

Решение O! хакатона

от команды FULLSTACK EXCEL

O! Хакатон. Кейс 2: автоматическое тегирование номеров отелей

Описание решения

ЯДРО РЕШЕНИЯ:

Использование современных высокопроизводительных подходов в рамках **NLP** и **машинного обучения**.

ОЧИСТКА ДАННЫХ

Очистка от мусорных символов с помощью регулярных выражений. **Эксперименты:** стемминг и лемматизация.

БАЗОВЫЕ МОДЕЛИ

В рамках единого пайплайна для каждого таргета по отдельности выбирались:

- 1. Векторайзеры для текста (**TF-IDF** или **count-vectors**);
- 2. Легкие модели (**логистическая регрессия, SVM**) и их **гиперпараметры** (*optuna*).

Эти модели хорошо закрывают “простые” таргеты на основе подбираемых ими ключевых слов

DEEP-LEARNING МОДЕЛИ

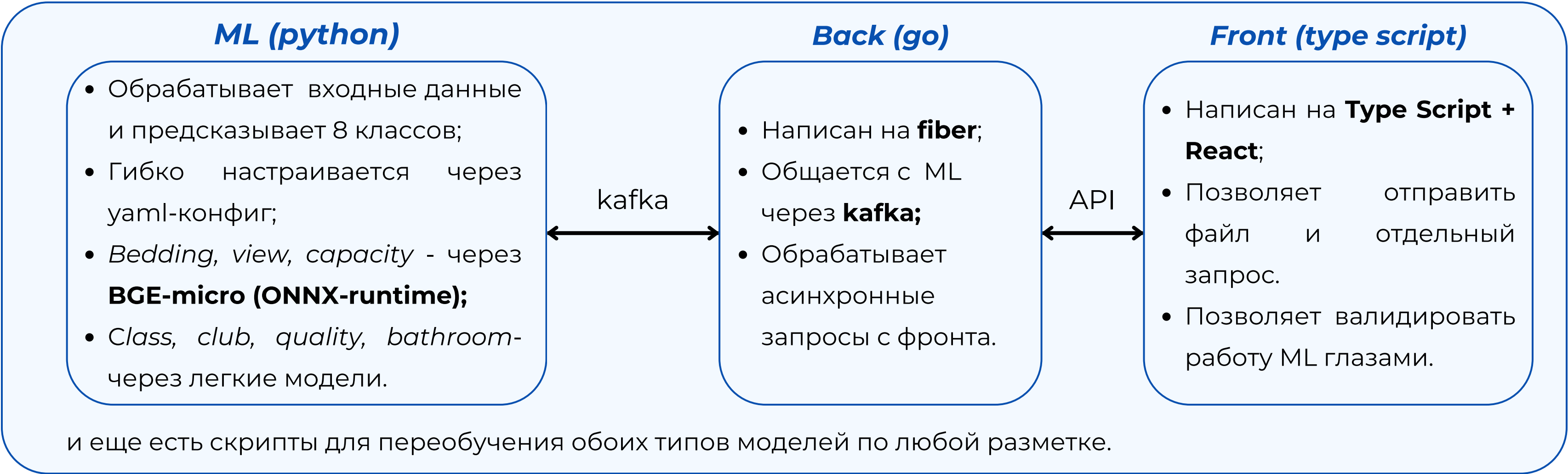
- 1. Для каждого таргета обучалась легковесная **bert-like** модели на базе **энкодера bge-micro** (дистилляция bge-small);
- 2. Проводились эксперименты по **multilabel классификации** (интересное продуктивное решение);
- 3. Инференс через **ONNX-runtime** ради графовой оптимизации и роста скорости;
- 4. **ostrovok-BGE-micro** закрывает более семантически сложные кейсы чем бейслайн.

Архитектура решения

ПОЧЕМУ МЫ?

Помимо самого алгоритма, для его легкой интеграции мы реализовали **API** на **GO (Fiber)** позволяющее работать как с csv файлами, так и со стандартными json-запросами.
А еще мы реализовали **UI-интерфейс**, в теории необходимый для удобной отладки работы алгоритма силами непрофессионалов.

Docker-Compose (можно поднять весь сервис за один *make up*)



Примеры интерфейса

ПОМНИТЕ

Мы говорили про реализованный интерфейс. Это он. Он позволяет вам оценить работу модели глазами. В рамках доработок планируется выводить ее. Полное видео лежит в репозитории и по QR коду.

