# Toxic Language Classification:
## *Cleaning Up Wikipedia*

Brian McMahon

12 March 2018

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



VIEWER DISCRETION ADVISED

Dataset from Alphabet-sponsored Kaggle Competition to <u>develop tools to "help improve online 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*

Severity: Low → High

**Toxic.** General taboo, rude language

**Obscene.** Profanity

**Insult.** Rude, degrading

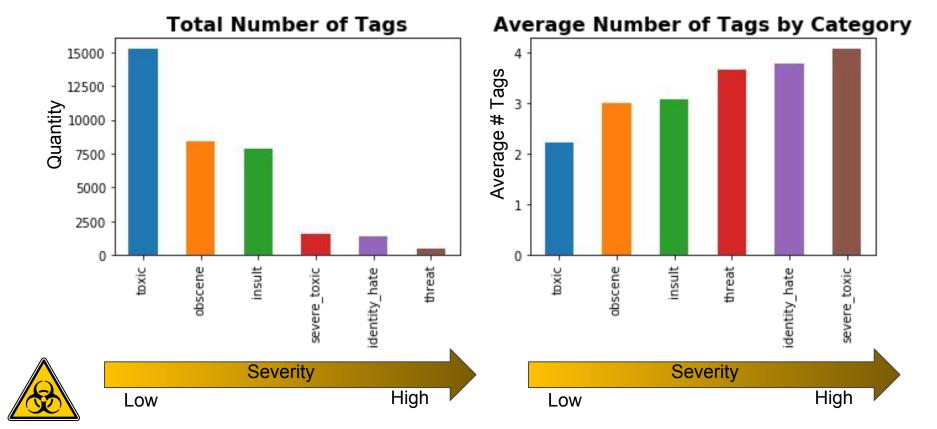**Severe Toxic.** Extreme language, statements not classified as Threat or Identity Hate

**Identity Hate.** Attacking race, religion, orientation
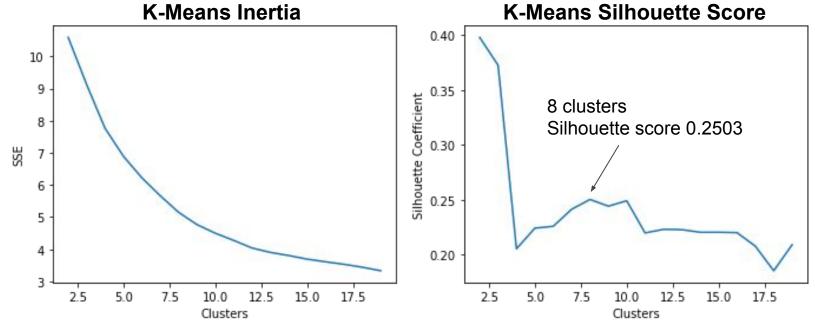
**Threat.** Death wishes, threats

# The Six Categories of Toxic Comments

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

# K Means Clustering: Toxic Words (All Rated > 0)

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



**K-Means Inertia**

**K-Means Silhouette Score**

8 clusters
Silhouette score 0.2503

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



Cluster 0 — **Severe Toxic**

Cluster 1 — **Toxic**

Cluster 2 — **Toxic**

Cluster 3 — **Identity Hate**

Cluster 4 — **Clean**

Cluster 5 — **Identity Hate**

Cluster 6 — **Insult**

Cluster 7 — **Obscene**

# 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 a challenge for machines to fully grasp!

Questions?

# APPENDIX

# 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*

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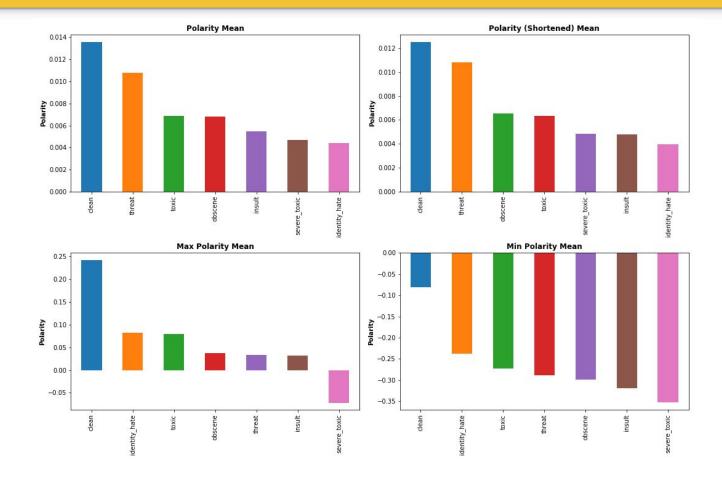
**NMF.** Matrix factorized into multiple matrices with no negative elements

