

# Toxic Language Classification: Cleaning Up Wikipedia

Brian McMahon

12 March 2018

This dataset contains text that may be considered profane, vulgar, or offensive.





## Why is this project important?

Dataset of wikipedia talk page comments from Alphabet-sponsored Kaggle Competition to <u>develop tools to "help improve online</u> <u>conversation"</u>

**Competition Objective** to Identify and classify comments labeled as belonging to one or more of <u>six different "toxic language" categories</u>

In **this presentation** we will conduct preliminary topic analysis on the dataset to draw distinctions amongst <u>the different categories</u> of toxic language



## The Six Category "Tags"

Dataset contains multi-tagged comments; "rating" refers to the number tags a comment has received

\_ | | | High

**Toxic.** General taboo, rude language

**Obscene.** Profanity

Insult. Rude, degrading

**Severe Toxic.** Extreme language, statements not classified as Threat or Identity Hate

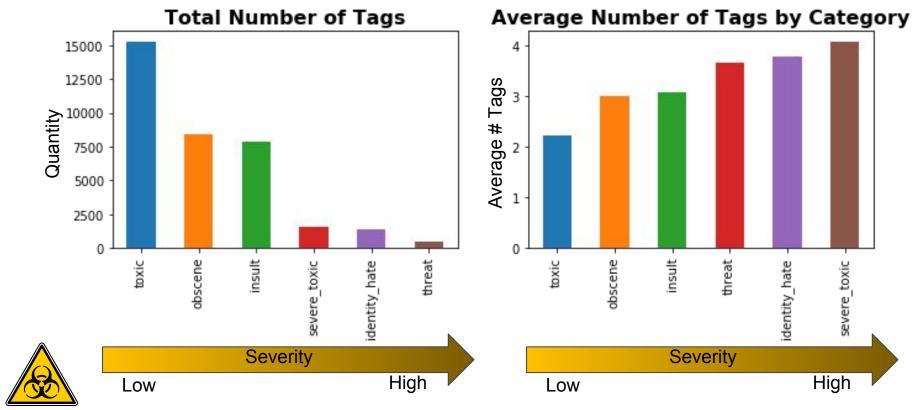
**Identity Hate.** Attacking race, religion, orientation

Threat. Death wishes, threats



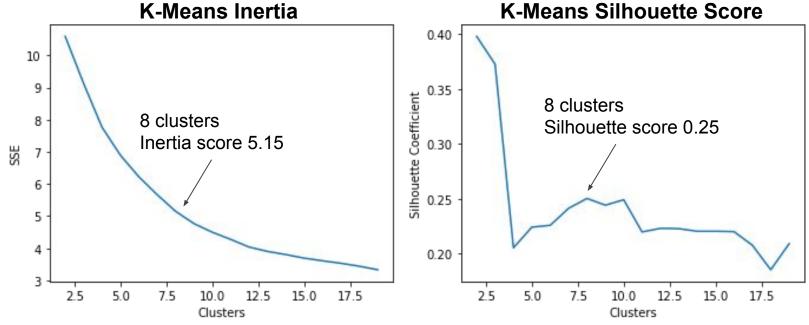
#### **Toxic Comment Mix**

Dataset contains multi-tagged comments; each comment can have from 0 ("clean") to 6 tags "Toxic" data filtered to comments with at least 1 tag



## K Means Clustering: Toxic

We can see some of the classifications within these clusters





## K Means Clustering: Toxic

#### We can see some of the classifications within these clusters











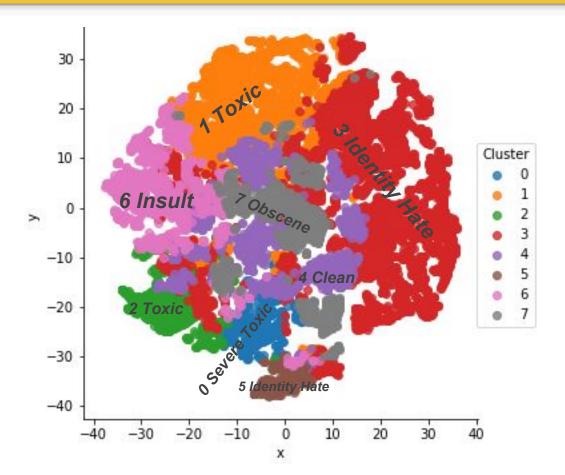








## T-SNE: Toxic





8

## Key Takeaways

Different platforms will have different thresholds for "toxic comments"

Approach toxic language from multiple angles:

- General inappropriate language (i.e., "bad words") may be acceptable on some forums
- Language which threatens or attacks a person's identity ("identity hate")
  does not always use "bad words" and must be recognized in context

