значение первой строки равно значению второй строки в колонке с таргетом x (категориальный признак), т.е. значения сдвигаются на один шаг вперед. Если ряд многомерный – состоящий из нескольких таргетов, то для каждого из них</p>	<p><b>Максимальное количество лагов</b></p>	<p>Указывается, на какое количество шагов может быть сдвинут временной ряд.</p>	<p>Создается таблица, график со смещенным временным рядом. Процесс создания признаков-лагов сохраняется в журнале преобразований, и далее отрабатывается при препроцессинге.</p>

	<p>передается общий массив признаков, с учетом лагов всех таргет-рядов. Такое действие является предварительным перед тем, как подавать данные в модель ИИ, чтобы у модели были не только фактические значения таргета, но и прогнозные.</p>				
<p><b>Преоброессинг текстовых данных</b></p>	<p>Алгоритмы машинного обучения не работают с «сырыми данными». Большая часть процесса – это подготовка текста, преобразование ее в вид, доступный для восприятия компьютером. В первую очередь выполняется <i>очистка</i> текста. Из текста удаляются бесполезные для машины данные – это большинство знаков пунктуации, особые символы, скобки, теги и т.д. Дальше наступает большой этап предварительной обработки – <i>преоброессинга</i>. Это приведение информации к виду, в котором она более понятна алгоритму. Методы преоброессинга:</p>	<table border="1"> <tr> <td data-bbox="840 667 1086 1407"> <p><b>Метод векторизации</b></p> </td> <td data-bbox="1086 667 1693 1407"> <p>1. <b>TF-IDF</b>. С англ. TF – term frequency (частота слова), IDF – inverse document frequency (обратная частота документа). Это мера, используемая для оценки важности слова в контексте документа, являющегося частью коллекции документов. Вес некоторого слова пропорционален частоте употребления этого слова в документе, и обратно пропорционален частоте употребления слова во всех документах коллекции. Метод TF-IDF используется в задачах <i>анализа текстов</i>, и представляет собой линейный классификатор с разреженными признаками, взвешенными по частоте. Этот метод выбирается по умолчанию.</p> <p>2. <b>Word2vec</b>. Принимает большой <i>текстовый корпус</i> в качестве входных данных, и сопоставляет каждому слову вектор, выдавая координаты слов на выходе. Сначала он генерирует словарь корпуса, а</p> </td> </tr> </table>	<p><b>Метод векторизации</b></p>	<p>1. <b>TF-IDF</b>. С англ. TF – term frequency (частота слова), IDF – inverse document frequency (обратная частота документа). Это мера, используемая для оценки важности слова в контексте документа, являющегося частью коллекции документов. Вес некоторого слова пропорционален частоте употребления этого слова в документе, и обратно пропорционален частоте употребления слова во всех документах коллекции. Метод TF-IDF используется в задачах <i>анализа текстов</i>, и представляет собой линейный классификатор с разреженными признаками, взвешенными по частоте. Этот метод выбирается по умолчанию.</p> <p>2. <b>Word2vec</b>. Принимает большой <i>текстовый корпус</i> в качестве входных данных, и сопоставляет каждому слову вектор, выдавая координаты слов на выходе. Сначала он генерирует словарь корпуса, а</p>	<p><i>Числовые векторы</i>, созданные на основе исходной текстовой информации, которые отражают <i>важность</i> использования каждого слова.</p>
<p><b>Метод векторизации</b></p>	<p>1. <b>TF-IDF</b>. С англ. TF – term frequency (частота слова), IDF – inverse document frequency (обратная частота документа). Это мера, используемая для оценки важности слова в контексте документа, являющегося частью коллекции документов. Вес некоторого слова пропорционален частоте употребления этого слова в документе, и обратно пропорционален частоте употребления слова во всех документах коллекции. Метод TF-IDF используется в задачах <i>анализа текстов</i>, и представляет собой линейный классификатор с разреженными признаками, взвешенными по частоте. Этот метод выбирается по умолчанию.</p> <p>2. <b>Word2vec</b>. Принимает большой <i>текстовый корпус</i> в качестве входных данных, и сопоставляет каждому слову вектор, выдавая координаты слов на выходе. Сначала он генерирует словарь корпуса, а</p>				

	<p>-приведение символов к одному регистру;          -<i>токенизация</i> – разбиение текста на токены. Так называют отдельные компоненты – слова, предложения или фразы;          -<i>лемматизация</i> – приведение слов к изначальным словоформам, часто с учетом контекста;          -удаление <i>стоп-слов</i> – артиклей, междометий и пр.;</p> <p>После подготовки на выходе получается набор подготовленных слов. Но алгоритмы работают с числовыми данными, поэтому из входящей информации создают <i>векторы</i> – представляют ее как набор числовых значений.</p> <p>На Платформе используются следующие методы векторизации – TF-IDF, и Word2vec.</p>		<p>затем вычисляет векторное представление слов, ‘обучаясь’ на входных тестах. Векторное представление основывается на контекстной близости: слова, встречающиеся в тексте с одинаковыми словами (а следовательно, имеющие схожий смысл), будут иметь близкие (по <i>косинусному расстоянию</i>) векторы. Полученные векторные представления слов используются для обработки естественного языка и машинного обучения.</p>	
<p><b>Кодирование целевого признака</b></p>	<p>Данная функция применяется, когда необходимо преобразовать <i>категориальный целевой признак в датасете</i> в числовое значение. Такое преобразование</p>	<p><b>Использован          ие GPU и          нейронной          сети</b></p>	<p>В данном поле устанавливается галочка, когда обучение модели ИИ происходит с использованием нейронной сети.</p>	<p>Закодированный целевой признак в датасете.</p>

	<p>выполняется перед подачей входных данных в алгоритм. Правила перевода категориальный признаков в числовые прописываются в <i>кодировщике</i>. Данная функция представляет собой первый тип кодирования – <b>Label Encoder</b>. Выполняется <i>порядковое кодирование</i> всех уникальных значений категориального признака: первое (выбранное каким-то образом) уникальное значение кодируется нулем, второе единицей, и так далее, последнее кодируется числом, равным количеству уникальных значений минус единица.</p>		
<p><b>Порядковое кодирование категориальных признаков</b></p>	<p>Отличие данной функции в том, что она выполняет преобразование <i>всех категориальных признаков датасета</i> в числовые значения. При этом кодировщик используется тот же, что и в предыдущей функции – Label Encoder, но кодируются признаки. Выполняется <i>порядковое кодирование</i> каждой категориальной переменной (кроме целевой).</p>	<p>Для данной функции не предусмотрен ручной ввод параметров.</p>	<p>Закодированные категориальные признаки в датасете.</p>



<p><b>Срез временного ряда по индексу</b></p>	<p>Позволяет создавать выборки данных за период времени, используя временные метки или временные диапазоны.</p>	<table border="1"> <tr> <td data-bbox="855 284 1108 384"> <p><b>Дата начала</b></p> </td> <td data-bbox="1108 284 1706 384"> <p>Дата начала среза</p> </td> </tr> <tr> <td data-bbox="855 384 1108 459"> <p><b>Дата окончания</b></p> </td> <td data-bbox="1108 384 1706 459"> <p>Дата окончания среза</p> </td> </tr> </table>	<p><b>Дата начала</b></p>	<p>Дата начала среза</p>	<p><b>Дата окончания</b></p>	<p>Дата окончания среза</p>	<p>Временной ряд после применения фильтра.</p>
<p><b>Дата начала</b></p>	<p>Дата начала среза</p>						
<p><b>Дата окончания</b></p>	<p>Дата окончания среза</p>						
<p><b>Фильтрация текстового шума</b></p>	<p>Данная функция позволяет очистить текст от шумов: из текста убираются знаки препинания, заглавные буквы (они заменяются на строчные) и стоп-слова (различные служебные части речи - союзы, предлоги, частицы и т.д.)</p>	<p>-</p>	<p>Текст без шумов</p>				
<p><b>Лемматизация текста</b></p>	<p>Лемматизация - это процесс приведения всех встречающихся форм слова к одной, нормальной словарной форме. В процессе лемматизации платформа использует словарь и морфологический анализ, чтобы привести слово к его канонической форме – т.н. «лемме», в итоге получается текст, состоящий из слов приведенных к единственному числу, мужскому роду, именительному падежу и инфинитиву (в зависимости от части речи).</p>	<p>-</p>	<p>Нормализованный текст</p>				

<p><b>Векторизация текста</b></p>	<p>Векторизация текста - это процесс преобразования слов в векторы (числа), которые являются «читаемым» форматом для алгоритмов машинного обучения.</p>	<table border="1"> <tr> <td data-bbox="840 284 1093 719"> <p><b>Метод векторизации</b></p> </td> <td data-bbox="1093 284 1686 719"> <p><b>1.TD IDF</b> - метод, используемый для оценки важности слова в контексте документа. Вес некоторого слова пропорционален частоте употребления этого слова в документе и обратно пропорционален частоте употребления слова во всех документах коллекции. <b>2.Word to Vec</b> - данный метод использует контекст, чтобы сформировать численные представления слов, в результате слова, используемые в одном и том же контексте, имеют похожие векторы.</p> </td> </tr> <tr> <td data-bbox="840 719 1093 850"> <p><b>Максимальная размерность текста</b></p> </td> <td data-bbox="1093 719 1686 850"> <p>Указывается примерное количество уникальных слов в тексте.</p> </td> </tr> <tr> <td data-bbox="840 850 1093 981"> <p><b>Количество признаков</b></p> </td> <td data-bbox="1093 850 1686 981"> <p>Указывается количество столбцов таблицы, которая получится в результате преобразования текста в числовой вид</p> </td> </tr> <tr> <td data-bbox="840 981 1093 1142"> <p><b>Сгенерировать тензор для GPU</b></p> </td> <td data-bbox="1093 981 1686 1142"> <p>Выбирается в случае если предполагается что дальше будет использоваться графический процессор. Тензор - это просто таблица особого вида</p> </td> </tr> </table>	<p><b>Метод векторизации</b></p>	<p><b>1.TD IDF</b> - метод, используемый для оценки важности слова в контексте документа. Вес некоторого слова пропорционален частоте употребления этого слова в документе и обратно пропорционален частоте употребления слова во всех документах коллекции. <b>2.Word to Vec</b> - данный метод использует контекст, чтобы сформировать численные представления слов, в результате слова, используемые в одном и том же контексте, имеют похожие векторы.</p>	<p><b>Максимальная размерность текста</b></p>	<p>Указывается примерное количество уникальных слов в тексте.</p>	<p><b>Количество признаков</b></p>	<p>Указывается количество столбцов таблицы, которая получится в результате преобразования текста в числовой вид</p>	<p><b>Сгенерировать тензор для GPU</b></p>	<p>Выбирается в случае если предполагается что дальше будет использоваться графический процессор. Тензор - это просто таблица особого вида</p>	
<p><b>Метод векторизации</b></p>	<p><b>1.TD IDF</b> - метод, используемый для оценки важности слова в контексте документа. Вес некоторого слова пропорционален частоте употребления этого слова в документе и обратно пропорционален частоте употребления слова во всех документах коллекции. <b>2.Word to Vec</b> - данный метод использует контекст, чтобы сформировать численные представления слов, в результате слова, используемые в одном и том же контексте, имеют похожие векторы.</p>										
<p><b>Максимальная размерность текста</b></p>	<p>Указывается примерное количество уникальных слов в тексте.</p>										
<p><b>Количество признаков</b></p>	<p>Указывается количество столбцов таблицы, которая получится в результате преобразования текста в числовой вид</p>										
<p><b>Сгенерировать тензор для GPU</b></p>	<p>Выбирается в случае если предполагается что дальше будет использоваться графический процессор. Тензор - это просто таблица особого вида</p>										
<p><b>Тесты на нормальность распределения</b></p>											
<p><b>Коэффициент асимметрии Skewness</b></p>	<p>Данный метод проверяет выборку на нормальность распределения путем расчета асимметрии данных. Если правый хвост асимметрии длиннее левого,</p>	<table border="1"> <tr> <td data-bbox="840 1241 1048 1337"> <p><b>Признаки</b></p> </td> <td data-bbox="1048 1241 1686 1337"> <p>Выбираются все признаки в датасете для расчета коэффициента асимметрии.</p> </td> </tr> </table>	<p><b>Признаки</b></p>	<p>Выбираются все признаки в датасете для расчета коэффициента асимметрии.</p>	<p>Словарь с данными.</p>						
<p><b>Признаки</b></p>	<p>Выбираются все признаки в датасете для расчета коэффициента асимметрии.</p>										



