



ModelOps в Loginom: управление моделями и экспериментами

Николай Паклин



- Николай Паклин – тимлид группы «Обучение»
- 19 лет в low-code аналитике и Loginom Company
- С командой запустил Loginom Skills, проект «Мастерская», библиотеки компонентов для Loginom
- Бизнес-тренер, преподаватель, автор книг и учебников по аналитике данных

О чем будем говорить

1. Что такое ModelOps
2. Когда и для чего важно внедрять ModelOps
3. Как управлять моделями и экспериментами в Loginom

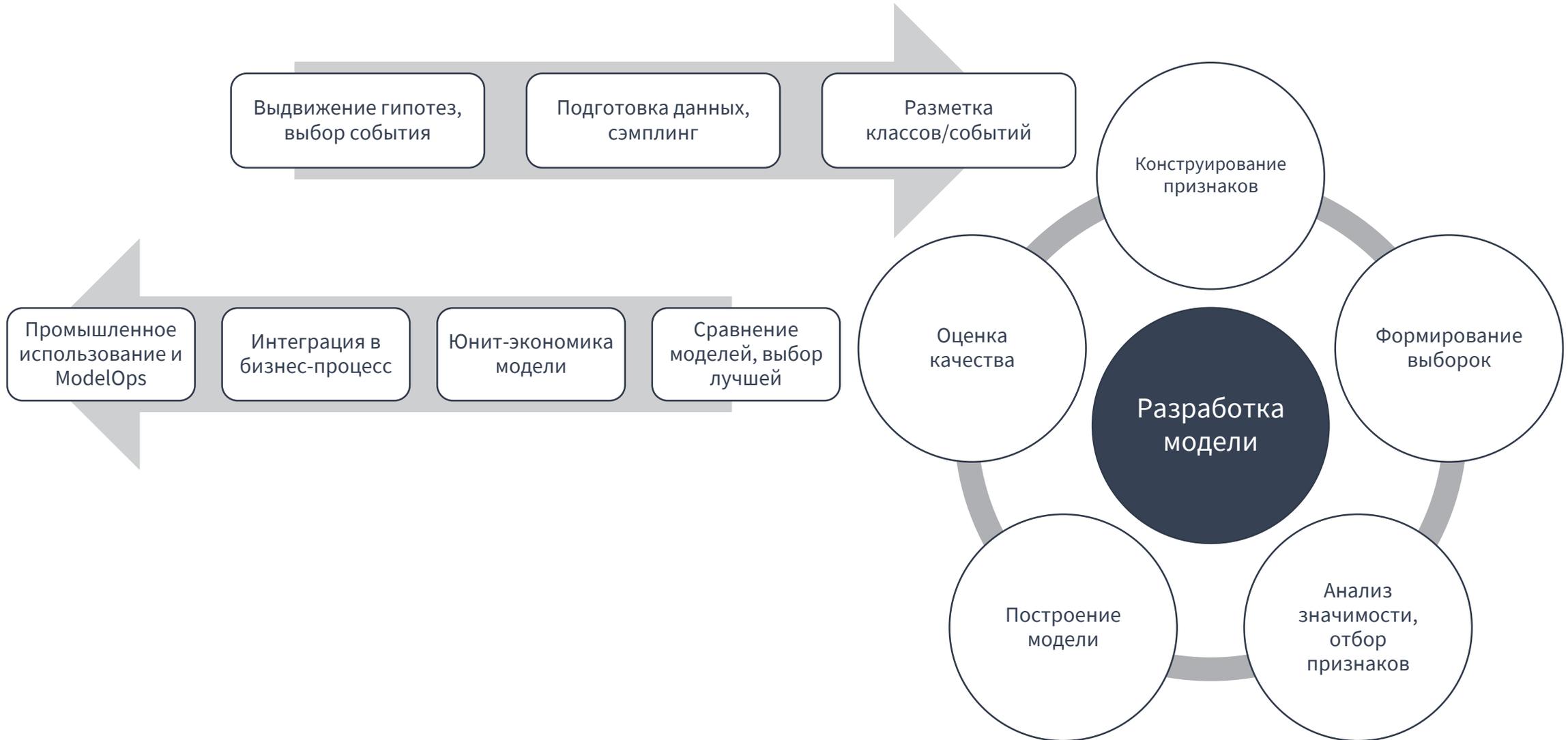
Операционализация аналитики данных

- **ModelOps** – набор практик, направленных на унификацию процессов разработки и развертывания аналитических моделей и решений *(часто идет речь о моделях machine learning - MLOps)*
- **DataOps** – набор практик для повышения качества и скорости предоставления данных, а также сокращения времени и усилий на создание и обслуживание конвейеров данных

Ключевые требования к бизнес-процессу

- В рамках одной команды важно обеспечить **воспроизводимость** вычислений
- Коллеги должны иметь возможность **повторить эксперименты** без особых усилий

Общая схема процесса построения ML-модели



Особенности процесса построения ML-модели

- Несколько версий датасетов
- Несколько моделей
- Трудность выбора лучшей модели сразу
- Проект может требовать переключаться между моделями в процессе эксплуатации

Типичные ситуации

1. Датасеты, на которых строилась модель удалены/утеряны/перезатерты
2. В папке много сценариев и не понятно, какой из них актуальный
3. Построено 20 моделей, но не ясно, где модель с лучшим ROC-AUC
4. Нужно передать коллеге исследовательскую задачу
5. Передача модели (сценария) в другой отдел для переноса в продакшн
6. Нужно понять, чем отличаются друг от друга две модели

Типичные проблемы

1. Различия в среде исполнения (разработка, тестирование, продакшн)
2. Неуправляемая случайность в данных или алгоритмах (Random Seed)
3. Проблемы с зависимостями со ссылками-пакетами, библиотеками языков программирования, зависимости в зависимостях

