

Toxic Comment Classification

By Minh Ngoc Pham
Data Science Career Track

Project Introduction

- ◆ Toxic Comment Classification Challenge

- ◆ Published in 2018

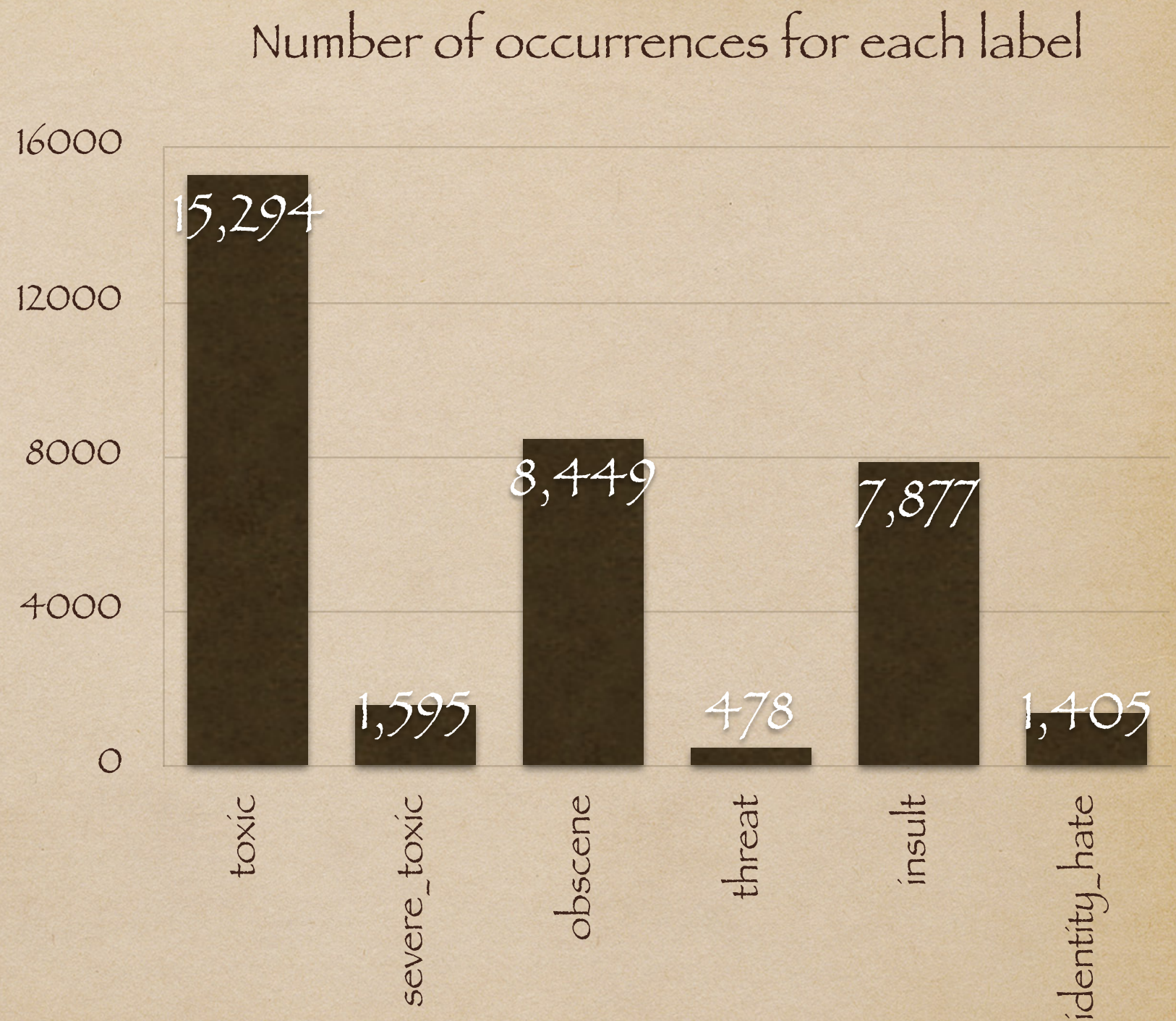
- ◆ 159,571 comments in training set