	то коэффициент положителен, иначе – отрицателен. Если распределение симметрично (в форме ‘колокола’), коэффициент равен нулю.				
<b>Тесты на стационарность временного ряда</b>					
<b>Тест Дики-Фуллера</b>	<p>Проверяется, является ли временной ряд <i>стационарным</i> – не влияют ли на него тренды и сезонность. Для такого ряда суммарные статистические данные согласованы по времени, например, <i>среднее значение</i> и <i>дисперсия наблюдений</i>.</p> <p>Стационарность влияет на легкость моделирования – часто требуется, чтобы временной ряд был стационарным, чтобы быть эффективным.</p>	<table border="1"> <tr> <td><b>Пороговое значение alpha</b></td> <td> <p>Задается пороговое значение <math>p</math> из теста Дики-Фуллера, с использованием которого интерпретируются результаты гипотез:</p> <ul style="list-style-type: none"> <li>• <i>Нулевая гипотеза</i> – временной ряд имеет единичный корень, то есть он нестационарный;</li> <li>• <i>Альтернативная гипотеза</i> – нулевая гипотеза отвергается, и предполагается, что временной ряд не имеет единичного корня, то есть он является стационарным.</li> </ul> <p>Значение <math>p</math> ниже порогового значения означает, что отвергается нулевая гипотеза и временной ряд <i>стационарный</i>. Значение <math>p</math> выше порогового значения означает, что подтверждается нулевая гипотеза и временной ряд <i>нестационарный</i>.</p> <p>Значение <math>p</math> задается в формате числа с плавающей точкой (float).</p> </td> </tr> </table>	<b>Пороговое значение alpha</b>	<p>Задается пороговое значение <math>p</math> из теста Дики-Фуллера, с использованием которого интерпретируются результаты гипотез:</p> <ul style="list-style-type: none"> <li>• <i>Нулевая гипотеза</i> – временной ряд имеет единичный корень, то есть он нестационарный;</li> <li>• <i>Альтернативная гипотеза</i> – нулевая гипотеза отвергается, и предполагается, что временной ряд не имеет единичного корня, то есть он является стационарным.</li> </ul> <p>Значение <math>p</math> ниже порогового значения означает, что отвергается нулевая гипотеза и временной ряд <i>стационарный</i>. Значение <math>p</math> выше порогового значения означает, что подтверждается нулевая гипотеза и временной ряд <i>нестационарный</i>.</p> <p>Значение <math>p</math> задается в формате числа с плавающей точкой (float).</p>	Результаты теста Дики-Фуллера.
<b>Пороговое значение alpha</b>	<p>Задается пороговое значение <math>p</math> из теста Дики-Фуллера, с использованием которого интерпретируются результаты гипотез:</p> <ul style="list-style-type: none"> <li>• <i>Нулевая гипотеза</i> – временной ряд имеет единичный корень, то есть он нестационарный;</li> <li>• <i>Альтернативная гипотеза</i> – нулевая гипотеза отвергается, и предполагается, что временной ряд не имеет единичного корня, то есть он является стационарным.</li> </ul> <p>Значение <math>p</math> ниже порогового значения означает, что отвергается нулевая гипотеза и временной ряд <i>стационарный</i>. Значение <math>p</math> выше порогового значения означает, что подтверждается нулевая гипотеза и временной ряд <i>нестационарный</i>.</p> <p>Значение <math>p</math> задается в формате числа с плавающей точкой (float).</p>				
<b>2.Группа «Машинное обучение»</b>					
<b>Валидация модели</b>	<p>На тестовой выборке данных (обычно это 20% датасета) проверяется правильность работы (предсказательная</p>	<table border="1"> <tr> <td><b>Метрика</b></td> <td> <p>Из списка выбирается название метрики для валидации. Для задачи <i>классификации</i>: Accuracy, F1, Precision, Recall, AUC_ROC.</p> </td> </tr> </table>	<b>Метрика</b>	<p>Из списка выбирается название метрики для валидации. Для задачи <i>классификации</i>: Accuracy, F1, Precision, Recall, AUC_ROC.</p>	Таблица со значением выбранной метрики, отражающей количество правильных ответов обученной модели на тестовой выборке
<b>Метрика</b>	<p>Из списка выбирается название метрики для валидации. Для задачи <i>классификации</i>: Accuracy, F1, Precision, Recall, AUC_ROC.</p>				

	способность) модели ИИ, построенной на основе машинного обучения.		Для задачи <i>регрессии</i> : RMSE, MAE, WMAPE.	данных (максимальное значение метрики равно 1).						
		–	Журнал преобразований над данными.							
		–	Обученная модель ИИ.							
<b>Прогноз модели</b>	Выполняется последовательность действий по прогнозированию будущих значений целевых признаков.	-		1.Точность прогноза. 2.Словарь с данными. 3.Датасет логирования.						
<b>Разделение датасета на обучающую и тестовую выборки</b>	Разделение выборки данных на две категории: для обучения модели ИИ, и для проверки результатов обучения.		<table border="1"> <tr> <td><b>Доля тестовой выборки в датасете</b></td> <td>Обычно на 80% датасета выполняется обучение модели, а на оставшихся 20% – ее валидация. Значение указывается в формате <i>0.2</i>.</td> </tr> <tr> <td><b>Перемешивать наблюдения перед разделением</b></td> <td>Необходимо установить галочку в поле, чтобы случайным образом поменять местами наблюдения в датасете. Не рекомендуется перемешивать временные ряды, т.к. наблюдения в них упорядочены и зафиксированы последовательно по времени.</td> </tr> <tr> <td><b>Разделять с учетом меток классов</b></td> <td>Выбирается, учитывать ли долю таргетов при разделении датасета. Используется только для задач классификации, когда объекты распределяются по категориям согласно определенным и заданным заранее признакам.</td> </tr> </table>	<b>Доля тестовой выборки в датасете</b>	Обычно на 80% датасета выполняется обучение модели, а на оставшихся 20% – ее валидация. Значение указывается в формате <i>0.2</i> .	<b>Перемешивать наблюдения перед разделением</b>	Необходимо установить галочку в поле, чтобы случайным образом поменять местами наблюдения в датасете. Не рекомендуется перемешивать временные ряды, т.к. наблюдения в них упорядочены и зафиксированы последовательно по времени.	<b>Разделять с учетом меток классов</b>	Выбирается, учитывать ли долю таргетов при разделении датасета. Используется только для задач классификации, когда объекты распределяются по категориям согласно определенным и заданным заранее признакам.	1.Отдельно обучающая и тестовая выборки. 2.Журнал преобразований.
<b>Доля тестовой выборки в датасете</b>	Обычно на 80% датасета выполняется обучение модели, а на оставшихся 20% – ее валидация. Значение указывается в формате <i>0.2</i> .									
<b>Перемешивать наблюдения перед разделением</b>	Необходимо установить галочку в поле, чтобы случайным образом поменять местами наблюдения в датасете. Не рекомендуется перемешивать временные ряды, т.к. наблюдения в них упорядочены и зафиксированы последовательно по времени.									
<b>Разделять с учетом меток классов</b>	Выбирается, учитывать ли долю таргетов при разделении датасета. Используется только для задач классификации, когда объекты распределяются по категориям согласно определенным и заданным заранее признакам.									

**Классификация** решает задачу разделения множества наблюдений (объектов) на группы, называемые *классами*, на основе анализа их формального описания. При классификации каждое наблюдение относится к определенной группе на основе некоторого качественного свойства. Пусть **X** – множество описаний

объектов,  $Y$  – конечное множество номеров/имен/меток классов. Существует неизвестная целевая зависимость отображения  $y^*: X \rightarrow Y$ , значения которой известны только на объектах обучающей выборки  $X^m = (x_1, y_1), \dots, (x_m, y_m)$ . Строится алгоритм, способный классифицировать произвольный объект  $x \in X$ .

<p><b>Логистическая Регрессия</b></p>	<p>Используется логистическая функция для моделирования зависимости выходной переменной <math>y</math> от набора входных переменных <math>x</math>, в случае когда первая является бинарной. Например, с помощью логистической регрессии можно оценивать вероятность наступления/или не наступления некоторого события. Предсказывается непрерывная переменная – коэффициент логистической регрессии, принимающий значение от 0 до 1:</p> <ul style="list-style-type: none"> <li>• если значение коэффициента больше порогового значения, то вероятность наступления события равна 1;</li> <li>• иначе вероятность наступления события равна 0.</li> </ul>	<table border="1"> <tr> <td data-bbox="833 432 1117 667"> <p><b>Коэффициент регуляризации</b></p> </td> <td data-bbox="1117 432 1688 667"> <p>Указывается значение строго больше нуля – положительное вещественное число, с помощью которого добавляется дополнительное ограничение к условию с целью предотвратить переобучение модели.</p> </td> </tr> <tr> <td data-bbox="833 667 1117 799"> <p><b>Порог классификации</b></p> </td> <td data-bbox="1117 667 1688 799"> <p>Значение вещественного типа от 0 до 1, определяющее принадлежность объекта к тому или иному классу.</p> </td> </tr> <tr> <td data-bbox="833 799 1117 1099"> <p><b>Флаг возврата вероятности при прогнозе</b></p> </td> <td data-bbox="1117 799 1688 1099"> <p>Возвращает вероятность или метки классов для дальнейшего прогноза после обучения модели. Используется для решения задач бинарной классификации, когда выходная переменная может принимать только два значения – решается вопрос о принадлежности объекта к одному из двух классов.</p> </td> </tr> <tr> <td data-bbox="833 1099 1117 1232"> <p><b>Оптимизация гиперпараметров</b></p> </td> <td data-bbox="1117 1099 1688 1232"> <p>Флаг подбора гиперпараметров. Флаг активируется, когда указывается несколько гиперпараметров.</p> </td> </tr> <tr> <td data-bbox="833 1232 1117 1396"> <p><b>Метрика для оптимизации</b></p> </td> <td data-bbox="1117 1232 1688 1396"> <p>Критерий остановки итераций. Настройка, позволяющая определить точность нахождения минимума функции ошибки.</p> </td> </tr> </table>	<p><b>Коэффициент регуляризации</b></p>	<p>Указывается значение строго больше нуля – положительное вещественное число, с помощью которого добавляется дополнительное ограничение к условию с целью предотвратить переобучение модели.</p>	<p><b>Порог классификации</b></p>	<p>Значение вещественного типа от 0 до 1, определяющее принадлежность объекта к тому или иному классу.</p>	<p><b>Флаг возврата вероятности при прогнозе</b></p>	<p>Возвращает вероятность или метки классов для дальнейшего прогноза после обучения модели. Используется для решения задач бинарной классификации, когда выходная переменная может принимать только два значения – решается вопрос о принадлежности объекта к одному из двух классов.</p>	<p><b>Оптимизация гиперпараметров</b></p>	<p>Флаг подбора гиперпараметров. Флаг активируется, когда указывается несколько гиперпараметров.</p>	<p><b>Метрика для оптимизации</b></p>	<p>Критерий остановки итераций. Настройка, позволяющая определить точность нахождения минимума функции ошибки.</p>	<p>1. Модель бинарной классификации.                  2. Словарь с данными.                  3. Точность модели.                  4. Журнал преобразований.</p>
<p><b>Коэффициент регуляризации</b></p>	<p>Указывается значение строго больше нуля – положительное вещественное число, с помощью которого добавляется дополнительное ограничение к условию с целью предотвратить переобучение модели.</p>												
<p><b>Порог классификации</b></p>	<p>Значение вещественного типа от 0 до 1, определяющее принадлежность объекта к тому или иному классу.</p>												
<p><b>Флаг возврата вероятности при прогнозе</b></p>	<p>Возвращает вероятность или метки классов для дальнейшего прогноза после обучения модели. Используется для решения задач бинарной классификации, когда выходная переменная может принимать только два значения – решается вопрос о принадлежности объекта к одному из двух классов.</p>												
<p><b>Оптимизация гиперпараметров</b></p>	<p>Флаг подбора гиперпараметров. Флаг активируется, когда указывается несколько гиперпараметров.</p>												
<p><b>Метрика для оптимизации</b></p>	<p>Критерий остановки итераций. Настройка, позволяющая определить точность нахождения минимума функции ошибки.</p>												