**Тезис: в определенный момент нам
понадобится повторить ранее
проведенный вычислительный
эксперимент**

Архитектурная схема использования ModelOps



Управление экспериментами и моделями: уровни зрелости



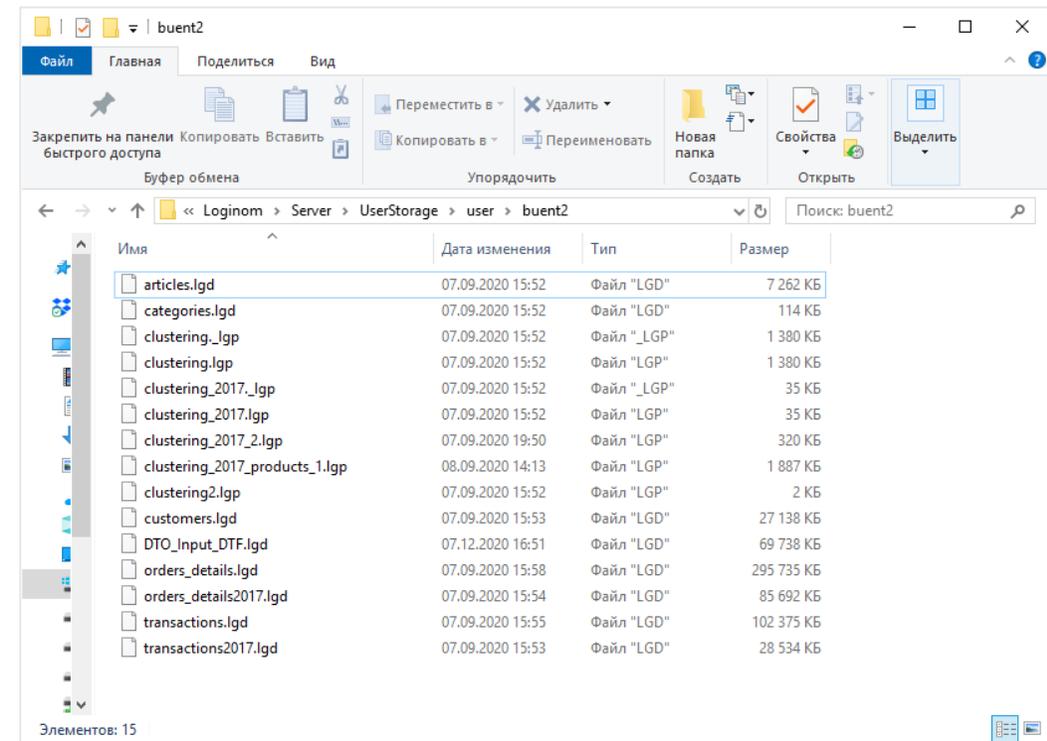
Я положил тебе в сетевую папку проект, там есть все, что требуется



Как же с этим разобраться...



Типичная ситуация при ручных процессах ModelOps

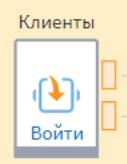


* Пример демонстрирует поиск аномалий алгоритмом LOF в наборе данных по продажам в случае, когда обучение велось на "чистом" наборе данных, то есть без аномалий. Эта задача еще называется "обнаружение новизны" (Novelty Detection).

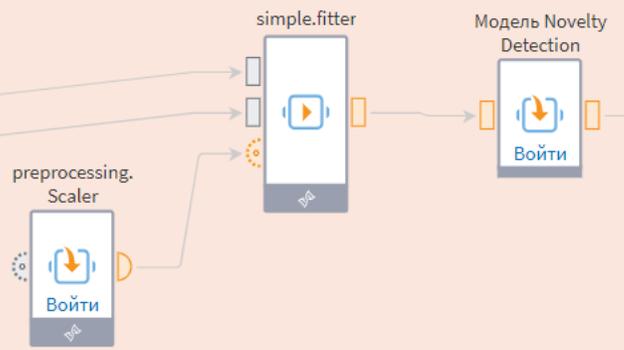
Для формирования такого "чистого" датасета мы удалили аномалии, найденные алгоритмом LOF на "грязном" датасете. Для проверки эффективности модели делается "скоринг" ранее найденных аномалий моделью обнаружения новизны.

Эта задача, когда обучение ведется на одном классе, называется полуконтролируемым обучением (Semi-Supervised Learning).

Исходную выборку мы разделили на 2 части: обучающую и тестовую. Найденные ранее аномалии находятся в тестовой выборке. Таким способом мы моделируем задачу Novelty Detection.



Построение модели с использованием библиотеки Loginom Python Kits. Метод LocalOutlierFactor из библиотеки sklearn. Используем режим Novelty = True. Предварительно делаем z-масштабирование выборки. 🤖 Число соседей n = 20, загрязненность (contamination) 1%.



В результате скоринга через модель обнаружения новизны мы получили 230 совпадений из 230. Такую модель можно использовать для скоринга новых записей, и относить или не относить их к аномальным.



Заметки и комментарии - не панацея, хотя помогают описать сценарий

- loginom-silver-kit
- Project
- Repository
- Issues 2
- Merge Requests 0
 - List
 - Labels
 - Milestones
- CI / CD
- Operations
- Wiki
- Snippets
- Settings

learn > loginom-silver-kit > Merge Requests

Open 0 Merged 30 Closed 2 All 32

Edit merge requests New merge request

Search or filter results... Last updated

Добавить компонент "Переменные в словарь JS"	!32 · opened 4 weeks ago by gusev add	MERGED 0 updated 2 weeks ago
Добавить компонент "Формула Хаверсина" + правки IV-отбора	!31 · opened 2 months ago by paklin add	MERGED 1 updated 2 months ago
Добавление сборщика PDF	!30 · opened 4 months ago by gusev update	MERGED 0 updated 4 months ago
"Релиз для 7.1.0"	!28 · opened 7 months ago by gusev update	MERGED 3 updated 7 months ago
Добавить компонент ParseJSON JS и Перцентиль N%	!27 · opened 8 months ago by gusev add	MERGED 3 updated 8 months ago
Релиз 3.0.2	!26 · opened 10 months ago by gusev update	MERGED 0 updated 10 months ago
Релиз 3.0.1	!25 · opened 1 year ago by gusev update	MERGED 0 updated 1 year ago
"Проверка на 7.0.1, исправление в Кластерные силуэты, Сравнение метрик"	!24 · opened 1 year ago by gusev bug update	MERGED 1 updated 1 year ago
Правки в кластерных силуэтах. Сравнение метрик. Компонент IF-объединение JS	!23 · opened 1 year ago by gusev bug	MERGED 0 updated 1 year ago
Добавление нормализации в компонент Кластерные силуэты	!22 · opened 1 year ago by alaeva update	MERGED 0 updated 1 year ago
Разделитель строк. Обновление библиотеки до 6.5.4"	!20 · opened 1 year ago by gusev bug add	MERGED 11 updated 1 year ago
Компонент Генератор списка	!18 · opened 2 years ago by gusev add bug loginom_silver_kit.lgp update	MERGED 9 updated 2 years ago
Компонент "Разметка событий"		MERGED 0

Git как репозиторий и система версионирования

Git для ModelOps – что не так?

1. Ограничения на размер файлов в репозитории
2. Нельзя отслеживать изменения в бинарных файлах
3. Не предназначен для отслеживания метрик модели и гиперпараметров алгоритма обучения
4. Ориентирован на разработчиков ПО, а не аналитиков данных
5. Плохо согласуется с парадигмой low-code

Требования к фреймворку:

