



Toxic Comments Classification

Bryan Goh

Why do people leave toxic comments?

1. Anonymity
2. Dehumanization
3. Lack of Real-Time Feedback
4. Mob Mentality

About 1 in 5 victims of online harassment say it happened in the comments section

BY MONICA ANDERSON

The Reuters news service [recently joined other media outlets](#) that have closed or restructured their online comments. [Popular Science](#) made headlines last year when it shut down comments. The Huffington Post ended [anonymous commenting](#) in late 2013, with its editor calling comments sections one of the “darkest places on the internet.” And the Chicago Sun-Times [suspended](#) comments in April until a new monitoring system could be developed.

A recent [Pew Research Center study](#) found that roughly one-in-five (22%) internet users that have been victims of online harassment reported that their last experience occurred in the comments section of a website. While social media sites (66%) were the most common place noted for harassment, comments sections were named more frequently than online gaming sites (16%) and discussion sites like reddit (10%).

Where People Witness Online Harassment

Among internet users who have witnessed online harassment, the % that mentioned each online environment to the open-ended question, “Can you describe what you have witnessed in your experiences of this kind?”

	% of total responses
Social Networking Site	15
General Online Comments/Comment Section	8
News Website/Blog/Article	8
Message Board/Online Forum/Chat Room	4
General Internet	4
General Website	2
Games Website	1
Email	1
Text Message	.1
Did Not Mention a Specific Online Environment	40%
Did Not Respond/No Answer	21%

Source: American Trends Panel (wave 4) Survey conducted May 30-June 30, 2014 N=1958.

Note: Responses do not add up to 100% because respondents could reference multiple online spaces.

PEW RESEARCH CENTER

Problem Statement

- Social media apps like Meta(Facebook), TikTok, etc, have risen to high level of popularity.
- Some depend on social media for their daily livelihood.
- More toxic comments → people may stop using that social media app → Affecting the platform's business.

		id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
44574	7715a8d2ca3fa10c		Screw You \n\nWhy don't you wikipedia jerks just go drop some bombs in the toilet? I don't care! mln	1	0	0	0	0	0

		id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
527	01625cc84c6ef15b		How do you know he is dead. Its just his plane that crashed. Jeezz, quit busting his nuts, folks.	0	0	1	0	0	0

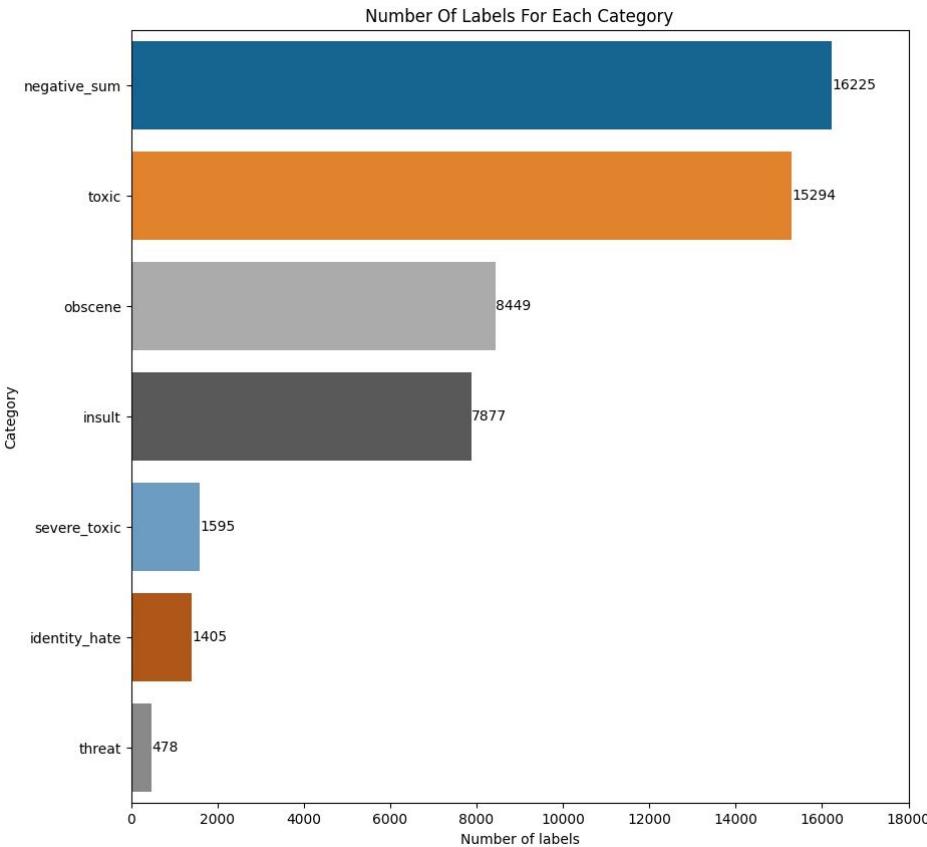
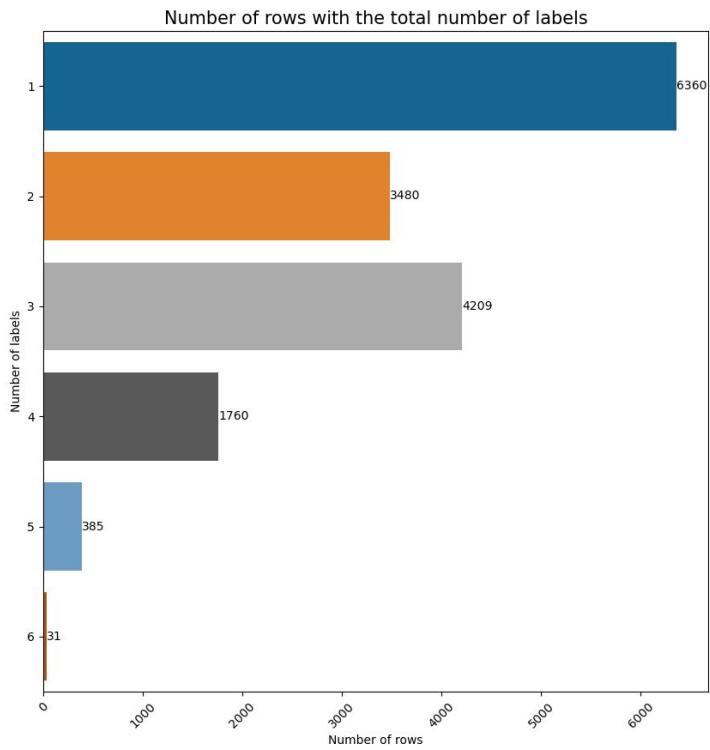
		id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
122729	908b347bfd7c9bb6		I'ma smack ya upside da head wit a shovel. \n\nI'm takin ya down, boi.	0	0	0	1	0	0