- ◆ 153,164 comments in test set

kaggle

- ◆ Highest count in toxic comments

- ◆ There are a lot of clean comments





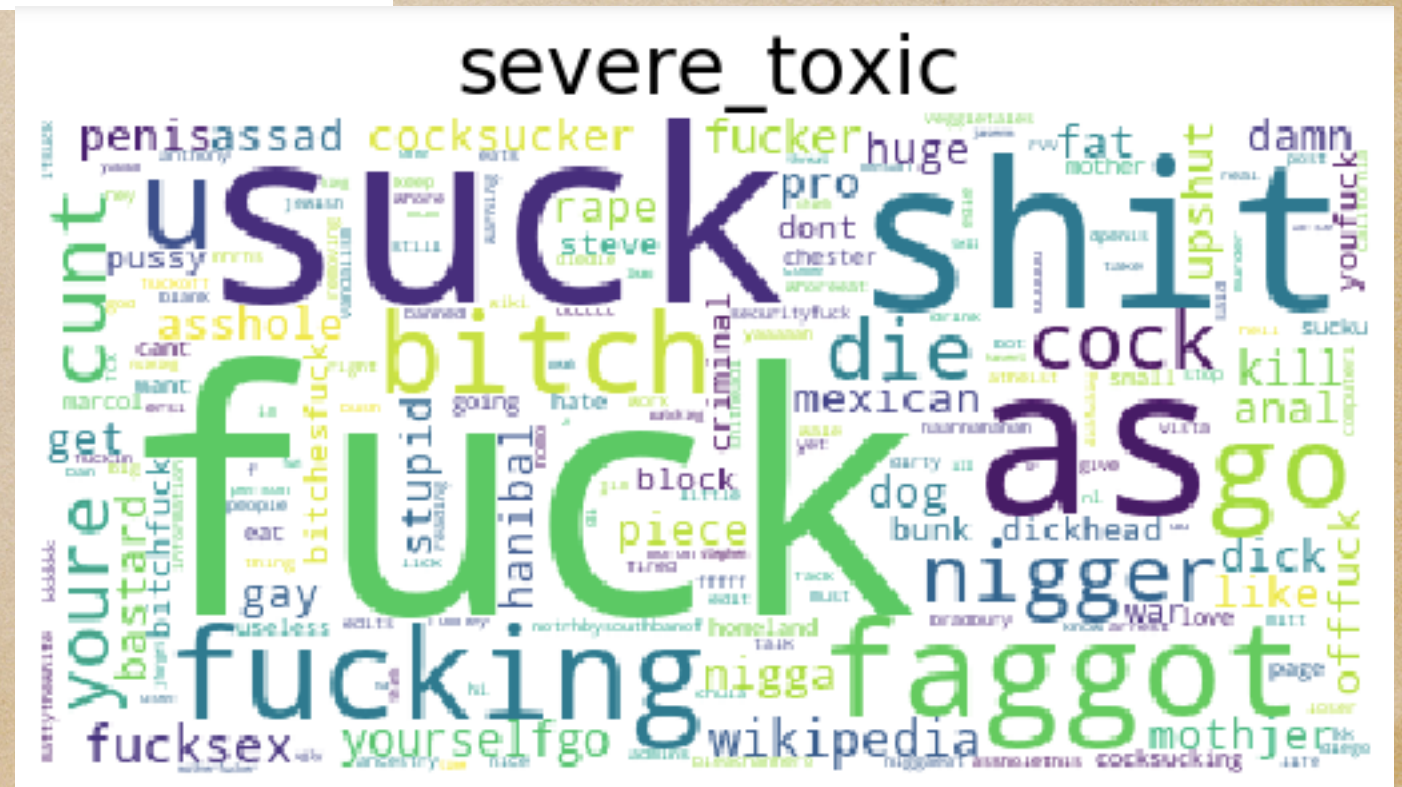
- ◆ 80.3% of the comments in the training test are clean
- ◆ Class imbalance issue

