

# Регрессия

Беляков Дмитрий

April 2, 2022

## 1 Описание задачи

Используем датасет YoutubeVideoDataset с базы данных Kaggle. Датасет содержит названия видео, его описание и одну из 6 категорий к которым это видео относится. Возможные категории: Путешествия, наука и технологии, еда, искусство и музыка, производство. Будем пытаться отгадать категорию по названию видео. Для этого используем структуру RandomForest из библиотеки sklearn. Будем использовать 4 способа анализа текста.

- Для каждого названия найдем сколько раз в нем встречается каждая буква алфавита.
- Для каждого названия найдем плотность количества каждой буквы алфавита
- Составим слоарь из 100 самых частых слов в датасете и найдем количество каждого из 100 слов в каждом названии
- Составим слоарь из 100 самых частых слов в датасете и найдем плотность каждого из 100 слов в каждом названии

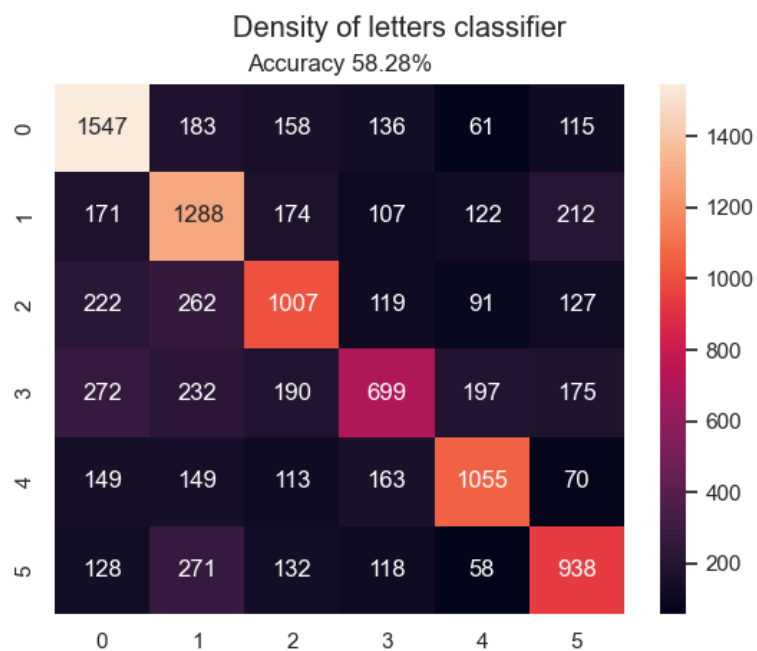
Под плотностью будем подразумевать количество каждой буквы/слова поделить на количество букв/слов в тексте. Используем 20 процентов д

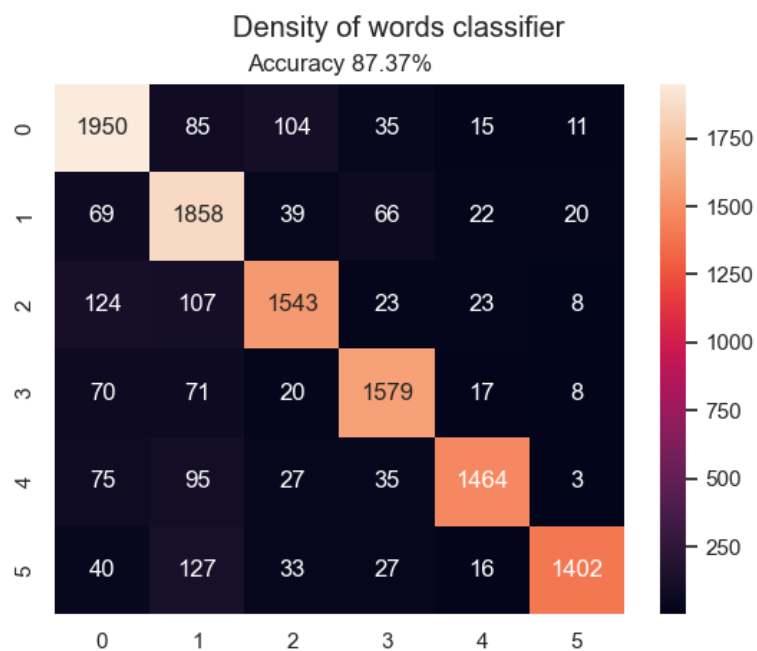
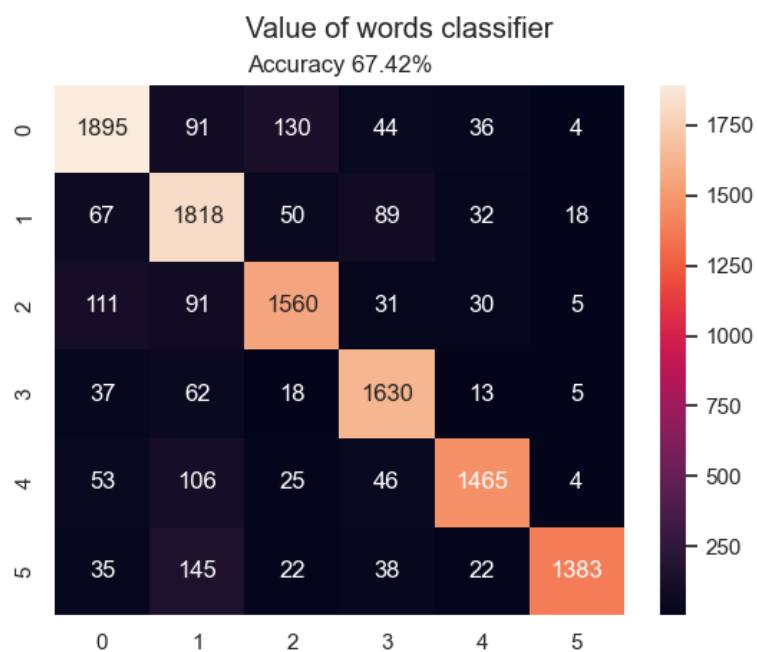
## 2 Описание решения

С помощью библиотеки WordCloud проиллюстрируем 100 самых частых слов во всем датасете



Выведем результаты с помощью тепловой карты





### 3 Вывод

Классификация по частоте слов показала впечатляющий результат почти в 90 процентов. Однако не менее впечатляющим оказалась классификация по количеству и плотности букв, так как при предполагаемой случайности распределения букв, показала аккуратность в 50 процентов, что выше чем 16 процентов если бы распределение было бы случайным