- Бесплатный в облаке, возможно, с ограничениями
- Работа в локальной сети
- Поддержка сторонних файловых хранилищ
- Подходит для систем класса low-code
- Низкий порог входа для аналитика
- Богатый Rest API



www.clear.ml

Модули ClearML

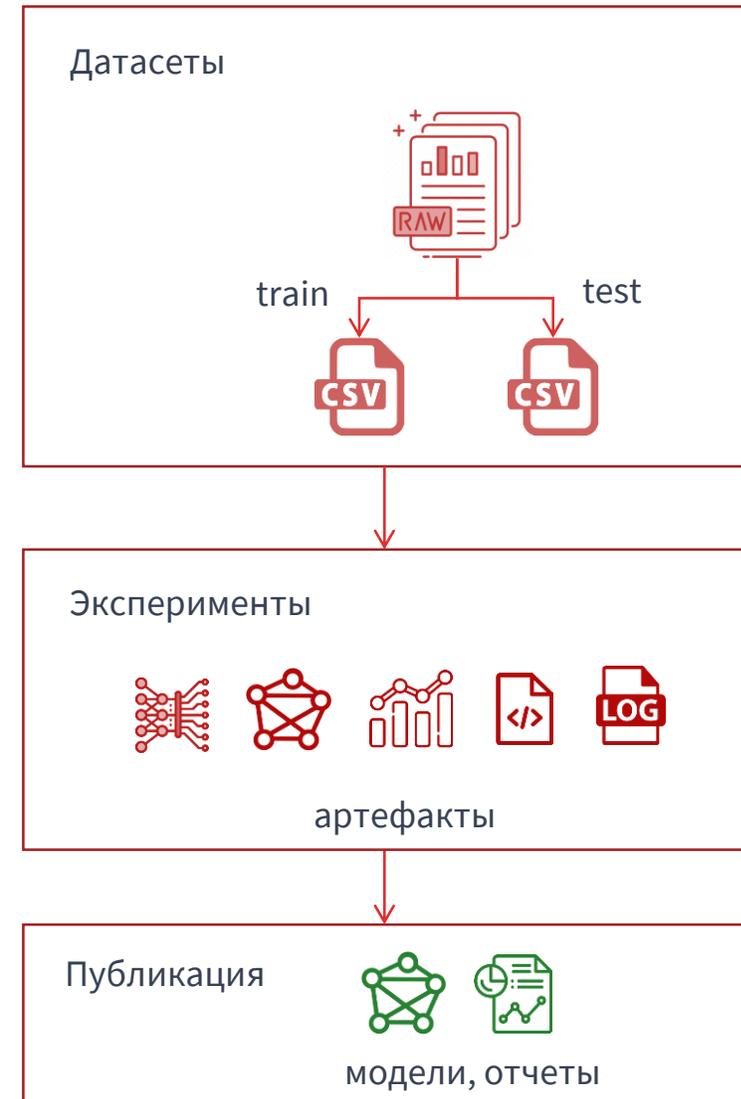


Реализация интеграции с Loginom – управление моделями и экспериментами



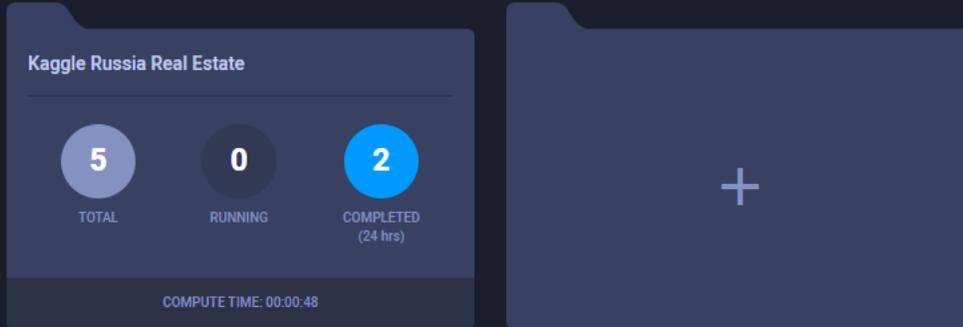
Терминология

1. Проект
2. Датасет, эксперимент
3. Модель, гиперпараметр
4. Метрика, артефакт
5. Отчет
6. Опубликованный объект





RECENT PROJECTS [VIEW ALL](#)



RECENT EXPERIMENTS

[MANAGE WORKERS AND QUEUES](#)

TYPE	TITLE	PROJECT	STARTED	UPDATED	STATUS
Training	XGBoost Regressor (15 признаков)	Kaggle Russia Real Estate	May 3 2024 16:41	May 3 2024 18:12	Published
Training	XGBoost Regressor (15 признаков)	Kaggle Russia Real Estate	May 3 2024 16:28	May 3 2024 16:28	Failed
Training	Линейная регрессия (15 признаков)	Kaggle Russia Real Estate	May 3 2024 15:05	May 3 2024 15:05	Completed
Data Processing	Kaggle-Russia-Real-Estate	Kaggle Russia Real Estate/.datasets/Kaggle-Russia-Real-Estate	May 3 2024 14:22	May 3 2024 14:23	Completed
Training	Линейная регрессия (14 признаков)	Kaggle Russia Real Estate	Apr 27 2024 17:10	Apr 27 2024 17:11	Completed

Веб-интерфейс ClearML: проекты и последние эксперименты

[Get Started](#)

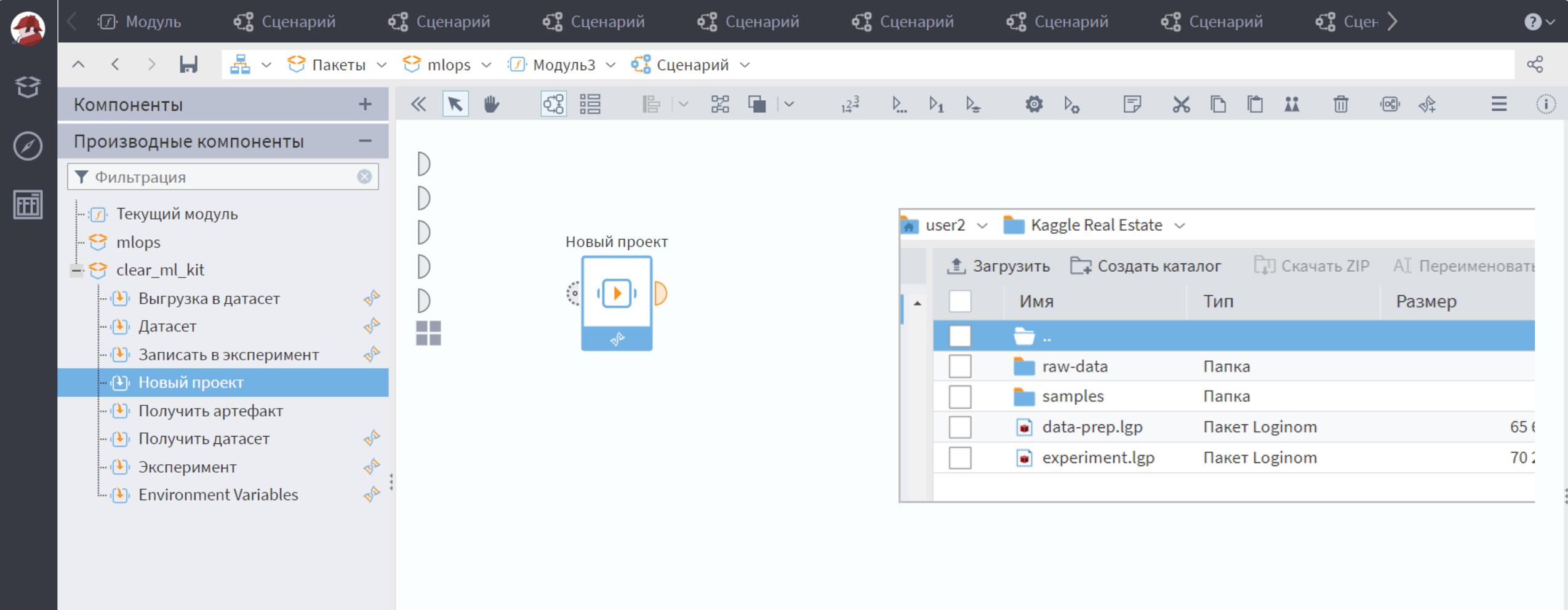
Технические подробности

Состав интеграции Loginom с ClearML:

