

Легковесный детектор для токсичного текстового контента

Михаил Мартинсон

29 сентября 2019 г.

Задача

В этом задании требовалось обучить легкий и быстрый классификатор для коротких текстов. Датасет состоит из коротких сообщений людей

Условие

Данные [Kaggle. Toxic Comment Classification Challenge](#)

Репозиторий решения github.com/MartinsonMichael

Основной файл решения [nn_model_simple_rnn.ipynb](#)

Инструкция по запуску

Предполагается наличие pip'a и питона.

Установить зависимости `pip install -r req.txt` в корне клонированного проекта.

Запустить `nn_model_simple_rnn.ipynb`

Код предобработки данных лежит в `view.ipynb`

1 Моя модель

Эксперименты

Так как в задании есть бонусная часть про интерпретируемость модели, вначале я решил использовать в качестве основной модель основанную на деревьях. Немного экспериментов с Random Forest лежат в `jupyter-notebook` [Trees.ipynb](#). В экспериментах видно увеличение точности по мере роста глубины деревьев. Но модель так и не достигнув точности второй, нейронной модели, стала весить критично много (906МБ для 1000 деревьев глубины 80).

Следующий и финальный эксперимент был с простой нейронной архитектурой.

Финальная модель

Модель состоит из byte pair encoder, который переводит строку текста в последовательность токенов. Выбор этого токенизатора обусловлен с одной стороны возможностью держать в словаре целые слова и строить эмбединги для них, а с другой стороны способностью токенизировать любое слово.

Далее для каждого токена строится эмбединг маленького размера, к которому применяется Dense слой для повышения размерности. Это позволяет не создавать много весов для эмбедингов, но в тоже время делать вычисляемый эмбединг слова не сильно маленьким.

Следующий шаг это два слоя двунаправленного LSTM слоя. Во многих работах по языковым моделям говорится, что нижние слои таких последовательных RNN слоев имеют тенденцию улавливать более низкоуровневые, текстовые фичи, а более высокие - более сложные семантические фичи. Для данной задачи хочется как раз более высокоуровневых фичей, но в тоже время нельзя делать модель слишком громоздкой, поэтому слоев два.

Далее еще один общий Dense слой и GlobalMaxPooling, чтобы привести выходы LSTM, которые могут быть разной длины для разных предложений, к вектору фиксированной размерности. И финально еще два Dense слоя выдающие матрицу $ANS = [7, 2]$ для одного текста. Где $ANS[i, 1], i \in \overline{0, 5}$ вероятность тексту получить тег `['toxic', 'severe_toxic', 'obscene', 'threat', 'insult', 'identity_hate']`

А $ANS[6, i], i \in \{0, 1\}$ - дополнительно добавленный таргет, являющийся 1, если хотя бы один из набора основных таргетов 1. Предсказывание этого распределения по-идее должно помочь модели выделять признаки токсичных сообщений.

Точность модели и Недостатки

Данные очень несбалансированные, $\sim 89\%$ данных вообще не содержат ни одной метки. Поэтому метрика accuracy, которую я считаю, почти ни о чем не говорит, хотя и является довольно высокой.

Лейбл	toxic	severe_toxic	obscene	threat	insult	identity_hate
accuracy	0.96	0.99	0.98	0.99	0.97	0.99
gosauc	0.84	0.74	0.86	0.5	0.77	0.5

Среднее gosauc 0.72

Как можно видеть, gosauc для **threat** и **identity_hate** составляет 0.5, что достигалось бы случайным классификатором. То есть модель так и не научилась выделять эти метки.

Улучшения классификатора и ускорение

Прямое следствие предыдущего пункта, улучшение классификации **threat** и **identity_hate**. Для этого можно потренировать модель только на сбалансированной подвыборке этих двух классов. Или сделать дополнительные слои в 'головах' предсказывающих распределения этих двух классов.

Почти все улучшения точности должны касаться либо обучения, либо архитектуры модели. Для архитектуры можно экспериментировать с токенизацией, так как для коротких текстов и легких классификаторов кажется важным как можно лучше образом токенизировать текст. Также можно изменить механизм получения эмбеддингов, например взять уже предобученные эмбеддинги. Еще для уменьшения модели можно заменить тяжелые и большие матрицы из Dense слоев на их приближения разложениями.

Если посмотреть на архитектуру нарисованную в юпитер-ноутбуке, то будет видно, что основной весовой вклад дают RNN слои. С ними связано возможное довольно сильное увеличение скорости. Так как в моделях работающих над последовательностями и имеющими внутри архитектуры RNN довольно долгой частью является работа рекуррентных блоков, так как их можно вычислять только строго последовательно. Поэтому возможно хорошим решением будет использование архитектур типа трансформера. В ней нет RNN блоков и она может хорошо параллеливаться на несколько процессоров. Но для этого придется полностью изменить архитектуру модели.

2 Баланс между качеством и скоростью

Глобально при выборе баланса между качеством и скоростью кажется хорошим решением, понять, что важнее, зафиксировать этот параметр и улучшать оставшийся. Отсюда мне кажется, что при выборе между легкой моделью, или оптимизацией большой, лучше выбрать легкую (конечно с достаточной capacity для решения задачи) и максимально хорошо ее обучить. Ибо при ускорении большой модели придется применить не меньше трюков, чем при обучении маленькой, но нужно еще и обучить большую.