		<p><b>Количество фолдов для оптимизации</b></p>	<p>Датасет делится на фолды – на указанное количество равных частей. При обучении модели каждый фолд становится валидационным один раз, при этом на остальных фолдах выполняется обучение. Каждый раз рассчитывается значение метрики. Затем рассчитывается <i>усредненная метрика</i>, которая характеризует точность модели.</p>	
<p><b>Модель Classifier</b></p>	<p><b>XGB</b> Алгоритм <b>XGBClassifier</b> анализирует связь между признаками и целевым признаком. На обучающей выборке модель обучается соотносить наблюдение к аномалиям, а на тестовой выборке выполняется валидация ответов обученной модели.</p>	<p><b>Глубина дерева</b></p>	<p>Заданное максимальное число разбиений в ветвях, по достижению которого обучение модели ИИ останавливается.</p> <p><b>Количество базовых моделей</b></p> <p>Определяет сколько независимых моделей будет работать над обучением.</p> <p><b>Порог классификации</b></p> <p>Значение от 0 до 1, указывающее на верхнюю границу вероятности причисления объекта к классу.</p> <p><b>Флаг возврата вероятности при прогнозе</b></p> <p>Возвращает вероятность или метки классов для дальнейшего прогноза после обучения модели. Используется для решения задач бинарной классификации, когда выходная переменная может принимать только два значения – решается вопрос о принадлежности объекта к одному из двух классов.</p>	<p>1. Таблица с матрицей ошибок. 2. Таблица верных и ошибочных прогнозов модели в разрезе классов. 3. Модель бинарной классификации.</p>

		<p><b>Оптимизация гиперпараметров</b></p>	<p>Флаг подбора гиперпараметров. Флаг активируется, когда указывается несколько гиперпараметров.</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Критерий остановки итераций. Настройка, позволяющая определить точность нахождения минимума функции ошибки.</p>	
		<p><b>Количество фолдов для оптимизации.</b></p>	<p>Датасет делится на фолды – на указанное количество равных частей. При обучении модели каждый фолд становится валидационным один раз, при этом на остальных фолдах выполняется обучение. Каждый раз рассчитывается значение метрики. Затем рассчитывается усредненная метрика, которая характеризует точность модели.</p>	
<p><b>Дерево решений для классификации</b></p>	<p>Предсказывает, к какому классу принадлежит объект из обучающего массива данных. Для этого строится дерево решений: древовидная структура, где моменты принятий решений соответствуют узлам, в узлах происходит ветвление процесса на ветки в зависимости от сделанного выбора, и конечные узлы (листья) – конечные результаты последовательного</p>	<p><b>Глубина дерева</b></p>	<p>Заданное максимальное число разбиений в ветвях, по достижению которого обучение модели ИИ останавливается.</p>	<p>1.Датасет с меткой класса, определяющей принадлежность объекта к одному из классов. 2.Модель ИИ, обученная классифицировать данные по заданным критериям. 3.Журнал преобразований над данными. 4.Словарь с переменными (описание модели, таблицы, графики) для отображения в интерфейсе Программы.</p>
		<p><b>Порог классификации</b></p>	<p>Значение порога определяет принадлежность объекта к одному из классов: к положительному – если порог выше указанного значения, к отрицательному – если порог ниже.</p>	
		<p><b>Флаг возврата вероятности при прогнозе</b></p>	<p>Возвращает вероятность или метки классов для дальнейшего прогноза после обучения модели.</p>	

	<p>принятия решений. В узлах, начиная с корневого, выбирается признак, значение которого используется для разбиения всех данных на два класса. Процесс продолжается до тех пор, пока не выполнится <i>критерий остановки</i> – дерево превысило заранее заданный «лимит роста» (достигнута глубина дерева). При этом разбиения выполняются таким образом, чтобы уменьшить выбранный критерий, например <i>энтропию</i> – степень неопределенности в разбиении на классы.</p>	<p><b>Оптимизация гиперпараметров</b></p>	<p>При выборе флага оптимизации не нужно вручную задавать глубину дерева, или можно задать несколько значений на выбор. Алгоритм подбирает глубину дерева из расчета получения максимального значения метрики.</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Метод, который рассчитывает точность обученной модели. Выбирается один из предлагаемых методов.</p>	
		<p><b>Количество фолдов для оптимизации</b></p>	<p>Датасет делится на фолды – на указанное количество равных частей. При обучении модели каждый фолд становится <i>валидационным</i> один раз, при этом на остальных фолдах выполняется обучение. Каждый раз рассчитывается значение метрики. Затем рассчитывается <i>усредненная метрика</i>, которая характеризует точность модели.</p>	
<p><b>Случайный лес для классификации</b></p>	<p>Строится множество решающих деревьев, и в задаче классификации принимается решение голосованием по большинству. Все деревья строятся независимо по схеме: 1.Выбирается подвыборка обучающей выборки и по ней строится дерево.</p>	<p><b>Глубина дерева</b></p>	<p>Максимальная глубина для деревьев решений.</p>	<p>**</p>
		<p><b>Количество деревьев</b></p>	<p>Число деревьев в «лесу».</p>	
		<p><b>Порог классификации</b></p>	<p>–</p>	

	<p>2.Для построения каждого расщепления в дереве просматривается максимальное количество случайных признаков.</p> <p>3.Выбирается наилучший признак и расщепление по нему (по заранее заданному критерию). Дерево строится, до достижения параметра, ограничивающего его высоту.</p> <p>Таким образом деревья обучаются не только на разных наборах данных, но и используют разные признаки для принятия решений – это создает некоррелированные деревья, которые и защищают друг друга от своих ошибок. Прогноз получается точнее, чем у любого отдельного дерева.</p>	<p><b>Флаг возврата вероятности при прогнозе</b></p>	<p>–</p>	
		<p><b>Оптимизация гиперпараметров</b></p>	<p>Необходимо активировать галочку в поле, чтобы подобрать гиперпараметры – глубину и количество деревьев. Гиперпараметры подбираются таким образом, чтобы получить максимальное значение метрики.</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Выбирается одна из предлагаемых метрик для оценки работы модели.</p>	
		<p><b>Количество фолдов для оптимизации</b></p>	<p>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</p>	
<p><b>Categorical Naive Bayes</b></p>	<p>Группа байесовских классификаторов позволяет определить к какому классу принадлежит объект на основе теоремы Байеса с допущением о независимости признаков. Категориальный наивный байесовский классификатор применяется для признаков</p>	<p><b>Параметр сглаживания Лапласа</b></p>	<p>Значение для аддитивного сглаживания Лапласа во избежание проблемы нулевой вероятности. Если равен нулю, то сглаживания нет.</p>	<p>1.Датасет с меткой класса, определяющей принадлежность объекта к одному из классов.</p> <p>2.Модель ИИ, обученная классифицировать данные по заданным критериям.</p> <p>3.Журнал преобразований над данными.</p> <p>4.Словарь с переменными (описание модели, таблицы,</p>
		<p><b>Априорные вероятности классов</b></p>	<p>Определяет, будут ли взяты в расчет априорные вероятности классов. Если не активирован, то применяются значения вероятностей для равномерного</p>	

	<p>с категориальным распределением.</p>		<p>распределения.</p>	<p>графики) для отображения в интерфейсе Программы.</p>
<p><b>Оптимизация гиперпараметров</b></p>	<p>Определяет, будет ли проводиться оптимизация гиперпараметров при обучении модели</p>	<p><b>Метрика для оптимизации</b></p>	<p>Выбирается одна из предлагаемых метрик для оценки работы модели.</p>	
<p><b>Количество фолдов для оптимизации</b></p>	<p>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</p>			
<p><b>Multinomial Naive Bayes</b></p>	<p>Мультиномиальный классификатор применяется для признаков с полиномиальным распределением. Пример: классификация текстов, где каждый текст представлен вектором слов (например, мешок слов или tf-idf).</p>	<p><b>Параметр сглаживания Лапласа</b></p>	<p>Значение для аддитивного сглаживания Лапласа во избежание проблемы нулевой вероятности. Если равен нулю, то сглаживания нет.</p>	
<p><b>Априорные вероятности классов</b></p>	<p>Определяет, будут ли взяты в расчет априорные вероятности классов. Если не активирован, то применяются значения вероятностей для равномерного распределения.</p>	<p><b>Оптимизация гиперпараметров</b></p>	<p>Определяет, будет ли проводиться оптимизация гиперпараметров при обучении модели</p>	
<p><b>Метрика для оптимизации</b></p>	<p>Выбирается одна из предлагаемых метрик для оценки работы модели.</p>			



		<table border="1"> <tr> <td><b>Количество фолдов для оптимизации</b></td> <td>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</td> </tr> </table>	<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.									
<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.												
<b>Complement Naive Bayes</b>	<p>Представляет собой вариант адаптации Multinomial Naive Bayes для датасетов с несбалансированными классами.</p> <p>Вместо вычисления вероятностей принадлежности объекта к конкретному классу для каждого класса вычисляются вероятности того, что объект им не принадлежит.</p> <p>Выбирается наименьшая вероятность "непринадлежности" к классу, так как это означает, что объект с наибольшей вероятностью принадлежит к данному классу.</p>	<table border="1"> <tr> <td><b>Параметр сглаживания Лапласа</b></td> <td>Значение для аддитивного сглаживания Лапласа во избежание проблемы нулевой вероятности. Если равен нулю, то сглаживания нет.</td> </tr> <tr> <td><b>Априорные вероятности классов</b></td> <td>Определяет, будут ли взяты в расчет априорные вероятности классов. Если не активирован, то применяются значения вероятностей для равномерного распределения.</td> </tr> <tr> <td><b>Оптимизация гиперпараметров</b></td> <td>Определяет, будет ли проводиться оптимизация гиперпараметров при обучении модели</td> </tr> <tr> <td><b>Метрика для оптимизации</b></td> <td>Выбирается одна из предлагаемых метрик для оценки работы модели.</td> </tr> <tr> <td><b>Количество фолдов для оптимизации</b></td> <td>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</td> </tr> </table>	<b>Параметр сглаживания Лапласа</b>	Значение для аддитивного сглаживания Лапласа во избежание проблемы нулевой вероятности. Если равен нулю, то сглаживания нет.	<b>Априорные вероятности классов</b>	Определяет, будут ли взяты в расчет априорные вероятности классов. Если не активирован, то применяются значения вероятностей для равномерного распределения.	<b>Оптимизация гиперпараметров</b>	Определяет, будет ли проводиться оптимизация гиперпараметров при обучении модели	<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.	<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.	<ol style="list-style-type: none"> <li>1.Датасет с <i>меткой класса</i>, определяющей принадлежность объекта к одному из классов.</li> <li>2.Модель ИИ, обученная классифицировать данные по заданным критериям.</li> <li>3.Журнал преобразований над данными.</li> <li>4.Словарь с переменными (описание модели, таблицы, графики) для отображения в интерфейсе Программы.</li> </ol>
<b>Параметр сглаживания Лапласа</b>	Значение для аддитивного сглаживания Лапласа во избежание проблемы нулевой вероятности. Если равен нулю, то сглаживания нет.												
<b>Априорные вероятности классов</b>	Определяет, будут ли взяты в расчет априорные вероятности классов. Если не активирован, то применяются значения вероятностей для равномерного распределения.												
<b>Оптимизация гиперпараметров</b>	Определяет, будет ли проводиться оптимизация гиперпараметров при обучении модели												
<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.												
<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.												
<b>Gaussian Naive Bayes</b>	Для значений признаков для каждого класса строится распределение Гаусса. В качестве значений		<ol style="list-style-type: none"> <li>1.Датасет с <i>меткой класса</i>, определяющей принадлежность объекта к одному из классов.</li> </ol>										