1. Библиотека компонентов **ClearML Kit**
2. Шаблоны сценариев

Для интеграции используется python-библиотека **clearml**





Компонент **Новый проект** создает папку с проектом, подпапки **raw-data** и **samples**, а также два шаблона сценария

Модули шаблона сценария **data-prep.lgp**

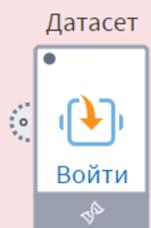




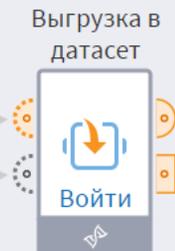
Сохранение в родительский датасет оригинальных файлов с данными (папка raw-data).

По умолчанию данный сценарий также сохраняется в эту папку (опция отключается установкой переменной "Добавить текущий сценарий" в FALSE).

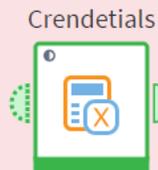
Имя проекта и датасет
Версия - опционально



Выгрузка оригинальных файлов в датасет
(родительский датасет)



Параметры авторизации к ClearML



Проверьте переменные



Компонент загружает все файлы из папки **raw-data** на сервер ClearML

Environment Variables • Быстрый просмотр

№	Метка	Значение
1	АВТОРИЗАЦИЯ	
2	CLEARML_WEB_HOST	https://app.clear.ml
3	CLEARML_API_HOST	https://api.clear.ml
4	CLEARML_FILES_HOST	https://files.clear.ml
5	Ключ доступа	PWD228F9DB0...
6	Секретный ключ	ecGEv32iMOhrPBXe0m9PX...
7	ОБЩИЕ ПАРАМЕТРЫ	
8	Каталог для загрузки артефактов	<null>
9	Список расширений	~lgp,.lck
10	Список файлов	<null>
11	Игнорировать	true



OPEN ARCHIVE

VERSIONS LIST SORTED BY

Kaggle-Russia-Real-Estate v 1.0.0 Final

original

Updated 23 days ago • Created by Nikolay Paklin

DETAILS

Kaggle-Russia-Real-Esta...

409.89 MB 23 days ago

VERSION INFO

Kaggle-Russia-Real-Estate v1.0.0

FINAL

ID 2ebd930e...

Parent -

Size 409.89 MB (origi...)

130.79 MB (comp...)

File count 4

Link count 0

FILES CHANGED

Added 4

Modified 0

Removed 0

Size 409.89 MB

Kaggle-Russia-Real-Estate v1.0.0 CONTENT PREVIEW CONSOLE

Tables

summary

id	time	geo_lat	geo_lon	region	building_type	level	levels	rooms	area	kitch
6050000	2018-02-19 20:00:21	59.8058084	30.376141	2661	1	8	10	3	82.6	10.8
8650000	2018-02-27 12:04:54	55.683807	37.297405	81	3	5	24	2	69.1	12
4000000	2018-02-28 15:44:00	56.29525	44.061637	2871	1	5	9	3	66	10
1850000	2018-03-01 11:24:52	44.996132	39.074783	2843	4	12	16	2	38	5
5450000	2018-03-01 17:42:43	55.918767	37.984642	81	3	13	14	2	60	10

Оригинальный или «сырой» датасет (тэг original)

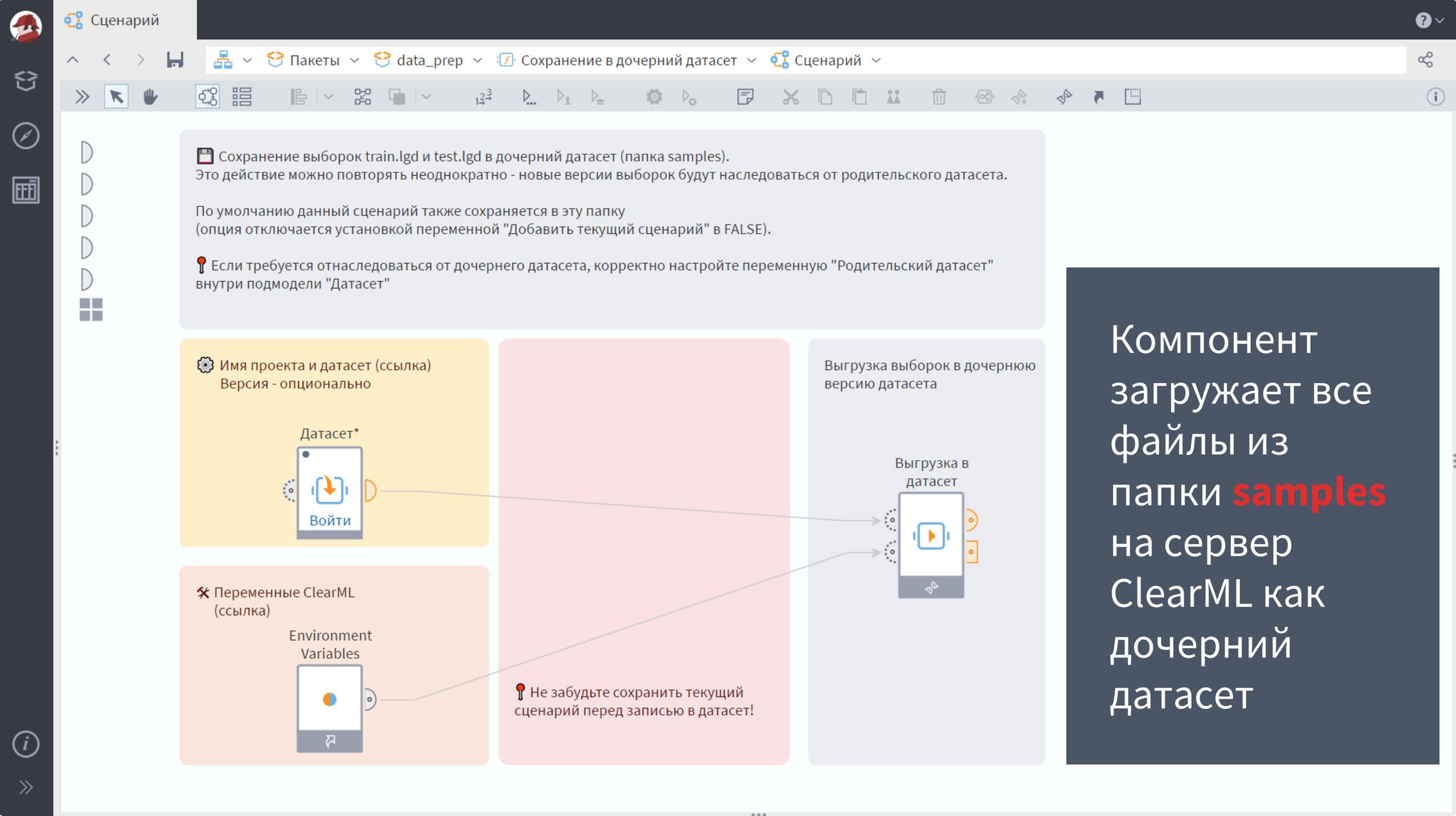


✎ Напишите сценарий подготовки обучающих выборок и выгрузите их в отдельную папку samples.