While models today do a reasonable job of identifying these issues, the ultimate decision of appropriateness remains subjective and is a challenge for machines to fully grasp



## Questions?

## **APPENDIX**

### **Techniques Utilized**

**TF-IDF.** Used to gauge the word importance within a collection or corpus, increases proportionally to number of times a word appears in the document and is offset by frequency of word in corpus

**NMF.** Matrix factorized into multiple matrices with no negative elements

**K-Means Clustering**. Determine optimized number of clusters by silhouette coefficient

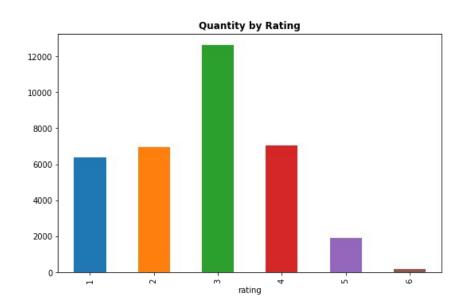
- Silhouette Score. A succinct graphical representation of how well each object lies within its cluster
- Inertia Score. Choose centroids that minimise the inertia, or within-cluster sum of squared criterion

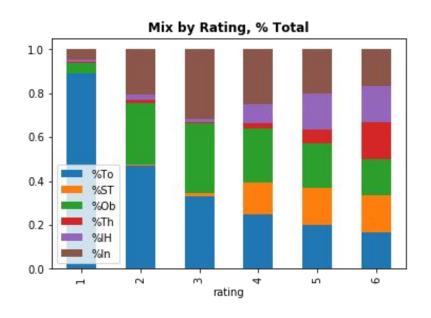
T-SNE. Embeds high-dimensional data into low-dimensional scatter plot

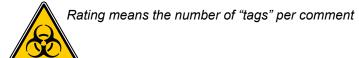


## The Six Categories of Toxic Comments

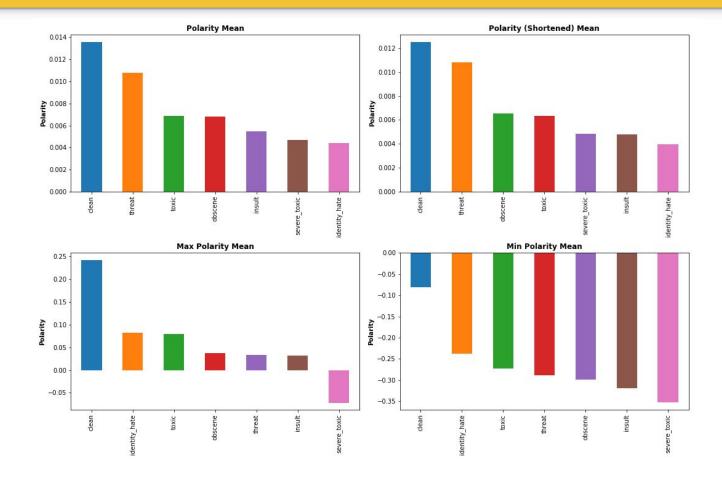
Dataset contains multi-tagged comments; "rating" refers to the number tags a comment has received





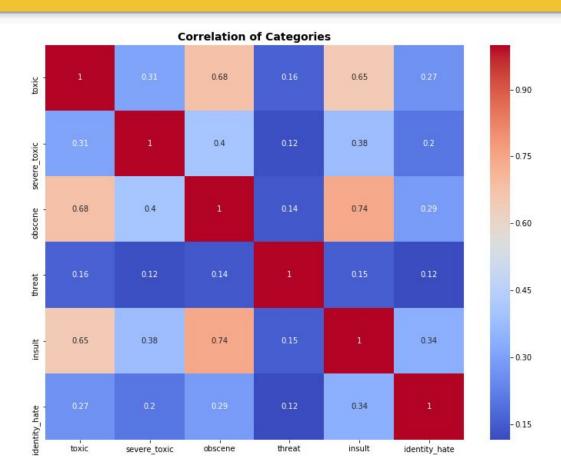


## **Polarity by Category**





## **Category Correlation**





## **Category Samples**

#### **Toxic**



Taboo, rude

#### Obscene



**Profanity** 

#### Insult



Rude, degrading

#### **Severe Toxic**



Extreme language

#### **Identity Hate**



Attacking race, religion and orientation

#### **Threat**



Threats of death

## Topic Modeling by Category: TF-IDF/NMF

TF-IDF/NMF - Frobenius Norm

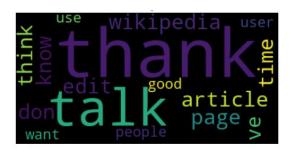
TF-IDF/NMF - Kullback-Leibler Divergence