	<p>правдоподобия для признаков берутся значения функции Гаусса из конкретного распределения (соответствующее признаку и классу).</p>	<p><b>Параметр сглаживания Лапласа</b></p>	<p>Значение для аддитивного сглаживания Лапласа во избежание проблемы нулевой вероятности. Если равен нулю, то сглаживания нет.</p>	<p>2. Модель ИИ, обученная классифицировать данные по заданным критериям. 3. Журнал преобразований над данными. 4. Словарь с переменными (описание модели, таблицы, графики) для отображения в интерфейсе Программы.</p>
		<p><b>Априорные вероятности классов</b></p>	<p>Определяет, будут ли взяты в расчет априорные вероятности классов. Если не активирован, то применяются значения вероятностей для равномерного распределения.</p>	
		<p><b>Оптимизация гиперпараметров</b></p>	<p>Определяет, будет ли проводиться оптимизация гиперпараметров при обучении модели</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Выбирается одна из предлагаемых метрик для оценки работы модели.</p>	
		<p><b>Количество фолдов для оптимизации</b></p>	<p>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</p>	
<p><b>Bernoulli Naive Bayes</b></p>	<p>Применяется для признаков с биномиальным распределением. Пример: классификация текстов, где каждый текст представлен вектором наличия слов из словаря (1 - есть слово, 0 - нет).</p>	<p><b>Параметр сглаживания Лапласа</b></p>	<p>Значение для аддитивного сглаживания Лапласа во избежание проблемы нулевой вероятности. Если равен нулю, то сглаживания нет.</p>	<p>1. Датасет с меткой класса, определяющей принадлежность объекта к одному из классов. 2. Модель ИИ, обученная классифицировать данные по заданным критериям. 3. Журнал преобразований над данными. 4. Словарь с переменными (описание модели, таблицы,</p>
		<p><b>Априорные вероятности классов</b></p>	<p>Определяет, будут ли взяты в расчет априорные вероятности классов. Если не активирован, то применяются значения</p>	

			вероятностей для равномерного распределения.	графики) для отображения в интерфейсе Программы.
		<b>Оптимизация гиперпараметров</b>	Определяет, будет ли проводиться оптимизация гиперпараметров при обучении модели	
		<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.	
		<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.	

**Кластеризация** – это задача группировки множества объектов на подмножества (кластеры) так, чтобы объекты одного кластера были более похожи друг на друга, чем на объекты других кластеров по какому-либо критерию. Относится к классу задач *обучения без учителя*.

<p><b>Алгоритм кластеризации DBSCAN</b></p>	<p>Алгоритм <b>DBScan</b> формирует группы коренных соседей/кластеры, объединяя точки, расположенные рядом. А точки, которые не попадают ни в одну из групп, отмечаются меткой -1 и приравниваются к аномалиям.</p>	<p><b>Датасет</b></p>	<p>Датасет с исходными данными.</p>	<p>1. Модель кластеризации.                  2. Выходной датасет, дополненный меткой кластера и/или флагом аномалии.                  3. Словарь, содержащий информацию (графики, таблицы, текст) для отображения в пользовательском интерфейсе.                  4. Журнал преобразований.                  **Для алгоритмов <i>кластеризации, регрессии, классификации</i> наборы выходных параметров идентичны, отличие заключается в их содержимом.</p>
		<p><b>Журнал преобразований</b></p>	<p>—</p>	
		<p><b>Радиус</b></p>	<p>Радиус в единицах расстояния, в рамках которого выполняется поиск потенциальных соседей (float/list/tuple).</p>	
		<p><b>Число соседей</b></p>	<p>Минимальное число ближайших соседей в указанном радиусе для формирования группы коренных соседей (int/list/tuple).</p>	
		<p><b>Метрика расстояния</b></p>	<p>Метрика расстояния (str/list): расстояние Евклида, косинусное расстояние. По умолчанию «Евклидово расстояние» — используется при кластеризации данных в текущем датасете, а также при отнесении нового объекта к кластеру.</p>	
		<p><b>Оптимизация гиперпараметров</b></p>	<p>Флаг подбора гиперпараметров. При значении «true» выполняется ручной ввод следующих гиперпараметров: радиус, число соседей, метрика расстояния. При значении «false» эти гиперпараметры подбираются автоматически.</p>	
		<p>* Параметры: <i>датасет</i> и <i>журнал преобразований</i> являются входными и выходными параметрами для всех алгоритмов.</p>		

<p><b>Изоляционный лес</b></p>	<p>Алгоритм поиска <i>аномалий (выбросов)</i> методом «Изоляционный лес»: Изолирует наблюдения, случайным образом выбирая объект, а затем случайным образом выбирая разделения между максимальным и минимальным значениями объекта. Разбиение представлено древовидной структурой, количество разбиений, необходимое для изоляции выборки, равно длине пути от корневого до конечного узла. Эта длина пути является мерой нормальности и функции принятия решений. Когда лес случайных деревьев создает более короткие пути для отдельных объектов, они, скорее всего, являются аномалиями.</p>	<table border="1"> <tr> <td data-bbox="846 280 1189 357"><b>Датасет</b></td> <td data-bbox="1189 280 1704 357">Датасет с исходными данными.</td> </tr> <tr> <td data-bbox="846 357 1189 456"><b>Журнал преобразований</b></td> <td data-bbox="1189 357 1704 456">–</td> </tr> <tr> <td data-bbox="846 456 1189 555"><b>Количество деревьев</b></td> <td data-bbox="1189 456 1704 555">По умолчанию устанавливается значение, равное 2</td> </tr> </table>	<b>Датасет</b>	Датасет с исходными данными.	<b>Журнал преобразований</b>	–	<b>Количество деревьев</b>	По умолчанию устанавливается значение, равное 2	<p>**</p>
<b>Датасет</b>	Датасет с исходными данными.								
<b>Журнал преобразований</b>	–								
<b>Количество деревьев</b>	По умолчанию устанавливается значение, равное 2								
<p><b>Кластеризация K-Means</b></p>	<p>Алгоритм кластеризации K-средних:                  1.Из исходного множества случайным образом выбирается K наблюдений, равное заданному количеству кластеров.                  2.Для каждого наблюдения определяется ближайший к нему центр кластера</p>	<table border="1"> <tr> <td data-bbox="846 1131 1133 1208"><b>Датасет</b></td> <td data-bbox="1133 1131 1704 1208">Датасет с исходными данными.</td> </tr> <tr> <td data-bbox="846 1208 1133 1307"><b>Журнал преобразований</b></td> <td data-bbox="1133 1208 1704 1307">–</td> </tr> <tr> <td data-bbox="846 1307 1133 1398"><b>Число кластеров</b></td> <td data-bbox="1133 1307 1704 1398">По умолчанию устанавливается значение, равное 2.</td> </tr> </table>	<b>Датасет</b>	Датасет с исходными данными.	<b>Журнал преобразований</b>	–	<b>Число кластеров</b>	По умолчанию устанавливается значение, равное 2.	<p>**</p>
<b>Датасет</b>	Датасет с исходными данными.								
<b>Журнал преобразований</b>	–								
<b>Число кластеров</b>	По умолчанию устанавливается значение, равное 2.								

	<p>(измеряется Евклидоваго расстояние до центра). Образуются начальные кластеры. 3.Вычисляются <i>центры тяжести кластеров</i> – вектора, элементы которых представляют собой среднее арифметическое значение признаков кластера. 4.Центры кластеров смещаются и объединяют вокруг себя наблюдения, пока центры и границы кластеров не перестанут изменяться.</p>	<table border="1"> <tr> <td data-bbox="846 252 1133 339"><b>Оптимизация гиперпараметров</b></td> <td data-bbox="1133 252 1697 339">Флаг подбора гиперпараметров.</td> </tr> </table>	<b>Оптимизация гиперпараметров</b>	Флаг подбора гиперпараметров.							
<b>Оптимизация гиперпараметров</b>	Флаг подбора гиперпараметров.										
<p><b>Агломеративная иерархическая кластеризация</b></p>	<p>Последовательно объединяет объекты во все более крупные подмножества, в результате образуется древовидная структура. Отдельные версии иерархии отличаются правилами вычисления расстояния между кластерами. Например, алгоритм средней связи на каждом шаге объединяет два ближайших кластера, рассчитывая среднюю арифметическую</p>	<table border="1"> <tr> <td data-bbox="846 882 1160 954"><b>Датасет</b></td> <td data-bbox="1160 882 1697 954">Датасет с исходными данными.</td> </tr> <tr> <td data-bbox="846 962 1160 1050"><b>Журнал преобразований</b></td> <td data-bbox="1160 962 1697 1050">–</td> </tr> <tr> <td data-bbox="846 1058 1160 1145"><b>Число кластеров</b></td> <td data-bbox="1160 1058 1697 1145">Задается оптимальное количество кластеров.</td> </tr> <tr> <td data-bbox="846 1153 1160 1374"><b>Метрика расчета расстояния</b></td> <td data-bbox="1160 1153 1697 1374">Используется при кластеризации данных в текущем датасете, а также при отнесении нового объекта к кластеру. Значения: евклидово расстояние, косинусная мера, расстояние городских кварталов, расстояние Чебышева.</td> </tr> </table>	<b>Датасет</b>	Датасет с исходными данными.	<b>Журнал преобразований</b>	–	<b>Число кластеров</b>	Задается оптимальное количество кластеров.	<b>Метрика расчета расстояния</b>	Используется при кластеризации данных в текущем датасете, а также при отнесении нового объекта к кластеру. Значения: евклидово расстояние, косинусная мера, расстояние городских кварталов, расстояние Чебышева.	<p>**</p>
<b>Датасет</b>	Датасет с исходными данными.										
<b>Журнал преобразований</b>	–										
<b>Число кластеров</b>	Задается оптимальное количество кластеров.										
<b>Метрика расчета расстояния</b>	Используется при кластеризации данных в текущем датасете, а также при отнесении нового объекта к кластеру. Значения: евклидово расстояние, косинусная мера, расстояние городских кварталов, расстояние Чебышева.										

	<p>дистанцию между всеми парами объектов.</p>	<p><b>Критерий связи для расчета расстояния</b></p>	<p>Правила вычисления расстояния между кластерами при каждой итерации. Значения: -алгоритм средней связи, -алгоритм одиночной связи или ближайшего соседа, -алгоритм полной связи или дальнего соседа, -метод минимума дисперсии Уорда. <i>Дисперсия</i> объединяет кластеры с минимальной общей внутрикластерной дисперсией после слияния, в качестве метрики расстояния используется евклидово расстояние. Минимальный использует самые близкие точки в обоих кластерах для расчета расстояния.</p>			
<p><b>Метод локтя K-Means</b></p>	<p>Метод локтя позволяет вычислить правильное значение <math>k</math> (количество кластеров) и повысить производительность модели. Вычисляется сумма квадратов расстояний между точками, и среднее арифметическое значение (Mean) – сумма элементов датасета, разделенная на их количество. Когда значение</p>	<p><b>Оптимизация гиперпараметров</b></p>	<p>Флаг подбора гиперпараметров.</p>	<p>**</p>		
		<table border="1"> <tr> <td data-bbox="853 1046 1099 1145"> <p><b>Число кластеров</b></p> </td> <td data-bbox="1099 1046 1691 1145"> <p>Задается оптимальное количество кластеров.</p> </td> </tr> </table>	<p><b>Число кластеров</b></p>	<p>Задается оптимальное количество кластеров.</p>		
<p><b>Число кластеров</b></p>	<p>Задается оптимальное количество кластеров.</p>					