✳ Настройте узел "Разбиение на множество" в соответствии с вашими задачами.



Аналитику требуется наполнить узлами подмодель **DataPrep**

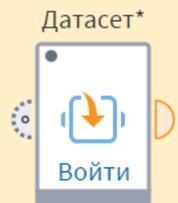


Сохранение выборок train.lgd и test.lgd в дочерний датасет (папка samples).
Это действие можно повторять неоднократно - новые версии выборок будут наследоваться от родительского датасета.

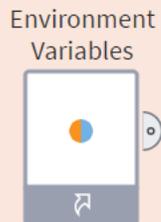
По умолчанию данный сценарий также сохраняется в эту папку
(опция отключается установкой переменной "Добавить текущий сценарий" в FALSE).

Если требуется отнаследоваться от дочернего датасета, корректно настройте переменную "Родительский датасет" внутри подмодели "Датасет"

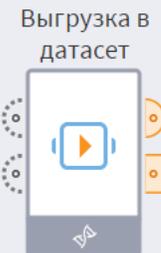
Имя проекта и датасет (ссылка)
Версия - опционально



Переменные ClearML
(ссылка)



Выгрузка выборок в дочернюю
версию датасета



Не забудьте сохранить текущий сценарий перед записью в датасет!

Компонент
загружает все
файлы из
папки **samples**
на сервер
ClearML как
дочерний
датасет

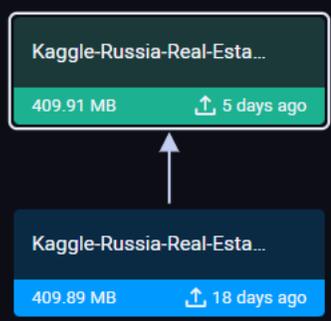


OPEN ARCHIVE

VERSIONS LIST SORTED BY

- Kaggle-Russia-Real-Estate v 1.0.2** Final
 - preprocessed с геофичей
 - Updated 5 days ago Created by Nikolay Paklin
- Kaggle-Russia-Real-Estate v 1.0.1** Final
 - preprocessed
 - Updated 18 days ago Created by Nikolay Paklin
- Kaggle-Russia-Real-Estate v 1.0.0** Final
 - original
 - Updated 18 days ago Created by Nikolay Paklin

DETAILS



VERSION INFO

Kaggle-Russia-Real-Estate v1.0.2
FINAL

ID: 0b107a96...

Parent: -

Size: 409.91 MB (original)
130.85 MB (compressed)

File count: 4

Link count: 0

FILES CHANGED

Added: 0

Modified: 1

Removed: 0

Size: 20.35 KB

Датасеты наследуются друг от друга, имеют версии и тэги

Kaggle-Russia-Real-Estate v1.0.2

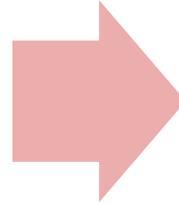
CONTENT PREVIEW CONSOLE

File Name (4 files)	File Size (total 409.91 MB)	Hash (SHA2)
raw-data/all_v2.csv	389.27 MB	46390bd7f9ac41375e7151f9691e83910c9621f...
raw-data/all_v2_1region.lgd	20.58 MB	ac7a2370361872b2c0771280b19b4d945ad6a9...
raw-data/data-prep.lgp	64.84 KB	63cf48e0ac753ec73753be7bceec7b5f7377a9c1...
raw-data/russia-real-estate-20182021_link.txt	71 B	ec552cafe63fcee88a535a25bde2673b81d0110...



Модули шаблона сценария **experiment.lgp**

Датасет



Эксперимент

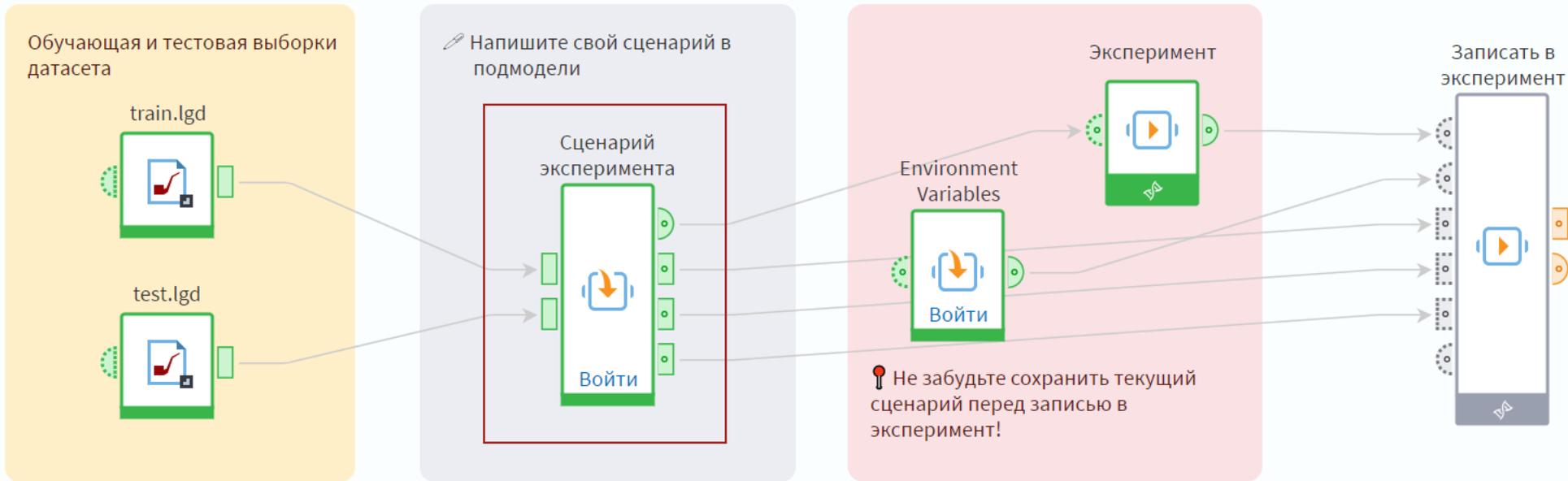
Загрузка локальных копий дочерних выборок **train.lgd** и **test.lgd** или проверка хэша уже имеющихся

Построение модели и запись в эксперимент:

- сценарий эксперимента
- среда эксперимента (версия Loginom, пакеты-ссылки, версии библиотек...)
- гиперпараметры
- метрики
- диаграмма