		id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
38357	666cac78892fd024		September 2014 (UTC)\nYes, your opinion. Quit acting like an aspie retard. 00:13, 14	0	0	0	0	1	0

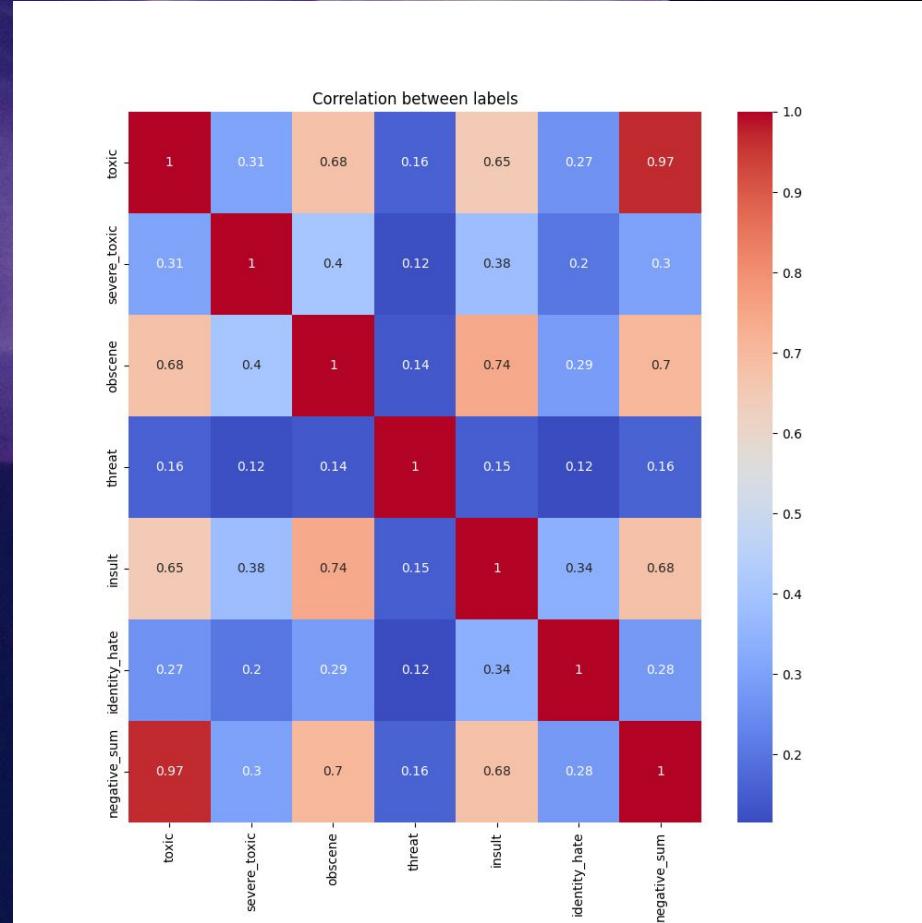
		id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
116799	7057dc91a3f1c800		"\n\n The LGBT Barnstar awarded because your sir/mam are a queer douche -"	0	0	0	0	0	1