	<p>к равно 1, сумма квадрата внутри кластера будет большой. По мере увеличения значения к сумма квадратов внутри кластера будет уменьшаться. Наконец будет построен график между значениями k и суммой квадрата внутри кластера. В момент, когда значение k резко уменьшится будет считаться оптимальным числом кластеров.</p>								
<p><b>Регрессия</b> – математическое выражение, отражающее связь между зависимой переменной y и независимыми переменными x. Алгоритмы регрессии используются для <i>контролируемого обучения</i> моделей ИИ – так называемого обучения «с учителем», когда данные размечаются для помощи в прогнозировании. Сопоставляя входные данные и полученные результаты на точность, модель постепенно обучается прогнозировать <i>числовые значения</i> целевых переменных.</p>									
<p><b>Линейная регрессия</b></p>	<p>Прогнозирует целевую переменную Y на основе одной или нескольких независимых переменных X. Для этого между X и Y строится линейная связь.</p>	<table border="1"> <tr> <td data-bbox="837 938 1178 1015"><b>Датасет</b></td> <td data-bbox="1178 938 1729 1015">Датасет с исходными данными.</td> </tr> <tr> <td data-bbox="837 1015 1178 1110"><b>Журнал преобразований</b></td> <td data-bbox="1178 1015 1729 1110">–</td> </tr> </table>	<b>Датасет</b>	Датасет с исходными данными.	<b>Журнал преобразований</b>	–	<p>**</p>		
<b>Датасет</b>	Датасет с исходными данными.								
<b>Журнал преобразований</b>	–								
<p><b>Дерево решений для регрессии</b></p>	<p>Предсказывает значение целевой переменной, изучая простые правила принятия решений, выведенные из характеристик данных. Представляет собой древовидный граф с узлами, где атрибут – вопрос, ребро – ответ на вопрос, а</p>	<table border="1"> <tr> <td data-bbox="837 1149 1178 1225"><b>Датасет</b></td> <td data-bbox="1178 1149 1729 1225">Датасет с исходными данными.</td> </tr> <tr> <td data-bbox="837 1225 1178 1321"><b>Журнал преобразований</b></td> <td data-bbox="1178 1225 1729 1321">–</td> </tr> <tr> <td data-bbox="837 1321 1178 1417"><b>Глубина дерева</b></td> <td data-bbox="1178 1321 1729 1417">Заданное максимальное число разбиений в ветвях, по достижению</td> </tr> </table>	<b>Датасет</b>	Датасет с исходными данными.	<b>Журнал преобразований</b>	–	<b>Глубина дерева</b>	Заданное максимальное число разбиений в ветвях, по достижению	<p>**</p>
<b>Датасет</b>	Датасет с исходными данными.								
<b>Журнал преобразований</b>	–								
<b>Глубина дерева</b>	Заданное максимальное число разбиений в ветвях, по достижению								



	<p>листья – фактический результат. Наблюдения классифицируются сверху вниз от корня до листьев.</p>		<p>которого обучение останавливается.</p>	
		<p><b>Оптимизация гиперпараметров</b></p>	<p>При выборе флага оптимизации не нужно вручную задавать глубину дерева, или можно задать несколько значений на выбор. Алгоритм подбирает глубину дерева из расчета получить максимальное значение метрики.</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Метод, который рассчитывает точность обученной модели. Выбирается один из предлагаемых методов.</p>	
		<p><b>Количество фолдов для оптимизации</b></p>	<p>Датасет делится на фолды – на указанное количество равных частей. При обучении модели регрессии каждый фолд становится валидационным один раз, при этом на остальных фолдах выполняется обучение. Каждый раз рассчитывается значение метрики. Затем рассчитывается <i>усредненная метрика</i>, которая характеризует точность модели.</p>	
<p><b>Случайный лес для регрессии</b></p>	<p>В отличие от предыдущего алгоритма здесь строится ансамбль решающих деревьев. При этом большое количество некоррелированных моделей (деревьев)</p>	<p><b>Датасет</b></p>	<p>Датасет с исходными данными.</p>	<p>**</p>
		<p><b>Журнал преобразований</b></p>	<p>–</p>	
		<p><b>Глубина дерева</b></p>	<p>–</p>	

	<p>превосходит любую из отдельных моделей.</p>	<p><b>Количество деревьев</b></p>	<p>–</p>	
		<p><b>Оптимизация гиперпараметров</b></p>	<p>Алгоритм подбирает гиперпараметры: глубину и количество деревьев из расчета получить максимальное значение метрики.</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Выбирается одна из предлагаемых метрик для оценки работы модели.</p>	
		<p><b>Количество фолдов для оптимизации</b></p>	<p>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</p>	
<p><b>Полиномиальная регрессия</b></p>	<p>Метод регрессионного анализа, в которой взаимосвязь между независимыми переменными <math>x</math> и зависимой переменной <math>y</math> моделируется как полином <math>n</math>-ой степени от <math>x</math>. Полиномиальная регрессия соответствует нелинейной зависимости между значением <math>x</math> и соответствующим условным средним <math>y</math>, обозначаемым <math>E(y x)</math>. В отличие от линейной регрессии моделирует нелинейно разделенные данные – более гибкая и может моделировать сложные взаимосвязи.</p>	<p><b>Степень полинома</b></p>	<p>Степень уравнения полиномиальной регрессии, которая определяет линию наилучшего соответствия. При неправильном выборе степени, модель может быть перенасыщена. Значение по умолчанию – 2.</p>	<p>1.Модель полиномиальной регрессии. 2.Словарь с переменными для отображения в интерфейсе. 3.Словарь с преобразованиями данных. 4.Выходной датасет.</p>
		<p><b>Только произведение</b></p>	<p>Если установить галочку в поле, то не выполняется возведение в степень, а только перемножение.</p>	
		<p><b>Оптимизация гиперпараметров</b></p>	<p>Нужно активировать галочку в поле, когда выбирается наиболее подходящая степень полинома из нескольких предложенных. А подбирается гиперпараметр так, чтобы получить наилучшее значение метрики.</p>	

		<p><b>Метрика для оптимизации</b></p>	<p>Значения на выбор: RMSE, MAE, WMAPE, где RMSE – среднеквадратическая ошибка, MAE – средняя абсолютная ошибка, а ошибка – разница между значениями, предсказанными моделью, и фактическими значениями переменной. Эти метрики используются для оценки работы модели регрессии – проверяют точность прогноза и измеряют величину отклонения от фактических значений.</p>	
<p><b>Метод опорных векторов для регрессии</b></p>	<p>В основе регрессии опорных векторов (с англ. <b>SVR</b> – Support Vector Regression) лежит поиск гиперплоскости, при которой риск в многомерном пространстве будет минимальным. SVR оценивает коэффициенты путем минимизации квадратичных потерь: считается сумма квадратов ошибок (между прогнозом и фактом), и к ней прибавляется штраф в виде произведения <i>коэффициента регуляризации</i> и суммы квадратов весов.</p>	<p><b>Количество фолдов для оптимизации</b></p>	<p>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</p>	<p>**</p>
		<p><b>Тип ядра</b></p>	<p>Функция ядра (<b>kernel</b>) может принимать значения: {'linear', 'poly', 'rbf', 'sigmoid'}.</p>	
		<p><b>Степень для ядра полинома</b></p>	<p>Если в качестве функции ядра используется полиномиальная функция ('poly'), которая является методом нелинейной регрессии, то зависимая переменная связана с независимыми переменными n-ой степени. В поле указывается степень этого ядра.</p>	
		<p><b>Коэффициент регуляризации</b></p>	<p>Мера степени наказания модели за каждую неверно спрогнозированную точку.</p>	
		<p><b>Оптимизация</b></p>	<p>Флаг подбора гиперпараметров.</p>	

	<p>*Вместо квадратичной функции используется кусочно-линейная, и задается отступ <i>eps</i> (по умолчанию, равная 0.1): Если разница между прогнозируемым и истинным значением меньше <i>eps</i> (прогнозное значение попадает в пространство гиперплоскости), модель не считает это за ошибку, иначе – берется модуль разницы.</p>	<table border="1"> <tr> <td data-bbox="853 248 1115 312"><b>гиперпараметров</b></td> <td data-bbox="1115 248 1688 312"></td> </tr> <tr> <td data-bbox="853 312 1115 408"><b>Метрика для оптимизации</b></td> <td data-bbox="1115 312 1688 408">Выбирается одна из предлагаемых метрик для оценки работы модели.</td> </tr> <tr> <td data-bbox="853 408 1115 539"><b>Количество фолдов для оптимизации</b></td> <td data-bbox="1115 408 1688 539">Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</td> </tr> </table>	<b>гиперпараметров</b>		<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.	<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.					
<b>гиперпараметров</b>													
<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.												
<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.												
<p><b>Байесовская гребневая регрессия</b></p>	<p>В основе метода лежит формула Байеса, которая дает возможность оценить вероятность событий эмпирическим путем.</p> <p><i>Гребневая</i> регрессия – один из методов снижения размерности. Для гребневой регрессии к функции потерь прибавляется параметр <b>lambda</b>, обозначающий размер штрафа. Чем меньше <b>lambda</b>, тем выше <i>дисперсия</i> и ниже <i>смещение</i>. Смещение – это погрешность оценки, возникающая в результате ошибочного</p>	<table border="1"> <tr> <td data-bbox="853 794 1115 927"><b>alpha_1, alpha_2</b></td> <td data-bbox="1115 794 1688 927">Допустимые максимальные расстояния графика регрессии до верхнего и нижнего доверительного интервала.</td> </tr> <tr> <td data-bbox="853 927 1115 1094"><b>lambda_1, lambda_2</b></td> <td data-bbox="1115 927 1688 1094">Размеры штрафов при выходе прогнозируемых значений за пределы верхнего и нижнего доверительного интервала.</td> </tr> <tr> <td data-bbox="853 1094 1115 1227"><b>Оптимизация гиперпараметров</b></td> <td data-bbox="1115 1094 1688 1227">Флаг подбора гиперпараметров.</td> </tr> <tr> <td data-bbox="853 1227 1115 1326"><b>Метрика для оптимизации</b></td> <td data-bbox="1115 1227 1688 1326">Выбирается одна из предлагаемых метрик для оценки работы модели.</td> </tr> <tr> <td data-bbox="853 1326 1115 1422"><b>Количество фолдов для</b></td> <td data-bbox="1115 1326 1688 1422">–</td> </tr> </table>	<b>alpha_1, alpha_2</b>	Допустимые максимальные расстояния графика регрессии до верхнего и нижнего доверительного интервала.	<b>lambda_1, lambda_2</b>	Размеры штрафов при выходе прогнозируемых значений за пределы верхнего и нижнего доверительного интервала.	<b>Оптимизация гиперпараметров</b>	Флаг подбора гиперпараметров.	<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.	<b>Количество фолдов для</b>	–	<p>**</p>
<b>alpha_1, alpha_2</b>	Допустимые максимальные расстояния графика регрессии до верхнего и нижнего доверительного интервала.												
<b>lambda_1, lambda_2</b>	Размеры штрафов при выходе прогнозируемых значений за пределы верхнего и нижнего доверительного интервала.												
<b>Оптимизация гиперпараметров</b>	Флаг подбора гиперпараметров.												
<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.												
<b>Количество фолдов для</b>	–												

	<p>предположения в алгоритме обучения. В результате большого смещения алгоритм может пропустить связь между признаками и выводом (недообучение). Дисперсия – это ошибка чувствительности к малым отклонениям в тренировочном наборе. При высокой дисперсии алгоритм может трактовать случайный шум в тренировочном наборе, а не желаемый результат (переобучение).</p>	<table border="1"> <tr> <td data-bbox="846 252 1115 316"><b>оптимизации</b></td> <td data-bbox="1115 252 1688 316"></td> </tr> </table>	<b>оптимизации</b>						
<b>оптимизации</b>									
<p><b>Метод ближайших соседей для регрессии</b></p>	<p><b>к- для</b> Для регрессии объекту присваивается среднее значение по <math>k</math> ближайшим к нему объектам, значения которых уже известны. Алгоритм применяется к выборке с большим количеством атрибутов (многомерной). Для этого перед применением определяется функция расстояния, классический вариант такой функции – <i>евклидова метрика</i>. Разные признаки могут иметь разный диапазон</p>	<table border="1"> <tr> <td data-bbox="846 914 1115 1042"><b>Количество ближайших соседей</b></td> <td data-bbox="1115 914 1688 1042">Число <math>k</math>, характеризующее количество соседей в кластере.</td> </tr> <tr> <td data-bbox="846 1042 1115 1249"><b>Тип веса для соседей</b></td> <td data-bbox="1115 1042 1688 1249">Задается одно из значений: ‘uniform’ (единый – всем признакам присваивается единый вес), или ‘distance’ (по расстоянию). Значение по умолчанию – единый.</td> </tr> <tr> <td data-bbox="846 1249 1115 1417"><b>Метрика расстояния</b></td> <td data-bbox="1115 1249 1688 1417">Задается одно из значений: ‘chebyshev’ (Чебышева), ‘euclidean’ (Евклидова), ‘cosine’ (Косинусное), ‘cityblock’ (Манхэттенское).</td> </tr> </table>	<b>Количество ближайших соседей</b>	Число $k$ , характеризующее количество соседей в кластере.	<b>Тип веса для соседей</b>	Задается одно из значений: ‘uniform’ (единый – всем признакам присваивается единый вес), или ‘distance’ (по расстоянию). Значение по умолчанию – единый.	<b>Метрика расстояния</b>	Задается одно из значений: ‘chebyshev’ (Чебышева), ‘euclidean’ (Евклидова), ‘cosine’ (Косинусное), ‘cityblock’ (Манхэттенское).	<p align="center">**</p>
<b>Количество ближайших соседей</b>	Число $k$ , характеризующее количество соседей в кластере.								
<b>Тип веса для соседей</b>	Задается одно из значений: ‘uniform’ (единый – всем признакам присваивается единый вес), или ‘distance’ (по расстоянию). Значение по умолчанию – единый.								
<b>Метрика расстояния</b>	Задается одно из значений: ‘chebyshev’ (Чебышева), ‘euclidean’ (Евклидова), ‘cosine’ (Косинусное), ‘cityblock’ (Манхэттенское).								