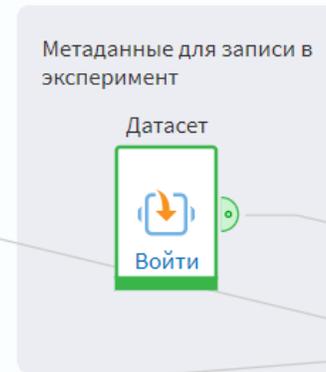
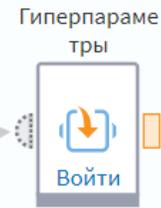
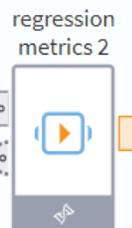
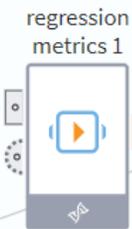
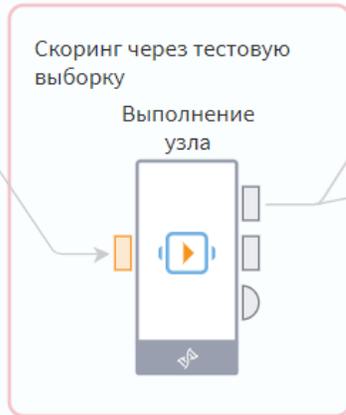
✎ Напишите сценарий эксперимента, используя файлы выборок train и test из папки samples.

✘ Сохраните сценарий и метаданные в эксперимент.



Аналитик наполняет узлами подмодель **Сценарий эксперимента**. Ряд показателей (версия Loginot, настройки узла-модели) собирается автоматически

Выйти из подмодели



Пример сценария эксперимента №1. Задача оценки стоимости недвижимости, датасет **Kaggle Russia Real Estate 2018-2021**, модель линейной ридж-регрессии



+ NEW EXPERIMENT

OPEN ARCHIVE



EXPERIMENTS LIST

SORTED BY

Linear regression (15 features) Completed

с геофичей Санкт-Петербург и область

Updated 10 days ago Created by Nikolay Paklin

Linear regression (14 features) Completed

Санкт-Петербург и область

Updated 16 days ago Created by Nikolay Paklin

Linear regression (14 features) Failed

Санкт-Петербург и область

Updated 16 days ago Created by Nikolay Paklin

Карточка эксперимента в ClearML: версия Loginom, датасет и гиперпараметры

COMPLETED

Linear regression (15 features) ID d2116cda...

с геофичей Санкт-Петербург и область

EXECUTION CONFIGURATION ARTIFACTS INFO CONSOLE SCALARS PLOTS

USER PROPERTIES

Properties	USER PROPERTIES
	Версия Loginom 7.1.5
	Редакция Loginom Standard

HYPERPARAMETERS

Датасет

Датасет	ДАТАСЕТ
	Датасет Id 0b107a9657854cb9bb950355f0950715
	Имя датасета Kaggle-Russia-Real-Estate

Модель

Модель	МОДЕЛЬ
	Алгоритм Линейная регрессия
	Нормализация False
	Регуляризация Ridge

Сводка



+ NEW EXPERIMENT

OPEN ARCHIVE



EXPERIMENTS LIST

SORTED BY

Linear regression (15 features) Completed

с геофичей Санкт-Петербург и область

Updated 10 days ago Created by Nikolay Paklin

Linear regression (14 features) Completed

Санкт-Петербург и область

Updated 16 days ago Created by Nikolay Paklin

Linear regression (14 features) Failed

Санкт-Петербург и область

Updated 16 days ago Created by Nikolay Paklin

Completed

Linear regression (15 features) ID: d2116cda...

с геофичей Санкт-Петербург и область

EXECUTION CONFIGURATION ARTIFACTS INFO CONSOLE SCALARS PLOTS

OTHER

experiment.lgp	https://files.clear.ml/Kaggle%20Real%20Estate/%25D0%259B%25D0%25B8%25D0%25BD%25D0%25B5%25D0%25B9%25D0%25BD%25D0%25B0%25D1%258F%20%25D1%2580%25D0%25B5%25D0%25B3%25D1%2580%25D0%25B5%25D1%2581%25D1%2581%25D0%25B8%25D1%258F%20%2815%20%25D0%25BF%25D1%2580%25D0%25B8%25D0%25B7%25D0%25BD%25D0%25B0%25D0%25BA%25D0%25BE%25D0%25B2%29.d2116cdacd534e5abcb84637eee6ac46/artifacts/experiment.lgp/experiment.lgp
FILE PATH	
FILE SIZE	66.91 KB
HASH	e8b83d15507446f38b9e3f423043d191ddd8062ddf71692f7858f550c896e7de

PREVIEW

```
experiment.lgp - 68.51 KB
```

Карточка эксперимента в ClearML: сценарий эксперимента (файл *.lgp)





+ NEW EXPERIMENT OPEN ARCHIVE

EXPERIMENTS LIST SORTED BY

Linear regression (15 features) Completed

с геофичей Санкт-Петербург и область

Updated 10 days ago Created by Nikolay Paklin

Linear regression (14 features) Completed

Санкт-Петербург и область

Updated 16 days ago Created by Nikolay Paklin

Linear regression (14 features) Failed

Санкт-Петербург и область

Updated 16 days ago Created by Nikolay Paklin

COMPLETED



Линейная регрессия (15 признаков)

ID d2116cda...

с геофичей Санкт-Петербург и область

EXECUTION CONFIGURATION ARTIFACTS INFO CONSOLE SCALARS PLOTS DEBUG SAMPLES



Summary

R ² (train)	MAPE (train)	R ² (test)	MAPE (test)
0.7191002859318962	0.2139850207170322	0.643921859032758	0.20300311764034515

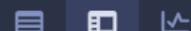
Карточка эксперимента в ClearML: метрики