EDA

Positive - 90%
Negative - 10%

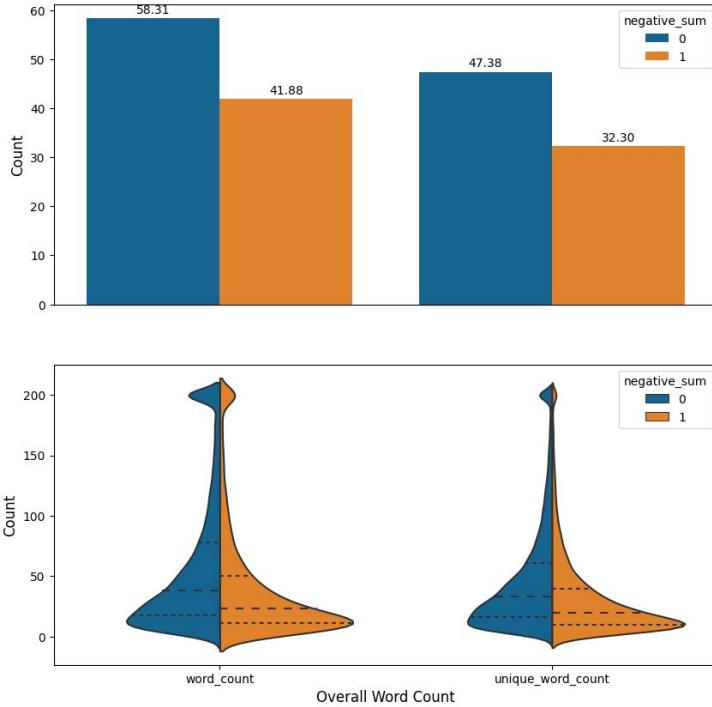


EDA

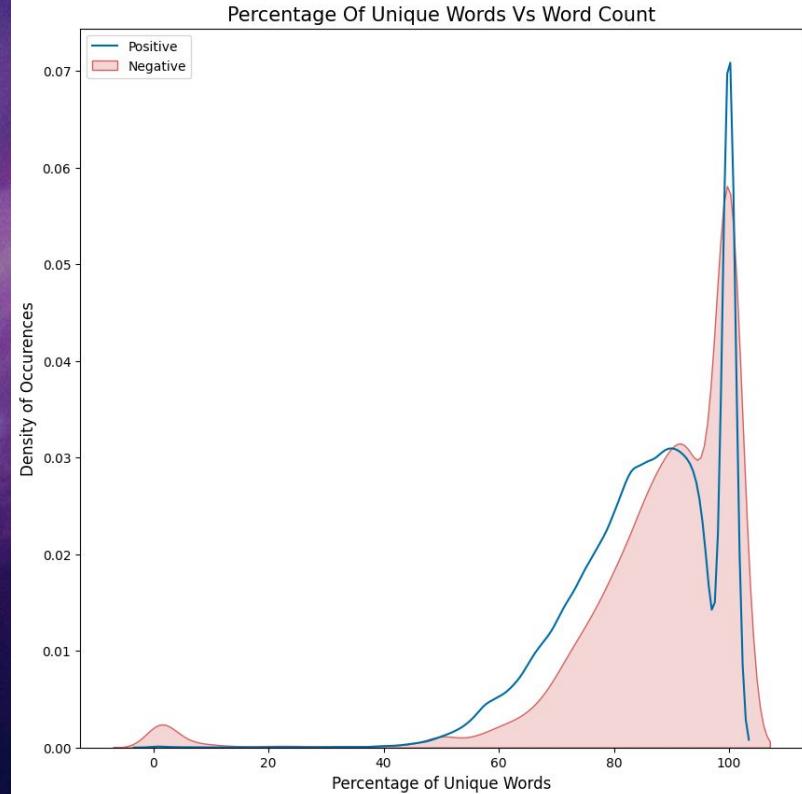


EDA

Average Word Count Vs Average Unique Word Count