	<p>представленных значений в выборке, поэтому выполняется <i>нормализация</i> данных. Некоторые значимые признаки могут быть важнее остальных, поэтому для каждого признака задается определенный <i>вес</i>. Алгоритм предполагает, что похожие наблюдения существуют в непосредственной близости: улавливается идея сходства (иногда называемого расстоянием или близостью) благодаря вычислению Евклидова расстояния между точками.</p>	<table border="1"> <tr> <td data-bbox="833 244 1115 416"></td> <td data-bbox="1115 244 1688 416">Значение по умолчанию – евклидово расстояние, когда вычисляется расстояние между всеми точками попарно.</td> </tr> <tr> <td data-bbox="833 416 1115 517"><b>Оптимизация гиперпараметров</b></td> <td data-bbox="1115 416 1688 517">Флаг подбора гиперпараметров.</td> </tr> <tr> <td data-bbox="833 517 1115 617"><b>Метрика для оптимизации</b></td> <td data-bbox="1115 517 1688 617">Выбирается одна из предлагаемых метрик для оценки работы модели.</td> </tr> <tr> <td data-bbox="833 617 1115 743"><b>Количество фолдов для оптимизации</b></td> <td data-bbox="1115 617 1688 743">Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</td> </tr> </table>		Значение по умолчанию – евклидово расстояние, когда вычисляется расстояние между всеми точками попарно.	<b>Оптимизация гиперпараметров</b>	Флаг подбора гиперпараметров.	<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.	<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.	
	Значение по умолчанию – евклидово расстояние, когда вычисляется расстояние между всеми точками попарно.										
<b>Оптимизация гиперпараметров</b>	Флаг подбора гиперпараметров.										
<b>Метрика для оптимизации</b>	Выбирается одна из предлагаемых метрик для оценки работы модели.										
<b>Количество фолдов для оптимизации</b>	Указывается, на сколько равных частей разбивается входной датасет при обучении модели.										
<p><b>Авторегрессия</b> – модель временного ряда, в которой ее текущее значение <i>линейно</i> зависит от предыдущих (ретроспективных) значений этого же ряда. <i>Линейная зависимость</i> означает, что текущее значение равно взвешенной сумме нескольких предыдущих значений ряда. Зная параметры модели и соответствующие <i>ретроспективные</i> значения временного ряда, можно предсказать его будущие значения. Основное назначение авторегрессионной модели – прогнозирование. Также с ее помощью можно производить анализ временных рядов – выявлять тенденции, сезонность, и другие особенности.</p>											
<p><b>ARIMA/ SARIMAX</b></p>	<p><b>Авторегрессионное интегрированное скользящее среднее</b> (с англ. ARIMA – autoregressive integrated moving average) используется при работе с временными рядами для более глубокого понимания данных, или предсказания</p>	<table border="1"> <tr> <td data-bbox="833 1126 1115 1262"><b>Число шагов для прогноза</b></td> <td data-bbox="1115 1126 1688 1262">Количество шагов, на которые модель будет предсказывать.</td> </tr> <tr> <td data-bbox="833 1262 1115 1430"><b>Порядок авторегрессии, p</b></td> <td data-bbox="1115 1262 1688 1430">Количество запаздывающих наблюдений, включенных в модель, также называется лаговый порядок. <i>P</i> помогает настроить линию для</td> </tr> </table>	<b>Число шагов для прогноза</b>	Количество шагов, на которые модель будет предсказывать.	<b>Порядок авторегрессии, p</b>	Количество запаздывающих наблюдений, включенных в модель, также называется лаговый порядок. <i>P</i> помогает настроить линию для	<p>**</p>				
<b>Число шагов для прогноза</b>	Количество шагов, на которые модель будет предсказывать.										
<b>Порядок авторегрессии, p</b>	Количество запаздывающих наблюдений, включенных в модель, также называется лаговый порядок. <i>P</i> помогает настроить линию для										

<p>будущих точек ряда. Упоминается как <b>ARIMA (p, d, q)</b>, где <i>p</i>, <i>d</i> и <i>q</i> – целые неотрицательные числа, характеризующие порядок для частей модели (соответственно – авторегрессионной, интегрированной и скользящего среднего).</p> <p><b>Авторегрессия.</b> Модель, использующая зависимую связь между наблюдением и некоторым количеством запаздывающих наблюдений.</p> <p><b>Интегрированный.</b> Использование разности необработанных наблюдений (например, вычитание наблюдения из наблюдения на предыдущем временном шаге), чтобы сделать временной ряд стационарным.</p> <p><b>Скользящая средняя.</b> Модель, в которой используется зависимость между наблюдением и остаточной ошибкой из модели скользящего среднего, применяемая к запаздывающим наблюдениям.</p>		прогнозирования серии. Чисто авторегрессионные модели напоминают линейную регрессию, где прогностическими переменными являются <i>p</i> числа предыдущих периодов.	
	<b>Порядок интегрирования</b> , <i>d</i>	Число обычных дифференцирований – количество раз, когда необработанные наблюдения различаются, также называется степенью различия. В модели ARIMA временные ряды преобразуются в <i>стационарные</i> (серии без тренда и сезонности), используя дифференцирование. Стационарный ряд – это когда среднее значение и дисперсия постоянны во времени.	
	<b>Порядок скользящего среднего</b> , <i>q</i>	Размер окна скользящей средней.	
	<b>Параметры модели SARIMAX:</b>		
	<b>Порядок авторегрессии</b>	–	
	<b>Порядок интегрирования</b>	Число сезонных дифференцирований.	
	<b>Порядок скользящего среднего</b>	–	
	<b>Сезонный</b>	Число наблюдений за сезон.	

	<p>Модель <b>SARIMAX</b> используется для временных рядов с учетом сезонности.</p>	<p><b>период</b></p>			
		<p><b>Оптимизация гиперпараметров</b></p>	<p>Флаг подбора гиперпараметров.</p>		
<p><b>Группа «Работа с текстами»</b></p>					
<p><b>Автореферирование текста</b></p>	<p>Данная функция представляет собой автоматический процесс выделения краткого содержания текста с помощью модели машинного обучения. На выходе получается датасет заданного объема, который можно представить в виде таблицы.</p>	<table border="1"> <tr> <td data-bbox="857 515 1016 647"> <p><b>Объем автореферата</b></p> </td> <td data-bbox="1016 515 1682 647"> <p>Максимальное количество символов в выходном результате</p> </td> </tr> </table>	<p><b>Объем автореферата</b></p>	<p>Максимальное количество символов в выходном результате</p>	<p>Таблица с кратким содержанием</p>
<p><b>Объем автореферата</b></p>	<p>Максимальное количество символов в выходном результате</p>				
<p><b>Группа «Управление моделями»</b></p>					
<p><b>Сохранение модели</b></p>	<p>Сохраняет модель по настроенному в системе пути и названию файла, а также сохраняет словарь с переменными. В этом словаре содержится отдельно список независимых переменных и список целевых признаков, с указанием выполненных над ними преобразований. Сохраняет шаг ресемплирования</p>	<table border="1"> <tr> <td data-bbox="857 959 1016 1043"> <p><b>Название модели</b></p> </td> <td data-bbox="1016 959 1682 1043"> <p>Пользователь задает название для обучаемой модели.</p> </td> </tr> </table>	<p><b>Название модели</b></p>	<p>Пользователь задает название для обучаемой модели.</p>	<p>Созданная модель сохраняется в пункте меню системы <b>Модели</b> – <b>&gt; Сохранённые модели</b></p>
<p><b>Название модели</b></p>	<p>Пользователь задает название для обучаемой модели.</p>				
<p><b>Классификация</b></p>					



<p><b>Сохранение модели классификации изображений</b></p>	<p>Функция предназначена для сохранения в системе модели классификации изображений.</p>	<p><b>Название модели</b></p>	<p>Пользователь задает название для обучаемой модели.</p>	<p>**</p>
<p><b>Обнаружение объектов</b></p>				
<p><b>Сохранение модели YOLO v5</b></p>	<p><i>* данный функционал находится в разработке, в текущей версии 2.3.3 применение функции недоступно.</i>                  Функция предназначена для сохранения в системе модели распознавания изображений «YOLO v5». Сохранив модель, можно создать на ее основе приложение для последующей интеграции со сторонними системами. Также обученная модель может использоваться повторно для анализа онлайн данных.</p>	<p><b>Название модели</b></p>	<p>Пользователь задает название для обучаемой модели.</p>	<p>**</p>
<p><b>Spark</b></p>				
<p><b>Сохранение модели Spark</b></p>	<p>Функция предназначена для сохранения в системе моделей, собранных с применением технологии Spark.</p>	<p><b>Название модели</b></p>	<p>Пользователь задает название для обучаемой модели.</p>	
<p><b>Группа «Глубокое обучение»</b></p>				

<p><b>Валидация модели классификации изображений</b></p>	<p>После того, как <i>модель нейронной сети</i> обучена, натренирована и для нее выбраны оптимальные гиперпараметры, необходимо проверить ее точность и адекватность. Для этого выполняется валидация <i>итоговой модели нейронной сети</i> на тестовой выборке.</p> <p>В качестве <i>входных данных</i> для функции используются:</p> <ul style="list-style-type: none"> <li>-тестовый датасет с изображениями;</li> <li>- обученная модель;</li> <li>-словарь с преобразованиями данных;</li> <li>-выбранная метрика валидации.</li> </ul> <p>Рассчитывается сколько изображений тестовой выборки попадают в каждую ячейку матрицы ошибок.</p> <p>Оценивается качество классификации.</p>	<p><b>Метрика</b></p> <p>Метрика, которая оценивает работу обученной модели нейронной сети. Применяется к типу данных – изображения.</p> <p>Чтобы оценить качество модели классификации используются следующие метрики:</p> <ol style="list-style-type: none"> <li><b>Accuracy</b> – оценивает долю правильных ответов модели.</li> <li><b>F1</b> (среднее гармоническое) – агрегированная функция, которая позволяет вместо точности и полноты использовать только один параметр качества классификации. Формула:             <math display="block">F_1 = \frac{2 \cdot precision \cdot recall}{precision + recall}</math> <p>Чем ближе F1 к 1, тем лучше.</p> <p>*В задачах, в которых точность и полнота не равноценны, применяется взвешенное значение <math>F_{\beta}</math>.</p> </li> <li><b>Precision</b>. Метрика, которая оценивает <i>точность модели</i>. Рассчитывается по формуле:             <math display="block">precision = \frac{TP}{TP + FP}</math> <p>здесь считается точность для класса 1 (для класса 0 считается аналогично).</p> </li> <li><b>Recall</b> – оценивает <i>полноту модели</i>:             <math display="block">recall = \frac{TP}{TP + FN}</math> </li> </ol> <p>Идеально, чтобы точность и полнота были равны 1 (100%).</p>	<p>1.Строится <i>матрица ошибок</i>. Пример для бинарной классификации:</p> <table border="1" data-bbox="1736 351 2128 574"> <thead> <tr> <th colspan="2" rowspan="2"></th> <th colspan="2">Истинный класс</th> </tr> <tr> <th>0</th> <th>1</th> </tr> </thead> <tbody> <tr> <th rowspan="2">0</th> <th>0</th> <td>TN</td> <td>FN</td> </tr> <tr> <th>1</th> <td>FP</td> <td>TP</td> </tr> </tbody> </table> <p>где столбцы – истинные классы, а строки – предсказанные классы.</p> <p>Обозначения:</p> <ul style="list-style-type: none"> <li>TN – true negative</li> <li>TP – true positive</li> <li>FN – false negative</li> <li>FP – false positive</li> </ul> <p>Например, ячейка <b>TP</b> означает, что объект действительно принадлежит классу 1, и для него предсказан класс 1. А <b>FN</b> – объект неправильно отнесли к классу 0, хотя он принадлежит к классу 1.</p> <p>2.Отображается значение выбранной метрики.</p>			Истинный класс		0	1	0	0	TN	FN	1	FP	TP
		Истинный класс														
		0	1													
0	0	TN	FN													
	1	FP	TP													