+ NEW EXPERIMENT

OPEN ARCHIVE



EXPERIMENTS LIST

SORTED BY

Linear regression (15 features) ✓ Completed

с геофичей Санкт-Петербург и область

Updated 10 days ago • Created by Nikolay Paklin

Linear regression (14 features) ✓ Completed

Санкт-Петербург и область

Updated 16 days ago • Created by Nikolay Paklin

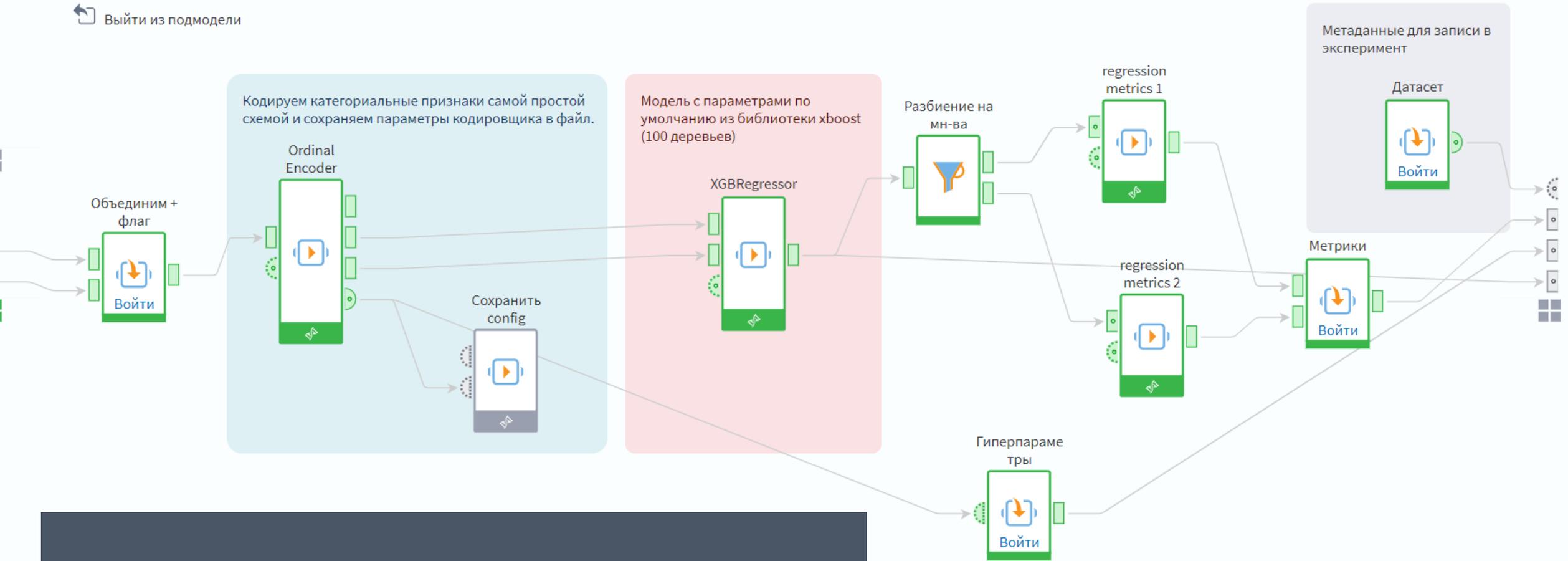
Linear regression (14 features) ✗ Failed

Санкт-Петербург и область

Updated 16 days ago • Created by Nikolay Paklin

Карточка эксперимента в ClearML: диаграмма рассеяния





Эксперимент №2. Модель **XGBoostRegressor**, построена в Loginom без кода на основе библиотеки КОМПОНЕНТОВ **Python Kits**



EXPERIMENTS + Values [dropdown] [up/down arrows] Hide Identical Fields [toggle] Type to search [input] [search icon] [dropdown]

Kaggle Russia Real Estate

XGBoost Regressor (15 признаков) ID 21212ee1...

с геофичей Санкт-Петербург и область

Published · 0 iterations · Last updated 10 days ago

▼ ДАТАСЕТ

Датасет Id	0b107a9657854cb9bb950355f0950715
Имя датасета	Kaggle-Russia-Real-Estate

> КОДИРОВЩИК

▼ МОДЕЛЬ

Алгоритм	XGBoost Regressor
----------	-------------------

▼ PROPERTIES

Версия Loginom	7.1.5
Редакция Loginom	Standard
Ссылка-пакет 1	loginom_clearml_kit ^1.0.0
Ссылка-пакет 2	loginom_silver_kit ^3.1.2
Ссылка-пакет 3	loginom_categoty_kit ^3.1.0
Ссылка-пакет 4	loginom_sklearn_kit ^3.1.0
Ссылка-пакет 5	data_prep ^1.0.0

Kaggle Russia Real Estate

Линейная регрессия (15 признаков) ID d2116cda...

с геофичей Санкт-Петербург и область

Completed ✓ · 0 iterations · Last updated 10 days ago

▼ ДАТАСЕТ

Датасет Id	0b107a9657854cb9bb950355f0950715
Имя датасета	Kaggle-Russia-Real-Estate

▼ МОДЕЛЬ

Алгоритм	Линейная регрессия
Нормализация	False
Регуляризация	Ridge

> СВОДКА

▼ PROPERTIES

Версия Loginom	7.1.5
Редакция Loginom	Standard

Сравнение экспериментов:
версии Loginom, датасеты и гиперпараметры



EXPERIMENTS



Graph



Group by

Metric + Variant

Horizontal Axis

Iterations

Smoothing



0

Exponential Moving Ave...

Find scalars



Hide all

Summary



MAPE (test)

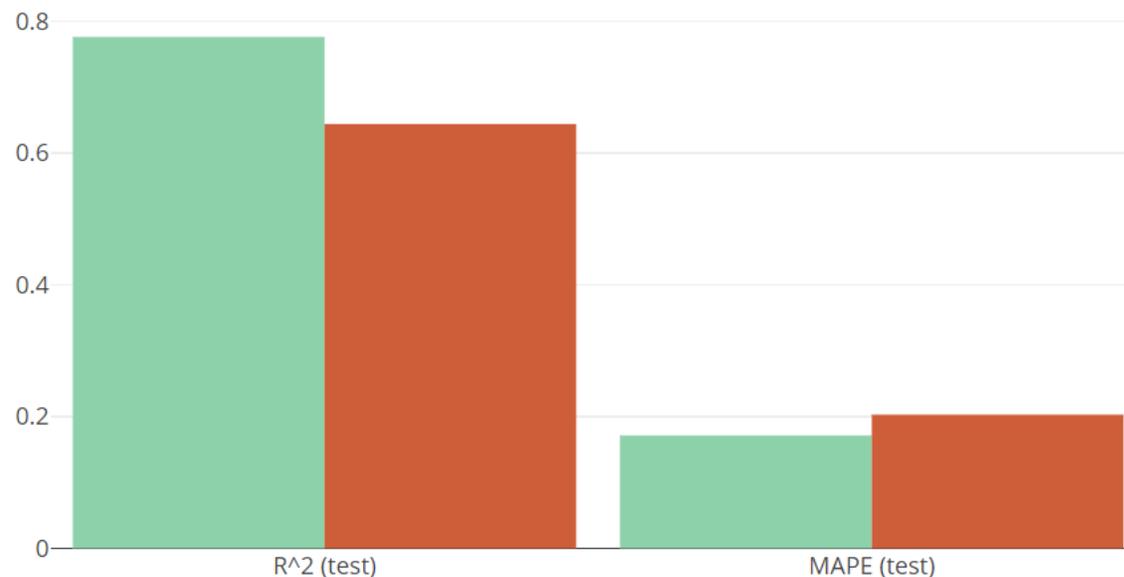


R² (test)



Сравнение экспериментов: метрики

Summary



Series



XGBoost Regressor (15 признаков)



Линейная регрессия (15 признаков)