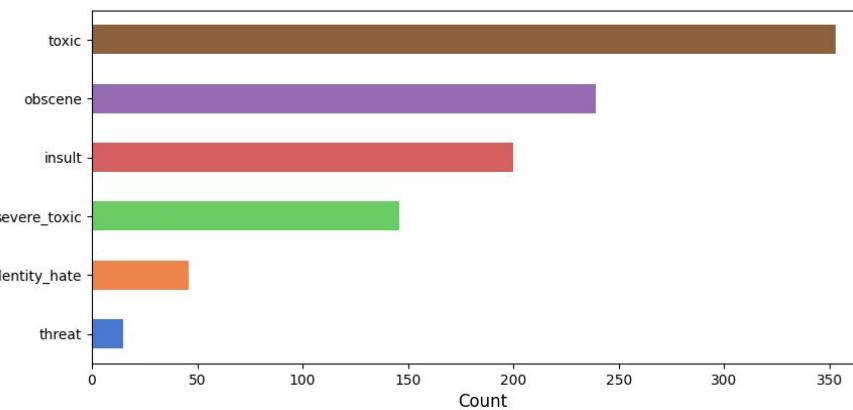
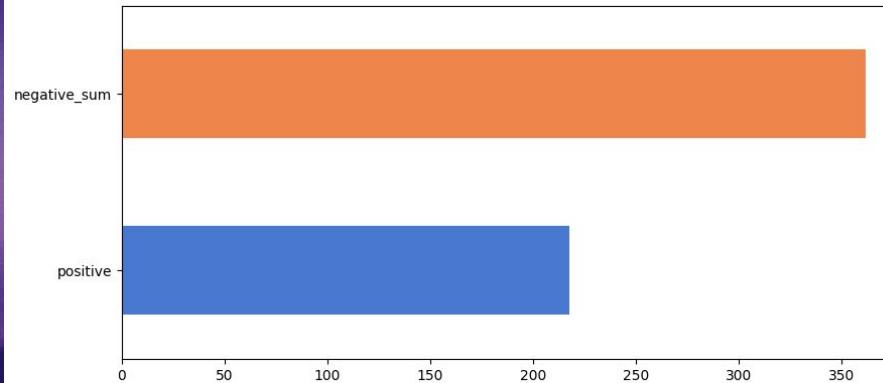


Percentage Of Unique Words Vs Word Count

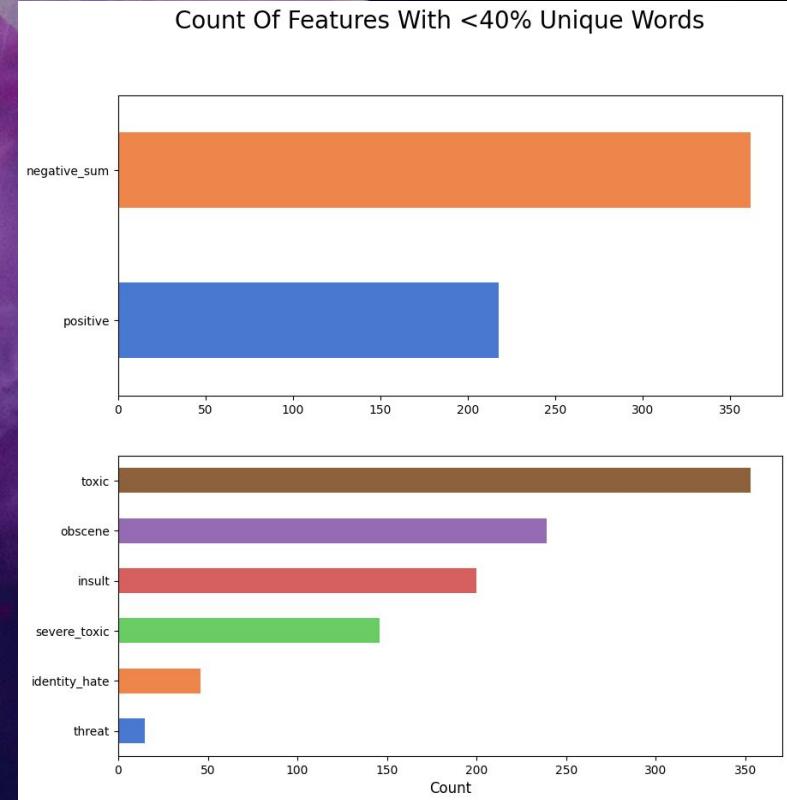
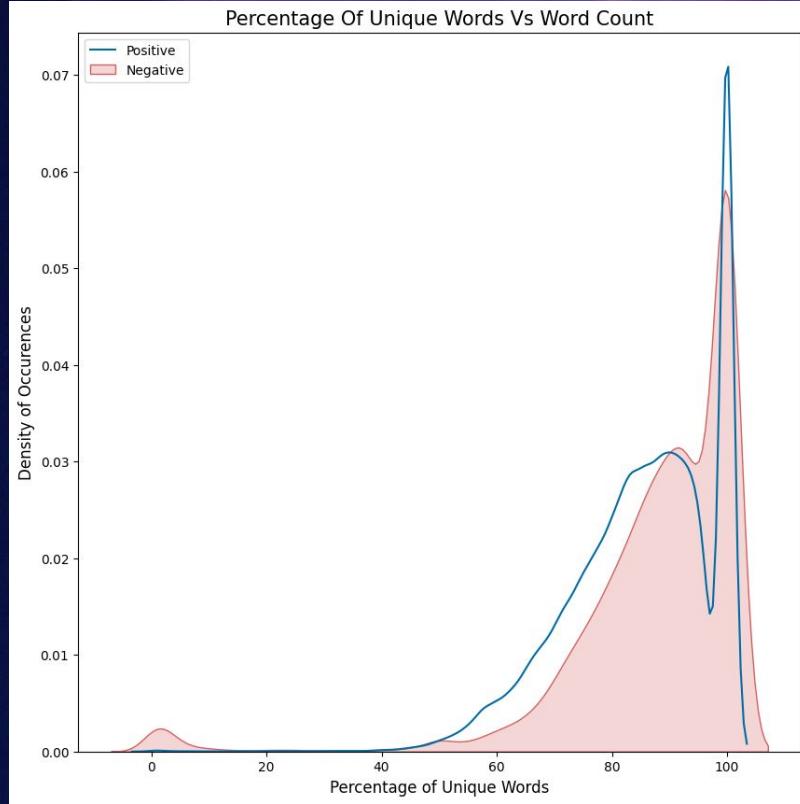


