

Тестовое задание для стажеров в команду Data Science IDP

Задание — написать решение по извлечению сущностей из документов (новостных текстов). Выполните задание в Jupyter Notebook. Ожидаемый результат: ipynb-файл с решением и всеми выводами ячеек, csv-файл с предсказаниями модели залитый в репозиторий на github.

В качестве датасета возьмите русские новости из [Balto-Slavic Natural Language Processing 2019 \(helsinki.fi\)](https://helsinki.fi/Balto-Slavic-Natural-Language-Processing-2019/Guidelines_20190122.pdf). Интересующие сущности: PER, ORG, LOC, EVT, PRO (см. [Guidelines_20190122.pdf \(helsinki.fi\)](https://helsinki.fi/Balto-Slavic-Natural-Language-Processing-2019/Guidelines_20190122.pdf)).

Достаточно использовать 9 документов про брекзит из предложенного организаторами семпла.

Пример одного документа:

ru-10

ru

2018-09-20

<https://rg.ru/2018/09/20/tereza-mej-rasschityvaet-usidet-v-sedle-do-zaversheniia-procedury-brexit.html>

Тереза Мэй рассчитывает усидеть в седле до завершения процедуры Brexit

Тем не менее, по сведениям британских СМИ, на предстоящей в конце сентября конференции партии тори противники Мэй навяжут ей серьезный бой, из которого не факт, что она выйдет победителем. Фаворит букмекеров в качестве возможного сменщика нынешнего премьера, бывший министр иностранных дел Британии Борис Джонсон намерен выступить с альтернативным докладом, который не оставит камня на камне от взглядов главы правительства на условия "брекзита". С точки зрения Джонсона, "Лондон обернул британскую конституцию поясом смертника и вручил детонатор Мишелю Барнье (главному переговорщику Брюсселя по "брекзиту". - От ред.)". С этой метафорой и предстоит сразаться на конференции главе правительства Альбиона.

...

Задание 1

Опишите задачу с точки зрения NLP. Что это за задача, какие классические методы существуют для ее решения? Как ее можно решать через LLM? Как обычно оценивают качество моделей в этой задаче?

Задание 2

Реализуйте чтение датасета в pandas DataFrame с обязательными колонками "document_id", "document_text", "entity", "gold_answer". Выведите шапку датафрейма.

Задание 3

Напишите функцию, которая принимает на вход строку датафрейма и выдает текст входного сообщения для LLM.

Задание 4

Получите ответы GigaChat для всех документов. Документов всего 9, поэтому сделать это можно вручную, пользуясь веб-интерфейсом GigaChat или ботом в VK или Телеграме. Не очищайте историю сообщений, чтобы потом продемонстрировать подлинность ответов на онлайн-собеседовании.

Внесите ответы GigaChat в датафрейм, сохраните его.

Задание 5

Реализуйте самостоятельно алгоритм для подсчета метрик `score_fn(gold: str, pred: str) → float`. Можно пользоваться только библиотеками `numpy`, `scipy`, `pandas`. Напишите юнит-тесты. Возможно ли ускорить вычисление функции через векторную реализацию? Поясните решение и обоснуйте, почему выбрали именно такую метрику.

Задание 6

Вычислите метрики для каждой строки в датафрейме. С агрегируйте результаты а) по каждой сущности, б) по каждому документу. Изобразите результаты на графиках. Какие выводы можно сделать?

Задание 7

Есть ли зависимость метрик от длины документа? Постройте графики, чтобы ответить на вопрос.

Задание 8

Проведите анализ ошибок. Когда модель чаще отвечает правильно, а когда ошибается? Предложите варианты, как повысить метрики.

Задание 9

Сделайте выводы по всему исследованию. Напишите, чему научились и что нового попробовали.