EXPERIMENTS +

Find plots

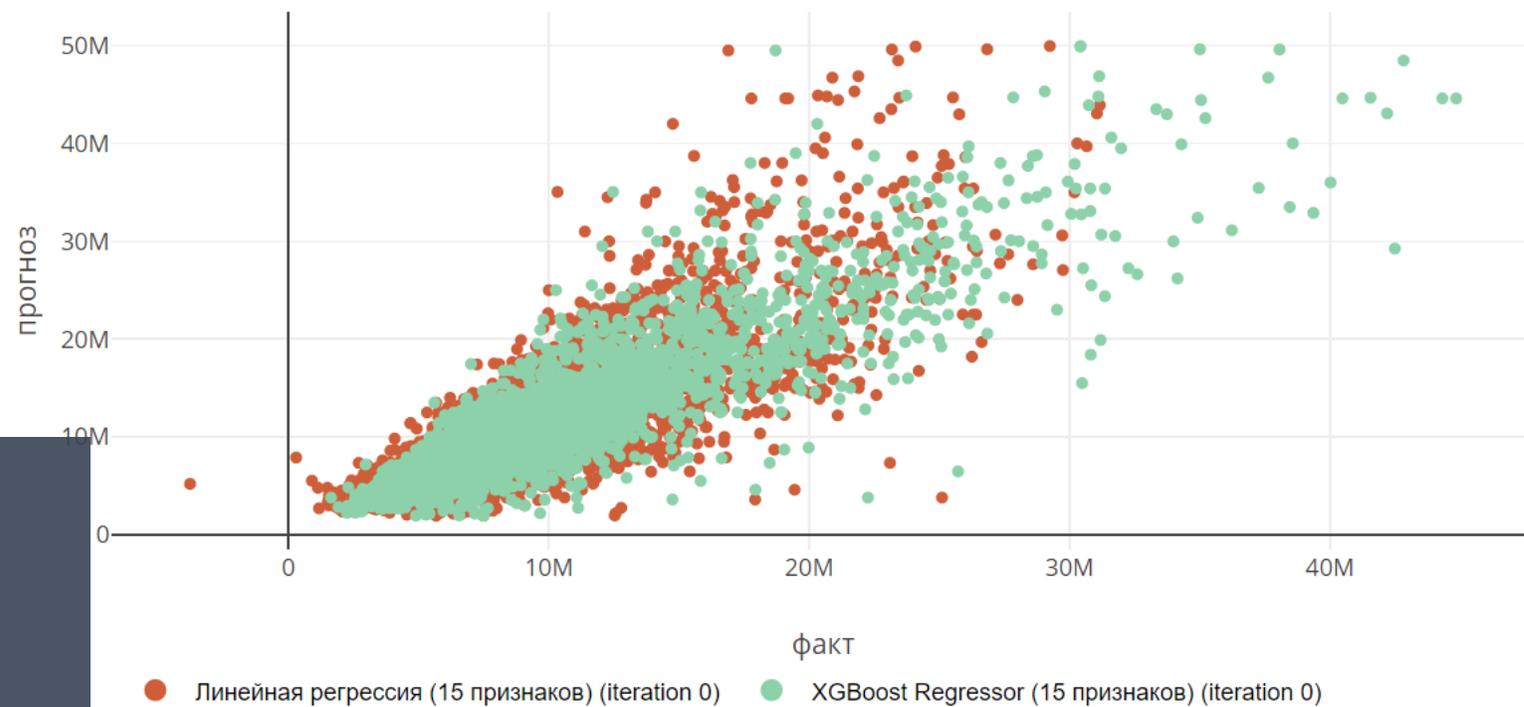


Hide all

Диаграмма рассеяния - Серия 1



Диаграмма рассеяния - Серия 1



Сравнение экспериментов: диаграммы рассеяния



RECENT ▾



OPEN ARCHIVE



NEW REPORT

Подготовка данных 📁



Kaggle Russia Real Estate • Nikolay Paklin

Updated 2 minutes ago

Draft

Описание последовательности действий по подготовке обучающей и тестовой выборок

[документация](#)

ML-модели 📁



Kaggle Russia Real Estate • Nikolay Paklin

Updated 2 days ago

Draft

Описание построенных ML-моделей

[документация](#)

В каждом проекте
создаются отчеты в
формате markdown

< Описание последовательности действий по подготовке обучающей и тестовой выборки

Description

Описание сценария

Импорт "сырых" данных

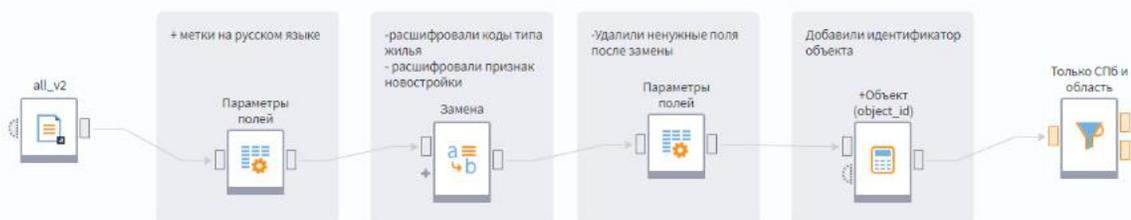
Датасет объемом ~5,5 млн. записей состоит из 540 тыс уникальных объектов с популярных порталов по продаже недвижимости в России из разных регионов. Набор данных содержит информацию о месторасположении дома, материале, из которого он построен (кирпичный, панельный, деревянный и т.д.), количестве этажей, площади квартиры и ее стоимости.

Набор данных состоит из 13 полей:

В датасете присутствуют выбросы и ошибки ввода данных.

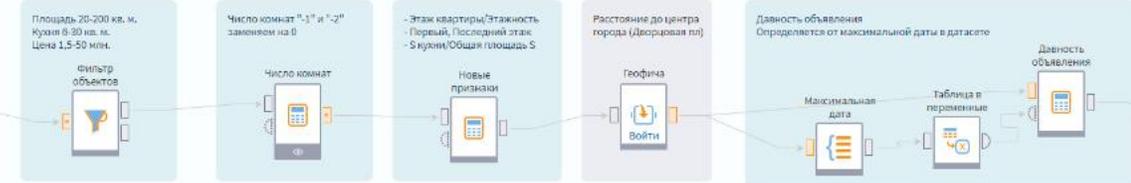
Сценарий подготовки данных представлен на рисунке ниже.

Сценарий преобразования исходного csv-файла all_v2.csv в Loginom Data файл (~5,5 млн. записей).
Датасет по недвижимости Russia Real Estate 2018-2021 можно скачать здесь: <https://www.kaggle.com/datasets/mrdaniilak/russia-real-estate-20182021>



Он вносит косметические изменения в "сырой" датасет и отфильтровывает записи по региону "СПб и область".

Подготовка выборки Для подготовки обучающей и тестовой выборки используется следующий сценарий.



В нем мы последовательно рассчитываем дополнительные показатели, но первый шаг - очистка от выбросов. Квартиры с общими площадями менее 20 кв. и более 200 метров признаются выбросами.

Пример отчета.
После публикации
его увидят все
пользователи



ЛЕТО

2024

- Релиз библиотеки компонентов **ClearML Kit**
- Вебинар по ModelOps

Мастерская Loginom Skills

Онлайн-воркшопы с экспертом

Аналитика
данных

Машинное обучение
и моделирование

Эффективный
зерокодинг



Мастерская
Loginom Skills

Проект **Мастерская Loginom Skills** –
воркшопы, мастер-классы, навыки.



https://t.me/loginom_skills

Спасибо за внимание!