EDA

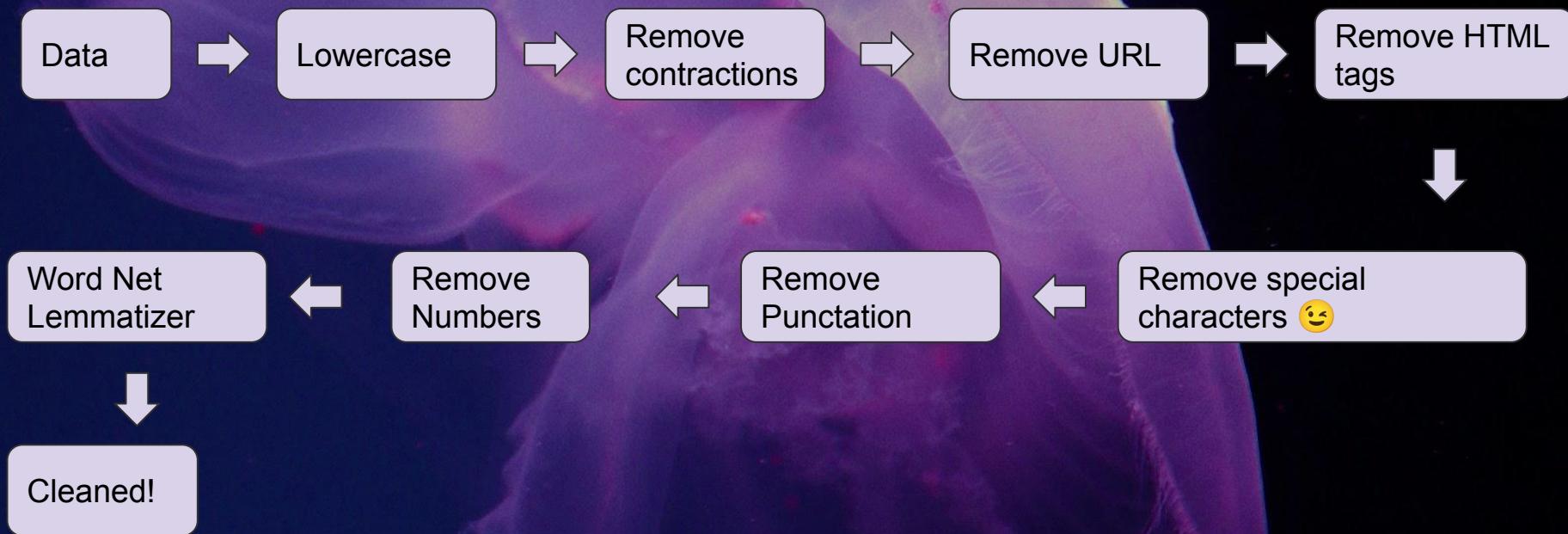
Count Of Features With <40% Unique Words