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

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

# Polarity by Category

# 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



**LDA**

**TF-IDF/NMF - Frobenius Norm**

**TF-IDF/NMF - Kullback-Leibler Divergence**

**Threat**

*Threats of death*

**Identity Hate**

*Attacking race, religion and orientation*

**All_Words Cluster 0**

**All_Words Cluster 1**

**All_Words Cluster 2**

**All_Words Cluster 3**

**All_Words Cluster 4**

**All_Words Cluster 5**

**All_Words Cluster 6**

**All_Words Cluster 7**

**All_Words Cluster 8**

**All_Words Cluster 9**

All_Words Clusters

Silhouette Score: All_Words

Optimize Number of Clusters: All_Words

# The Six Categories of Toxic Comments

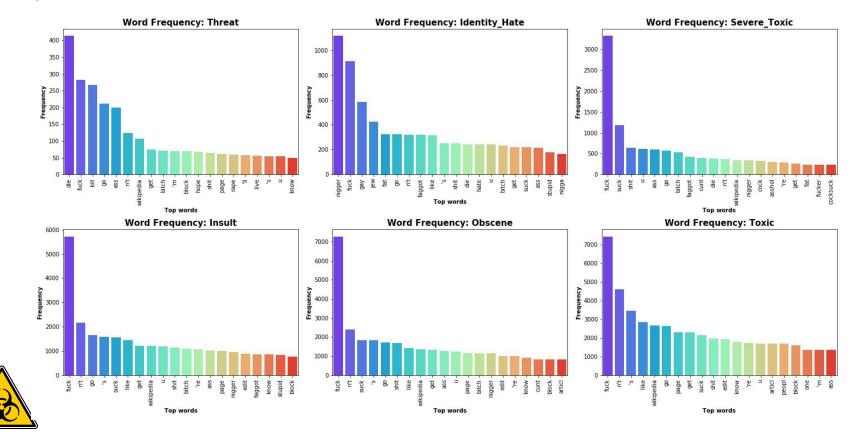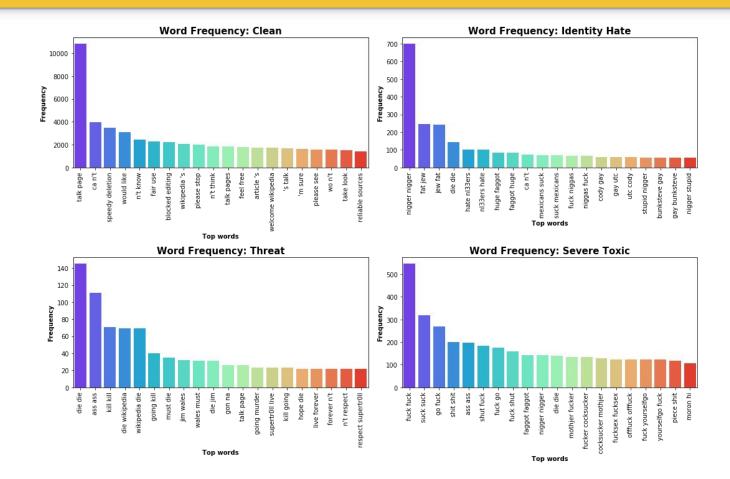*Dataset contains multi-tagged comments; "rating" refers to the number tags a comment has received*



*Rating means the number of "tags" per comment*

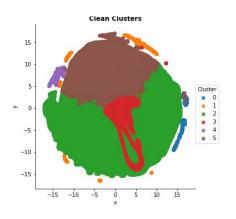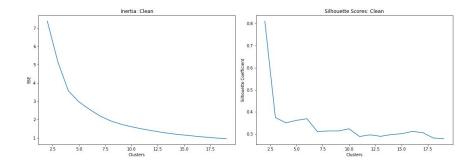Correlation of Categories

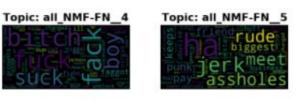# Word Frequency: Singles

*Only stopwords filtered*

# Word Frequency: Bigrams

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

# LDA WordClouds By Class



Clean

Toxic

Obscene

Insult

Severe Toxic

Identity Hate

Threat