

# Отчет

Штыков Павел<sup>1</sup>

shtykov.pa@gmail.com

15 июня 2022

## 1 Анализ исходных датасетов

Скрипт с кодом анализа исходных датасетов: **part1\_analysis.py**.

Анализ заключался в составлении списка всевозможных типов сущностей, их примеров и проверки на вложенность. Результаты работы лежат в папке **res/par1**.

Было выявлено, что в датасете Restaurant8k существует небольшое количество предложений, сущности в которых накладываются. При обучении они были отброшены. Также в датасете Restaurant8k есть "битые" предложения (с поломанным форматом json'a например). При обучении это было учтено.

## 2 Объединение датасетов

Скрипт с кодом анализа исходных датасетов: **part2\_merge.py**.

В двух данных датасетах логично совпадают только два типа сущностей: *date* и *time*. Соответственно в датасете ATIS все типы относящиеся к дате и времени были обобщены до указанных выше *date* и *time* (например: *return\_date.day\_number* → *date*). В датасете Restaurant8k сущности *date* и *time* присутствуют изначально.

Кроме этих двух "простых" сущностей в объединенный датасет были добавлены сущности *first\_name*, *last\_name* и *people* из Restaurant8k (соответственно из данного датасета в объединенный датасет были перенесены все сущности). Такой выбор был сделан в силу того, что сущности описывающие человека очень общие и могут встретиться почти в любом домене.

Из датасета ATIS дополнительно была взята сущность *city*, т.к. она также является крайне общей. Кроме того название нескольких английских городов встречается и в датасете Restaurant8k.

Все сущности, связанные с самолетами, были отброшены. Сущности связанные с административным делением США (имя штата, почтовый код и т.д.) также были отброшены, т.к. датасет Restaurant8k о Великобритании.

В файле **tils/entities\_map.json** можно найти полный словарь соответствий между старыми и новыми типами сущностей.

В файле **res/part2/union\_stat.json** можно найти все реализации всех типов сущностей в объединенном датасете.

## 3 Обучение моделей

Для обучения были выбраны следующие две модели: Flair Embeddings и W2NER.

В качестве критериев отбора моделей был хороший репозиторий на GitHub (количество звезд больше 100) и высокое качество в бенчмарках PaperWithCode. На CoNLL2003 датасете модели достигают качества 93.09 и 93.07 соответственно.

В качестве метрик были взяты стандартные (и уже реализованные в соответствующих моделях): precision, recall и F1-score (macro averaging).

В таблице 1 представлены лучшие результаты моделей на тестовой выборке. Видно, что обе модели достигли хороших результатов, однако W2NER показала себя заметно лучше.

Модель	Recall	Precision	F1-score
Flair	0.9515	0.9477	0.9495
W2NER	0.9698	0.9576	0.9637

Таблица 1: Лучшие результаты на тестовой выборке

Рассмотрим примеры работы моделей на тестовых данных.

#### Flair:

1. on [april first]<sub>date</sub> i need a flight going from [phoenix]<sub>city</sub> to [san diego]<sub>city</sub>
2. a table for [4]<sub>date</sub> for the [28th of April]<sub>date</sub> please
3. the booking is for [9 people]<sub>people</sub>
4. Thanks for your help, that is all i needed

В целом Flair предсказывает метки хорошо. Во 2-ом предложении есть ошибка с токеном «4», но такая ошибка понятна.

#### W2NER:

1. on [april first]<sub>date</sub> i need a ticket from [tacoma]<sub>city</sub> to [san jose]<sub>city</sub> departing [before 7 am]<sub>time</sub>
2. [8:00 pm]<sub>time</sub> will work for the [12]<sub>people</sub> invited.
3. I would to make a reservation [in three weeks]<sub>date</sub>, for [9 people]<sub>people</sub>
4. I want to book a table for [me and my wife]<sub>people</sub>.
5. Gluten free options are not important, but I want a cheap price range. I don't care much about price range, unless it is cheap or expensive.

W2NER справилась лучше с отдельно стоящим числом, классифицировав его как *people* (2-ое предложение). Качество остальных предсказаний также высокое.

## 4 Не успел

- При объединение нескольких сущностей в одну в датасете ATIS иногда между сущностями встречаются предлоги или артикли. Они разбивают сущность на две части. Возможно, правильное было бы объединить их в одну сплошную сущность.
- Не реализован скрипт для использования уже обученной модели W2NER.
- Не построены кривые обучения (по логам было заметно, что W2NER сходиться быстрее, однако на одну эпоху уходит больше времени)