EDA



Data Cleaning



Word Clouds

Word Cloud for Negative Comments



Word Cloud for Positive Comments



Word Clouds

Word Cloud for Toxic Comments



Word Cloud for Insult Comments



Word Cloud for Obscene Comments



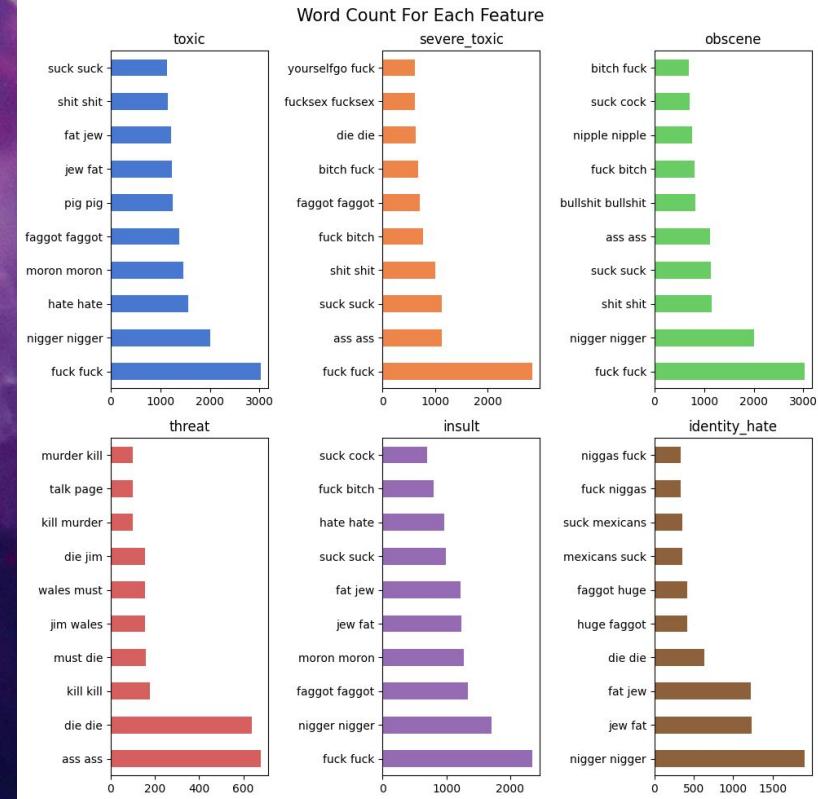
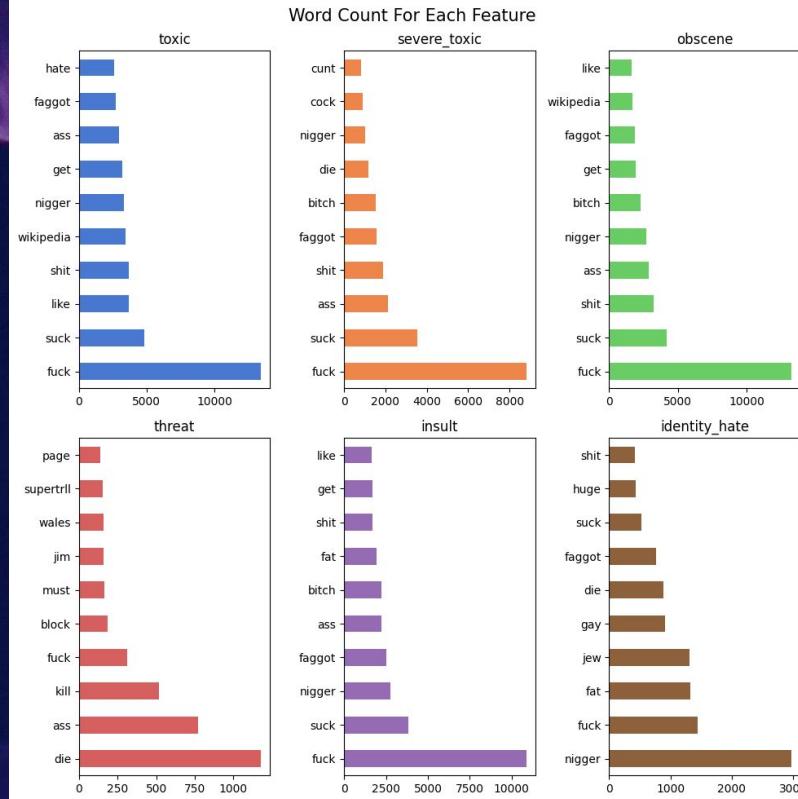
Word Cloud for Severe_Toxic Comments