		<p>5. <b>AUC_ROC</b>. Для анализа качества модели применяется ROC-анализ: строится ROC-кривая, которая наиболее часто используется для представления результатов бинарной классификации. Классов два: один называется классом с положительными исходами, второй – с отрицательными исходами. ROC-кривая показывает зависимость количества верно классифицированных положительных примеров от количества неверно классифицированных отрицательных примеров. Чем выше показатель AUC_ROC, тем качественнее классификатор, при этом значение 0,5 демонстрирует непригодность выбранного метода классификации. AUC_ROC = 0,9-1,0 означает отличное качество модели.</p>							
<b>Классификация</b>									
<p><b>Классификация (табличные данные)</b></p>	<p>Типы решаемых задач: анализ тональности, классификация текста по категориям, распознавание речи и многое другое. Рассмотрим функцию на примере задачи по отнесению документов к определенной категории на основании его содержания. Процесс классификации осуществляется с помощью применения методов машинного обучения, в частности <i>сверточных нейронных сетей</i>. Задача</p>	<table border="1"> <tr> <td data-bbox="853 927 1160 1058"><b>Количество эпох</b></td> <td data-bbox="1160 927 1688 1058">Параметр, который показывает сколько раз <i>модель</i> подвергается воздействию обучения.</td> </tr> <tr> <td data-bbox="853 1058 1160 1225"><b>Размер мини-батча</b></td> <td data-bbox="1160 1058 1688 1225">Количество обучающих примеров за одну итерацию. Под примерами имеются в виду <i>наблюдения</i> – строки в табличных данных.</td> </tr> <tr> <td data-bbox="853 1225 1160 1353"><b>Метрика для обучения</b></td> <td data-bbox="1160 1225 1688 1353">Метрика «Accuracy» (точность) показывает долю правильных ответов алгоритма</td> </tr> </table>	<b>Количество эпох</b>	Параметр, который показывает сколько раз <i>модель</i> подвергается воздействию обучения.	<b>Размер мини-батча</b>	Количество обучающих примеров за одну итерацию. Под примерами имеются в виду <i>наблюдения</i> – строки в табличных данных.	<b>Метрика для обучения</b>	Метрика «Accuracy» (точность) показывает долю правильных ответов алгоритма	**
<b>Количество эпох</b>	Параметр, который показывает сколько раз <i>модель</i> подвергается воздействию обучения.								
<b>Размер мини-батча</b>	Количество обучающих примеров за одну итерацию. Под примерами имеются в виду <i>наблюдения</i> – строки в табличных данных.								
<b>Метрика для обучения</b>	Метрика «Accuracy» (точность) показывает долю правильных ответов алгоритма								

<p>классификации текстов применима в решении следующих задач: борьба с массовой рассылкой рекламы, распознавание тональности текстов, сортировка документов и т.д.</p> <p>Задача определяется следующим образом: пусть существует конечное множество категорий, на вход алгоритма подается конечное количество документов, и есть целевая функция, которая определяет соответствие для каждой пары (документ, категория). Задача состоит в нахождении этой функции, называемой классификатором.</p> <p>Строится многослойная нейронная сеть, состоящая из слоев:</p> <ul style="list-style-type: none"> <li>-<i>входной</i>, на который поступают входные признаки;</li> <li>-<i>скрытый</i>, на котором рассчитываются промежуточные результаты;</li> <li>- <i>выходной</i>, на котором выводятся окончательные значения, вычисленные по гипотезе.</li> </ul>	<p><b>Алгоритм градиентного спуска</b></p> <p>Метод нахождения минимального значения <i>функции потерь</i>. Алгоритмы на выбор: <i>SGD, Adam</i>, и др.</p>						
	<p><b>Шаг градиентного спуска</b></p> <p>Параметр, регулирующий скорость обучения модели. Значения – 0.001 или 0.1.</p>						
	<p><b>Функция потерь</b></p> <p>Выбирается одна из функций, в зависимости от задачи: бинарная классификация или многоклассовая.</p>						
	<p><b>Добавить слой</b></p> <p>Задается одно из трех значений: Conv2D, Flatten, Dense (последний – полносвязный слой).</p>						
	<p><b>Перемешивать выборку перед обучением</b></p> <p>Необходимо установить галочку в поле, чтобы случайным образом поменять местами наблюдения, соответствующие строкам в таблице</p>						
	<p><b>Порог классификации</b></p> <p>Для бинарной классификации значение по умолчанию 0.5, для многоклассовой – параметр не заполняется</p>						
	<p><b>Флаг возврата вероятности при прогнозе</b></p> <p>Возвращает вероятность или метки классов для дальнейшего прогноза после обучения модели. Для бинарной классификации может использоваться одно поле с метками 0 и 1, обозначающими принадлежность к тому или иному классу. Для многоклассовой классификации используется несколько полей –</p>						

	<p><i>*Сверточными</i> искусственные нейронные сети называются из-за специальной архитектуры, с наличием операций сверки.</p>		<p>каждое поле соответствует отдельному классу (0, 1, 2 и т.д.), и записываются вероятности (от 0 до 1), с которыми наблюдения принадлежат классам</p>	
		<p><b>Оптимизация гиперпараметров</b></p>	<p>Алгоритм подбирает гиперпараметры</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Выбирается одна из предлагаемых метрик для оценки работы модели.</p>	
		<p><b>Количество фолдов для оптимизации</b></p>	<p>—</p>	
<p><b>Классификация изображений</b></p>	<p>Алгоритм решения задачи классификации:          1. Берется <i>тренировочная выборка</i> – набор изображений с известными значениями целевого признака <math>Y</math>. Нейронная сеть должна восстановить зависимость между нецелевыми признаками и целевым.          2. Задаются основные параметры нейронной сети.          3. Выписываются выражения для вероятностей принадлежности наблюдения к тому или иному классу (<math>Y = 0, 1, 2, 3, \dots</math> и т.д.).</p>	<p><b>Количество эпох</b></p>	<p>Это гиперпараметр, который определяет сколько раз <i>алгоритм обучения</i> будет обрабатывать весь <i>набор обучающих данных</i>. То есть <i>эпоха</i> – одна итерация в процессе обучения, включающая предъявление всех примеров из обучающего множества</p>	<p>**</p>
		<p><b>Метрика для обучения</b></p>	<p>Метрика «Accuracy» (точность) показывает долю правильных ответов алгоритма</p>	
		<p><b>Алгоритм градиентного спуска</b></p>	<p>Метод нахождения минимального значения <i>функции потерь</i>. Минимизация любой функции означает поиск самой глубокой впадины в этой функции. Функция используется, чтобы контролировать ошибку в прогнозах</p>	

	<p>4. По тренировочной выборке составляется функция потерь.</p> <p>5. Функция потерь <math>L(w)</math> содержит вхождения весов нейронной сети. Относительно этих переменных находится <i>точка минимума функции <math>L(w)</math></i>.</p> <p>6. Точка минимума определяет оптимальные веса нейронной сети.</p> <p>7. Весам нейронной сети присваиваются найденные оптимальные значения. Пусть изображение <math>A</math> не принадлежит тренировочной выборке. Объект <math>A</math> прогоняется через нейронную сеть и на выходе получают вероятности – они и являются предсказаниями для объекта <math>A</math> (по максимальной вероятности определяется принадлежность объекта к классу).</p>		<p>модели. Поиск минимума означает получение наименьшей возможности ошибки или повышение точности модели. Точность увеличивается перебором <i>учебных</i> данных при настройке параметров модели (весов и смещений). Суть алгоритма – процесс получения наименьшего значения ошибки. Аналогично это можно рассматривать как спуск во впадину в попытке найти самое низкое значение ошибки.</p> <p>Можно выбрать один из следующих алгоритмов: <i>SGD</i> (Stochastic gradient descent, с англ. Стохастический градиентный спуск), <i>Adam</i>, и др. Алгоритм <b>Adam</b> является <i>модифицированным</i>, в нем также выполняется минимизация функции потерь. Рассчитываются векторы <i>частных</i> в текущей точке функции, и определяются координаты следующей точки. Частные производные вычисляются, чтобы определить, какой был вклад в ошибку по каждому весу.</p>	
		<p><b>Шаг градиентного спуска</b></p>	<p>Параметр, регулирующий скорость обучения модели – насколько быстро функция потерь спускается к своему минимуму (скорость спуска/поиска). Выбирается значение: 0.001, или 0.1.</p>	
		<p><b>Функция потерь</b></p>	<p>Функция потерь находится в центре нейронной сети. Она используется для расчета ошибки между <i>реальными</i> и</p>	

			<p><i>полученными</i> ответами. Главная цель – минимизировать эту ошибку. Или: <b>максимизировать вероятность принадлежности к истинному классу для каждого объекта из тренировочной выборки.</b> Она также может зависеть от таких переменных, как веса и смещения, где <i>смещения</i> – это веса, добавленные к скрытым слоям.</p> <p>Выбирается одна из функций, в зависимости от задачи: бинарная классификация или многоклассовая.</p>	
		<p><b>Добавить слой</b></p>	<p>Основой алгоритмов распознавания изображений являются сверточные нейронные сети. Для их построения используются три главных типа слоев: сверточный слой, слой подвыборки и полносвязный слой. Соответственно пользователь задает одно из трех значений: Conv2D, Flatten, Dense (последний – полносвязный).</p> <p>В сверточных нейронных сетях одно изображение является одним наблюдением. Таким образом, исходное изображение преобразуется, слой за слоем, от начального значения пикселя до итоговой оценки класса. Слои, идущие до полносвязного, являются средствами предобработки изображения, и используются для выделения различных признаков, которые затем подаются на вход классификатору.</p>	

		<b>Порог классификации</b>	<b>Важно!</b> Параметр заполняется только для задачи бинарной классификации, значение по умолчанию 0.5. Для множественной классификации это поле остается пустым.	
<b>Регрессия</b>				
<b>Регрессия (табличные данные)</b>	Для обучения нейронной сети данные делятся на части меньшего размера, загружают их по очереди и обновляют веса нейросети в конце каждого шага, подстраивая их под данные.	<b>Количество эпох</b>	Указывается количество эпох для обучения модели. Одна эпоха – весь датасет прошел через нейронную сеть в прямом и обратном направлении только один раз. Так как одна эпоха слишком велика для компьютера, датасет делят на партии – <i>батчи</i> . С увеличением числа эпох, веса нейронной сети изменяются все большее количество раз. Кривая с каждым разом лучше подстраивается под данные, переходя последовательно из плохо обученного состояния в оптимальное. Если вовремя не остановиться, то может произойти переобучение.	**
		<b>Размер мини-батча</b>	Общее число тренировочных объектов, представленных в одном батче. Нельзя пропустить через нейронную сеть разом весь датасет. Поэтому делим данные на пакеты, сету или партии.	