- ◆ High number of comment with 0 label
- ◆ There are only 31 labels classified as all of the 6 labels

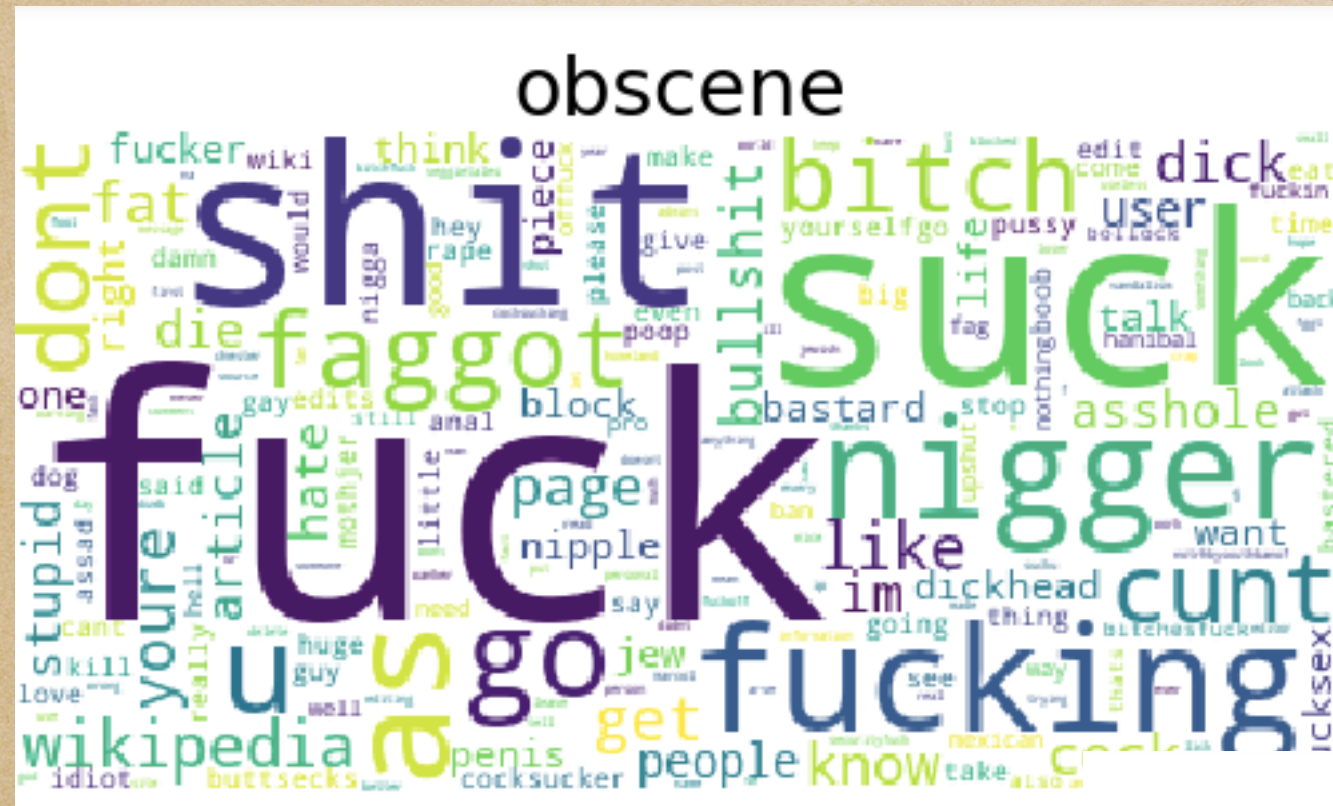
Comment label number



Word Cloud



Word Cloud



Word Cloud