Word Cloud for Identity_Hate Comments

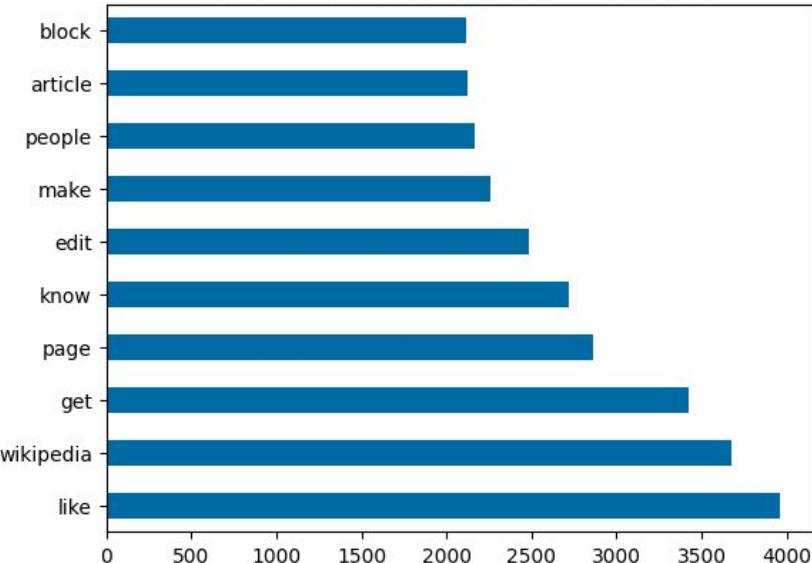


EDA

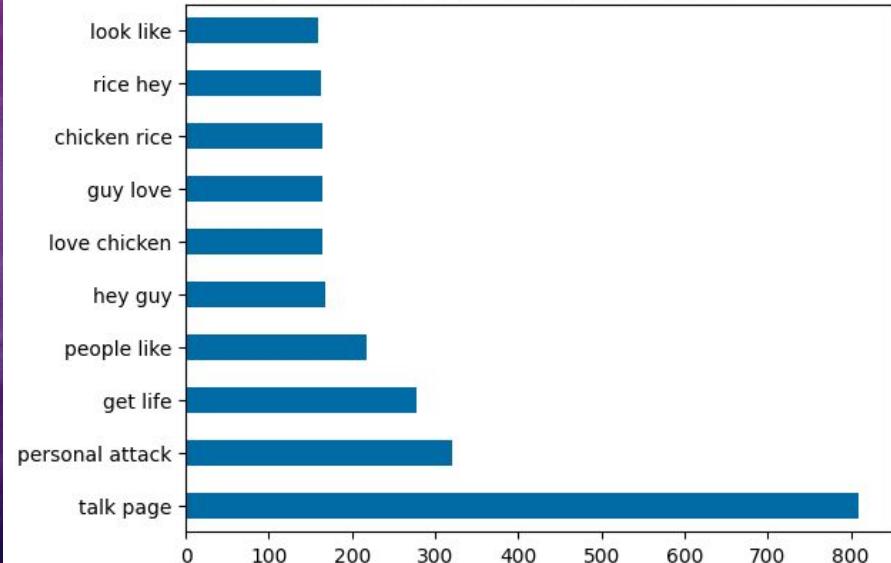


EDA