			<i>Итерации</i> – это число батчей, необходимое для завершения одной эпохи.	
		<b>Метрика для обучения</b>	Для регрессии: ['MAE', 'MAPE', 'MSE'].	
		<b>Алгоритм градиентного спуска</b>	Алгоритм итеративной оптимизации, используемой в машинном обучении для получения более точного результата. <i>Градиент</i> показывает скорость убывания или возрастания функции. <i>Спуск</i> говорит о том, что мы имеем дело с убыванием. Алгоритм <b>итеративный</b> , процедура проводится несколько раз, чтобы добиться оптимального результата. На каждом шаге результат получается лучше.	
		<b>Шаг градиентного спуска</b>	Скорость обучения модели алгоритмом градиентного спуска.	
		<b>Функция потерь</b>	Функция, которая используется для оптимизации алгоритма машинного обучения. Значение, вычисленное такой функцией, называется 'потерей'. Потери регрессии рассчитываются путем прямого сравнения выходного и истинного значения. Самая популярная функция для регрессионных моделей – это среднеквадратическая ошибка, MSE. Функция потерь определяет, как именно выходные данные связаны с исходными.	

			<p>По сути вычисляется насколько хорошо работает модель – сравнивается то, что модель прогнозирует, с фактическим значением. Сохраняется функция потерь, которая может эффективно наказывать модель, пока та обучается на тренировочных данных.</p>	
		<p><b>Добавить слой</b></p>	<p>Из списка выбирается дополнительный слой для нейросети.</p>	
		<p><b>Перемешивать выборку перед обучением</b></p>	<p>Необходимо установить галочку в поле, чтобы случайным образом поменять местами наблюдения.</p>	
		<p><b>Оптимизация гиперпараметро в</b></p>	<p>Если установить галочку в поле, то алгоритм выберет наилучшие гиперпараметры для создания модели из списка предложенных</p>	
		<p><b>Метрика для оптимизации</b></p>	<p>Выбирается одна из предлагаемых метрик для оценки работы модели.</p>	
		<p><b>Количество фолдов для оптимизации</b></p>	<p>Указывается, на сколько равных частей разбивается входной датасет при обучении модели.</p>	
<p><b>Обнаружение объектов</b></p>				
<p><b>YOLOv5</b></p>	<p>* <i>данный функционал находится в разработке, в текущей версии 2.3.3 применение функции недоступно.</i> Датасет должен быть разделен на две папки: <i>train</i></p>	<p><b>Размер мини-батча</b></p>	<p>Указывается количество изображений, которое одновременно подается на вход YOLO. Например, если задать размер 2,</p>	<p>В результате получаются изображения с обозначением детектированных объектов и значениями <i>confidence</i>. Где <i>confidence</i> – число от 0 до 1, характеризующее ‘уверенность’</p>

	<p>(тренировочная выборка) и <i>val</i> (валидационная). В каждой папке лежат еще две папки: <i>images</i> (изображения) и <i>labels</i> – папка с текстовыми файлами, содержащими метки объектов на этих изображениях в формате YOLO.</p> <p>Для этой функции предварительно выполняется разметка изображений на тренировочной и валидационной выборках. Пользователь с помощью ‘<i>bounding box</i>’ отмечает объекты на изображениях. Алгоритм по точкам объектов находит функцию их обнаружения. На валидационной выборке проверяется точность обученной модели.</p> <p>После запуска функции на выходе получается обученная нейронная сеть, результаты обучения которой сохраняются в БД.</p>	<p><b>Количество эпох</b></p>	<p>за один подход подается два изображения.</p> <p>Это гиперпараметр, который определяет сколько раз <i>алгоритм обучения</i> будет обрабатывать весь <i>набор обучающих данных</i>. То есть <i>эпоха</i> – одна итерация в процессе обучения, включающая предъявление всех примеров из обучающего множества. Например, для выборки в десять изображений и размера мини-батча два, эпоха равна пяти прохождениям.</p>	<p>модели в том, что детектирован объект или детектирован объект определенного класса. Еще один параметр <i>conf-thres</i> позволяет установить пороговое значение для <i>confidence</i> модели. Все объекты, <i>confidence</i> которых ниже этого значения не считаются объектами.</p> <p>Также отображаются: описание модели, графики обучения модели, матрица ошибок на валидационных данных. Где описание модели содержит информацию об оптимизаторе, тренировочном и валидационном датасетах, а также параметры обучения.</p>
<p><b>Отправка уведомлений</b></p>				
<p><b>Отправка уведомлений</b></p>	<p>Функция предназначена для осуществления отправки</p>		<p>Оповещение в телеграм канал</p>	

	<p>уведомлений в настроенный канал телеграм или по электронной почте (не реализовано в текущей версии). Данная функция выступает в паре с блоком шлюз, где необходимо задать условия, при которых будет отправлено уведомление.</p>	<table border="1"> <tr> <td data-bbox="835 244 1061 411"><b>Канал уведомлений</b></td> <td data-bbox="1061 244 1688 411">Выбирается ранее настроенный в разделе Администрирование -&gt; Уведомления канал, на который будет осуществляться отправка сообщений.</td> </tr> </table>	<b>Канал уведомлений</b>	Выбирается ранее настроенный в разделе Администрирование -> Уведомления канал, на который будет осуществляться отправка сообщений.							
<b>Канал уведомлений</b>	Выбирается ранее настроенный в разделе Администрирование -> Уведомления канал, на который будет осуществляться отправка сообщений.										
<p><b>Spark.</b> Группа функций для фреймворка Apache Spark. Названия функций дублируются с теми, что были описаны ранее, разница заключается в использовании другого модуля машинного обучения в Apache Spark <i>(в следующих версиях Платформы планируется сделать все функции универсальными)</i>.</p>											
<p><b>Сохранение датасета Spark в CSV</b></p>	<p>Функция распределяет входные данные в несколько файлов в одну директорию. Для этого выбираем: куда сохранить данные, как их назвать, и по необходимости можем подгрузить новую порцию данных.</p>	<table border="1"> <tr> <td data-bbox="835 655 1149 826"><b>Путь до директории для датасета</b></td> <td data-bbox="1149 655 1688 826">Выбирается путь до папки, в которую будут сохраняться данные.</td> </tr> <tr> <td data-bbox="835 826 1149 997"><b>Название датасета</b></td> <td data-bbox="1149 826 1688 997">В этом поле задается название для датасета. По умолчанию датасеты создаются с названиями формата <i>pySpark.csv</i>.</td> </tr> <tr> <td data-bbox="835 997 1149 1262"><b>Добавить данные к датасету</b></td> <td data-bbox="1149 997 1688 1262">Если преобразованные данные необходимо сохранять не в виде отдельного файла, а добавить к уже существующему и загруженному на платформе, необходимо установить галочку у данного признака. По умолчанию файл перезаписывается.</td> </tr> <tr> <td data-bbox="835 1262 1149 1418"><b>Название датасета для валидации</b></td> <td data-bbox="1149 1262 1688 1418">Указывается название, с которым будет сохранен датасет для валидации при активации параметра «Сохранить датасет для валидации»</td> </tr> </table>	<b>Путь до директории для датасета</b>	Выбирается путь до папки, в которую будут сохраняться данные.	<b>Название датасета</b>	В этом поле задается название для датасета. По умолчанию датасеты создаются с названиями формата <i>pySpark.csv</i> .	<b>Добавить данные к датасету</b>	Если преобразованные данные необходимо сохранять не в виде отдельного файла, а добавить к уже существующему и загруженному на платформе, необходимо установить галочку у данного признака. По умолчанию файл перезаписывается.	<b>Название датасета для валидации</b>	Указывается название, с которым будет сохранен датасет для валидации при активации параметра «Сохранить датасет для валидации»	<p>Таблица в формате csv с датасетом. Сохраняется в раздел данные</p>
<b>Путь до директории для датасета</b>	Выбирается путь до папки, в которую будут сохраняться данные.										
<b>Название датасета</b>	В этом поле задается название для датасета. По умолчанию датасеты создаются с названиями формата <i>pySpark.csv</i> .										
<b>Добавить данные к датасету</b>	Если преобразованные данные необходимо сохранять не в виде отдельного файла, а добавить к уже существующему и загруженному на платформе, необходимо установить галочку у данного признака. По умолчанию файл перезаписывается.										
<b>Название датасета для валидации</b>	Указывается название, с которым будет сохранен датасет для валидации при активации параметра «Сохранить датасет для валидации»										

		<p><b>Сохранить датасет для валидации</b></p> <p>В процессе работы пайплайна, исходный вид набора данных данных теряется, поэтому его нужно передать из блока "Загрузка данных" в конец пайплайна в блок сохранени. Датасет для валидации это и есть нетронутый набор данных в первоначальном виде, к нему только добавляется столбец с результатами.</p>	
		<p><b>Загрузка датасета для валидации в БД</b></p> <p>Позволяет загрузить датасет для вализации напрямую в базу данных ClickHouse</p>	
<b>Косинусное расстояние</b>	<p>На вход функция получает новые данные для анализа (датасет в формате csv), обученную модель, и числовой вектор. Выполняется поиск объектов, наиболее схожих с заданным вектором, и в качестве меры схожести используется косинусное расстояние - расстояние между значениями во входном векторе и значениями выбранных столбцов в наблюдениях.</p>	–	Таблица «Косинусное расстояние»
<b>Выбор признаков и целевых признаков</b>	<p>Аналогично стандартной функции.</p>		

<p><b>Разделение датасета на обучающую и тестовую выборки</b></p>	<p>Аналогично стандартной функции.</p>				
<p><b>Валидация модели</b></p>	<p>Аналогично стандартной функции.</p>				
<p><b>Прогноз модели</b></p>	<p>Аналогично стандартной функции.</p>				
<p><b>Порядковое кодирование признаков</b></p>	<p>Порядковое кодирование - это метод преобразования категориальных данных в цифровой вид. Применяются, когда в датасете существуют НЕ числовые признаки, которые заданы словами и для дальнейшего анализа их нужно преобразовать в числа. Порядковое кодирование позволяет пронумеровать признаки по порядку.</p>	<table border="1"> <tr> <td data-bbox="853 555 1093 691"> <p><b>Выбранные признаки</b></p> </td> <td data-bbox="1093 555 1688 691"> <p>Указываются признаки, над которыми необходимо провести операцию порядкового кодирования.</p> </td> </tr> </table>	<p><b>Выбранные признаки</b></p>	<p>Указываются признаки, над которыми необходимо провести операцию порядкового кодирования.</p>	
<p><b>Выбранные признаки</b></p>	<p>Указываются признаки, над которыми необходимо провести операцию порядкового кодирования.</p>				
<p><b>Нормализация признаков</b></p>	<p>Нормализация - это приведение числовых признаков к единой шкале. Бывает что числовой признак имеет минимальное и максимальное значение в очень широком диапазоне и это плохо для машинного обучения. Например, есть числовой признак, чье минимальное значение равно 0,001, а максимальное</p>	<p>—</p>	<p>—</p>		

	<p>- 100000, нормализация преобразовывает их к диапазону от 0 до 1, то есть 0.001 становится 0, а 100000 становится 1, значения между ними также преобразуются, 50 000 станет примерно равным 0.5. Данная функция позволяет оптимизировать дальнейшие вычисления.</p>				
<p><b>Модель градиентного бустинга Spark для бинарной классификации</b></p>	<p>Градиентный бустинг представляет собой ансамбль деревьев решений. В основе данного алгоритма лежит итеративное обучение деревьев решений с целью минимизировать функцию потерь. Благодаря особенностям деревьев решений градиентный бустинг способен работать с категориальными признаками, справляться с нелинейностями. Бустинг – это метод преобразования слабообученных моделей в хорошообученные. В бустинге каждое новое дерево обучается на модифицированной версии исходного датасета.</p>	<table border="1" data-bbox="853 655 1688 855"> <tr> <td data-bbox="853 655 1093 855"> <p><b>Количество базовых моделей</b></p> </td> <td data-bbox="1093 655 1688 855"> <p>Указывается число последовательных итераций по оптимизации модели, то есть количество попыток уменьшить ошибку.</p> </td> </tr> </table>	<p><b>Количество базовых моделей</b></p>	<p>Указывается число последовательных итераций по оптимизации модели, то есть количество попыток уменьшить ошибку.</p>	<p>—</p>
<p><b>Количество базовых моделей</b></p>	<p>Указывается число последовательных итераций по оптимизации модели, то есть количество попыток уменьшить ошибку.</p>				

## 19. Лист изменений

Таблица 19.1 – Лист изменений в версии 1.0.0

Наименование раздела	Содержание изменения	Обоснование
-	-	-