LDA

Clean







Discussions of article edits









**Toxic** 

Derogatory language

## Topic Modeling by Category: TF-IDF/NMF

LDA



Threats of death

**Identity Hate** 



Attacking race, religion and orientation

## TF-IDF/NMF - Frobenius Norm



TF-IDF/NMF - Kullback-Leibler Divergence







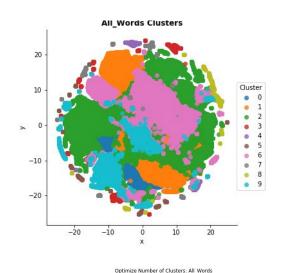
ass

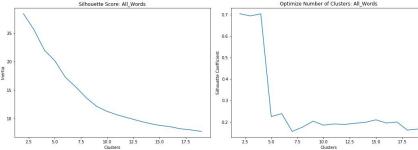
## K Means Clustering: All Words (Toxic + Clean)













## K Means Clustering: Clean Words (All Rated == 0)

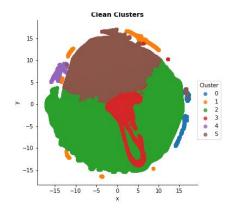


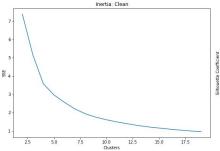


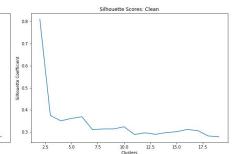








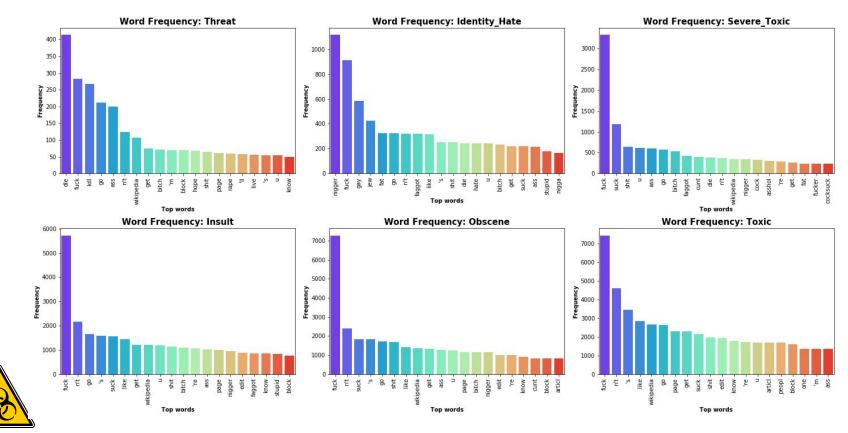




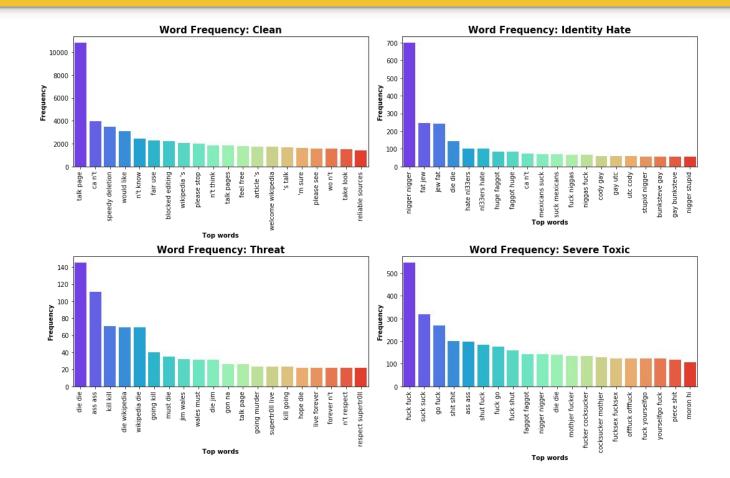


## Word Frequency: Singles

#### Only stopwords filtered



### Word Frequency: Bigrams





## Topic Modeling, 6 Topics from Rating > 0

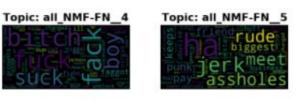
All tagged as at least one form of toxic language 6 categories, topical modeling also roughly ties to these































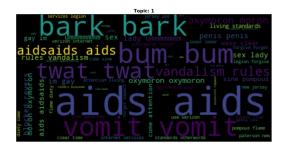




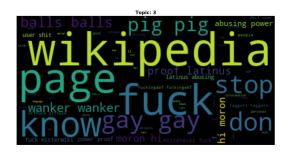


## LDA WordClouds By Toxic Bigram Topic













## LDA WordClouds By Class

Clean



Toxic



Obscene



Insult



Severe Toxic



**Identity Hate** 



**Threat** 





## LDA, 6 Topics