ОСТРОВОК
ХАКАТОН

Гость

Генерация по строке

Генерация по файлу

О нас

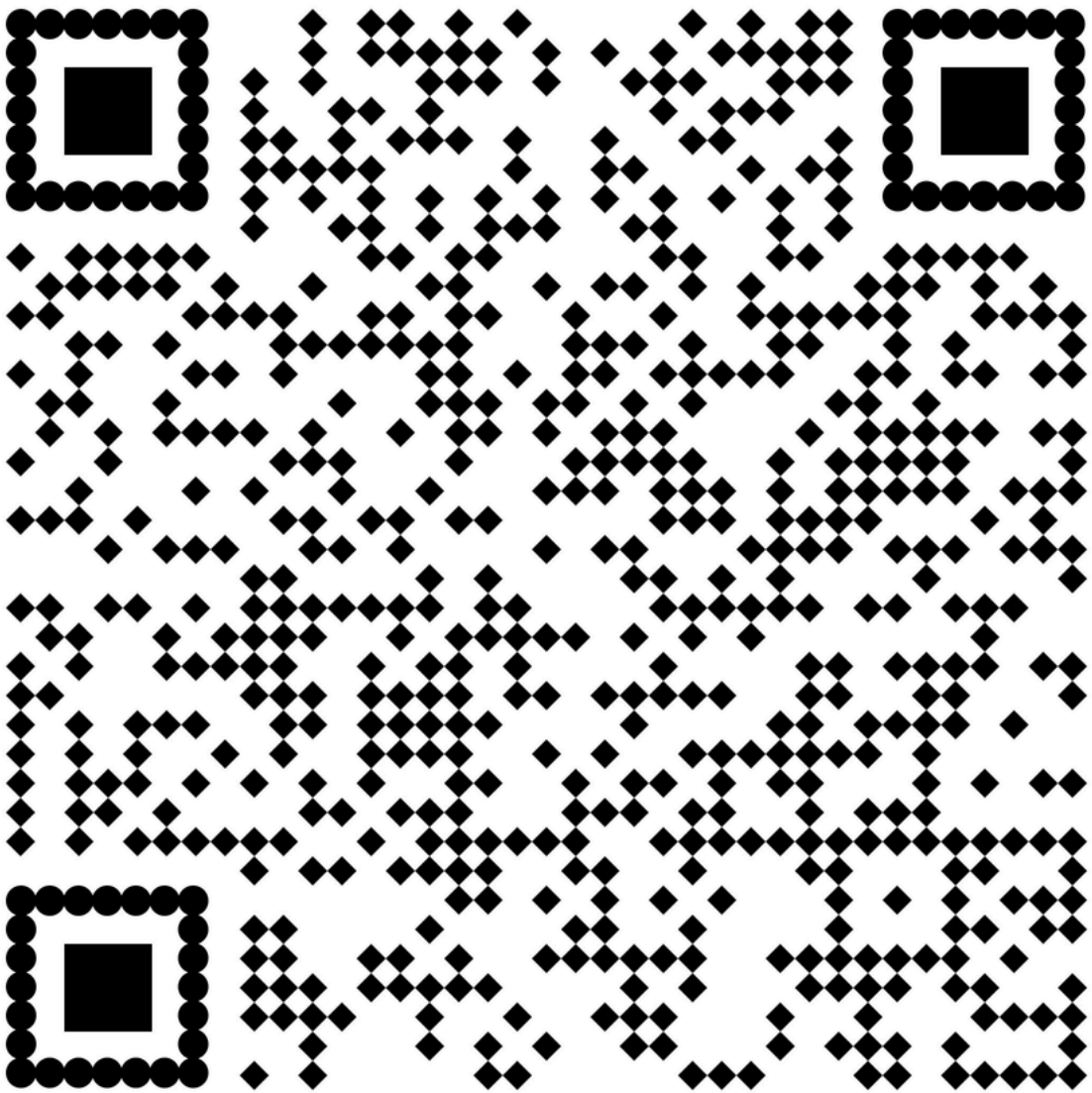
Результат

СКАЧАТЬ ТАБЛИЦУ

Таблица

rate_name	class	quality	bathroom	bedding	capacity	club	bedrooms	balcony	view
executive suite 1 bedroom mountain view 1 king bed	suite	executive	private bathroom	double/double-or-twin	double	not club	1 bedroom	undefined	mountain view
standard double bed sea view room sea	room	standard	private bathroom	double/double-or-twin	double	not club	undefined	undefined	sea view
twin double room executive sea view twin double room executive sea view sea	room	executive	private bathroom	double/double-or-twin	double	not club	undefined	undefined	sea view
standard standard double room 2 adults	room	standard	private bathroom	double/double-or-twin	double	not club	undefined	undefined	undefined
quadruple room with balcony and sea view	room	standard	private bathroom	undefined	quadruple	not club	undefined	undefined	sea view
family room bunk beds all inclusive	room	standard	private bathroom	bunk bed	double	not club	undefined	undefined	undefined
small	room	standard	private bathroom	undefined	undefined	not club	undefined	undefined	undefined
twin superior room twin superior room twin superior room 1 twin	room	superior	private bathroom	double/double-or-twin	double	not club	undefined	undefined	undefined

← 1 →



Примеры работы системы

4

Показательно

Визуально сложные примеры того, как наша модель отработала на грязных данных.

`rate_name`	`dirty['...']`	`predict`
Standard Single Bed in 8-Bed Dormitory Room (En Suite Share Bathroom)	suite (class)	dorm
Double Deluxe Flexible Room Only 200\$	run-of-house	room

Пояснение:
En Suite Bathroom - совмещенный санузел
Flexible room - тип зонирования и планировки

Точность

98.5-99.9%

на отложенной выборке в зависимости от колонки

СПАСИБО ЗА ВНИМАНИЕ!

Команда:

Игитов Максим - team-lead, DS. **НИУ ИТМО**

Богодист Всеволод - ML-lead, MLE. **НИУ ИТМО**

Грушевский Георгий - MLE, MLops. **НИУ ИТМО**

Вичук Артем - back-end, dev-ops. **НИУ ИТМО**

Горошко Андрей - front-end, ux-ui. **СПБГУ**