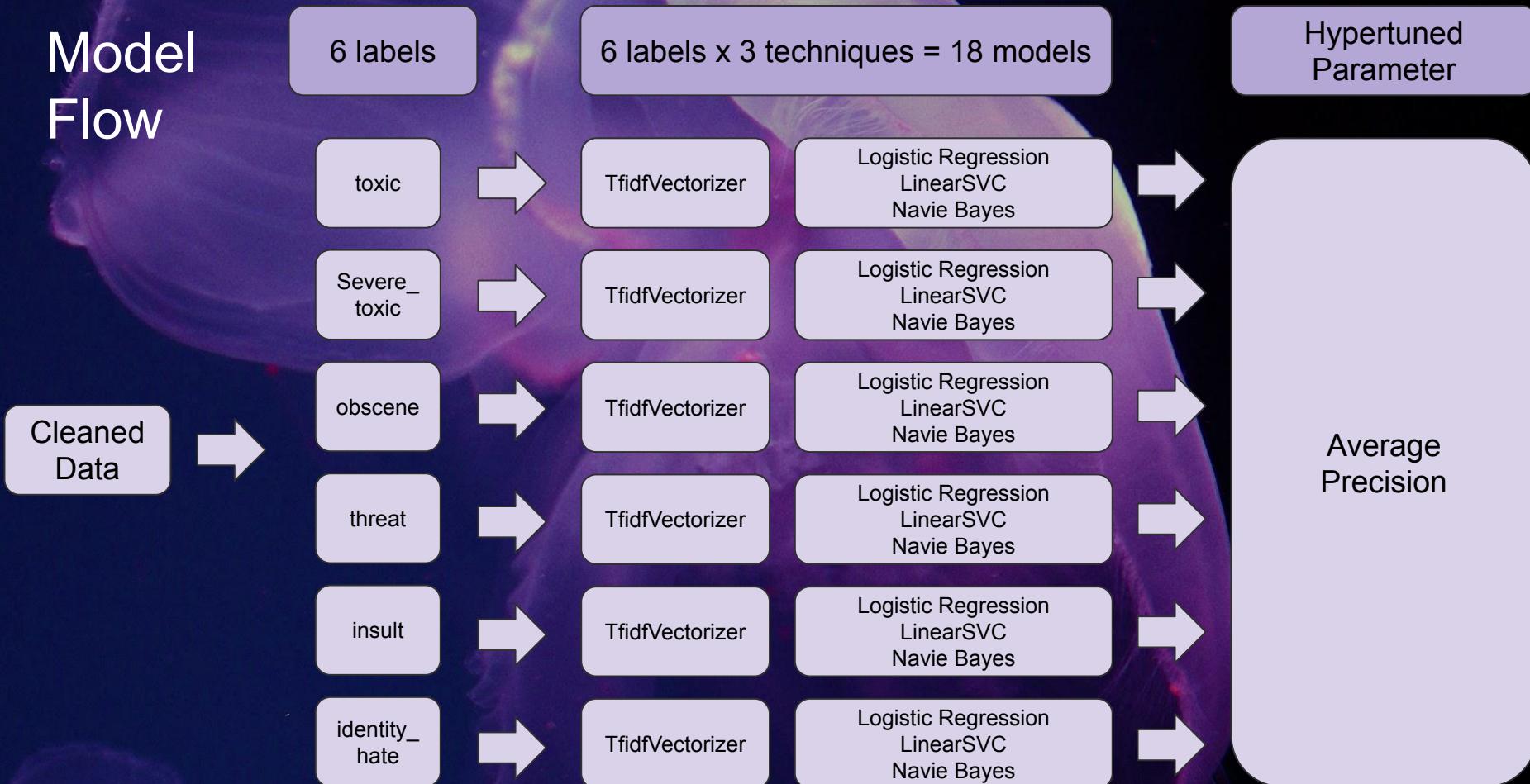
Word Count For Positive Comments



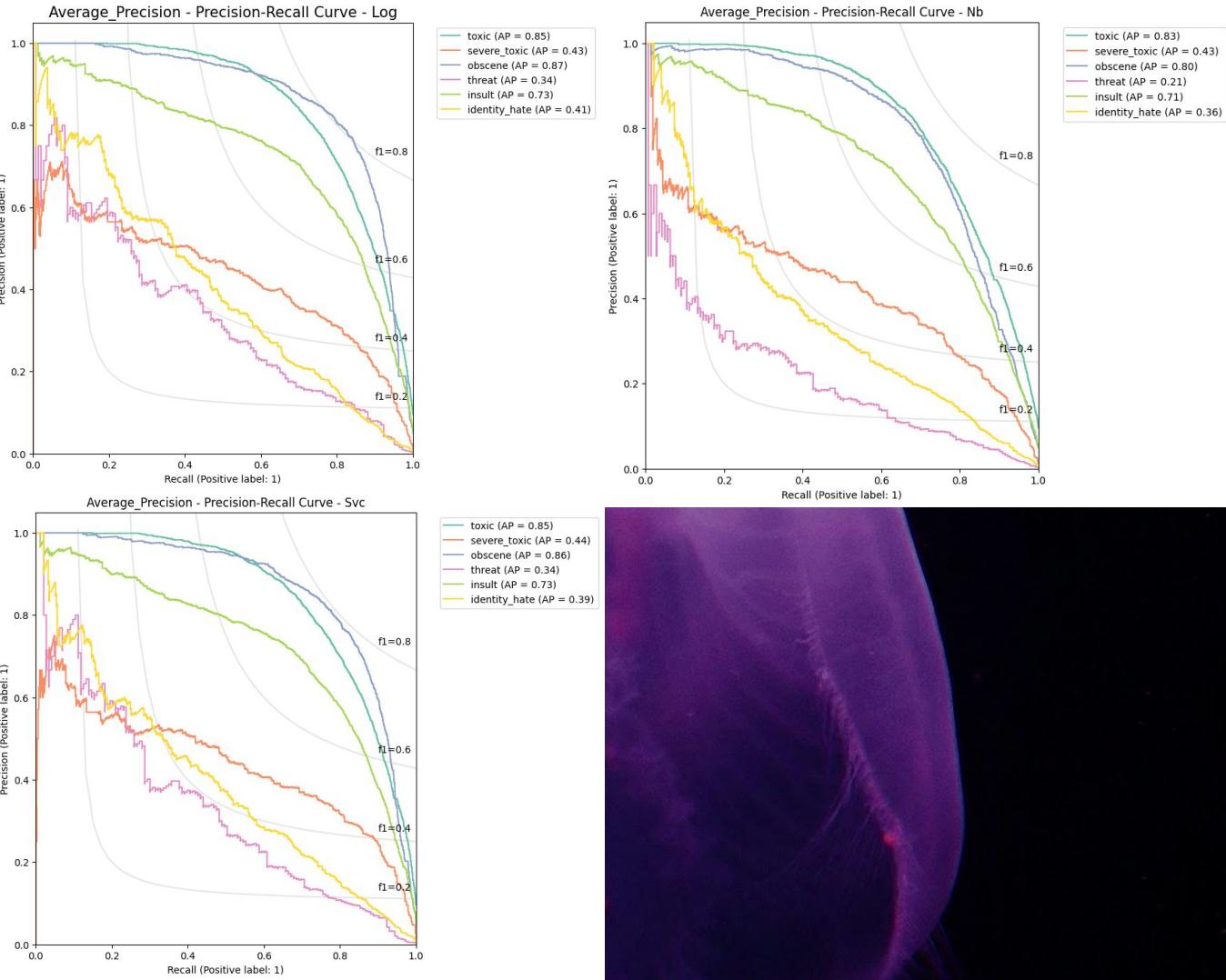
Word Count For Positive Comments - Bigrams



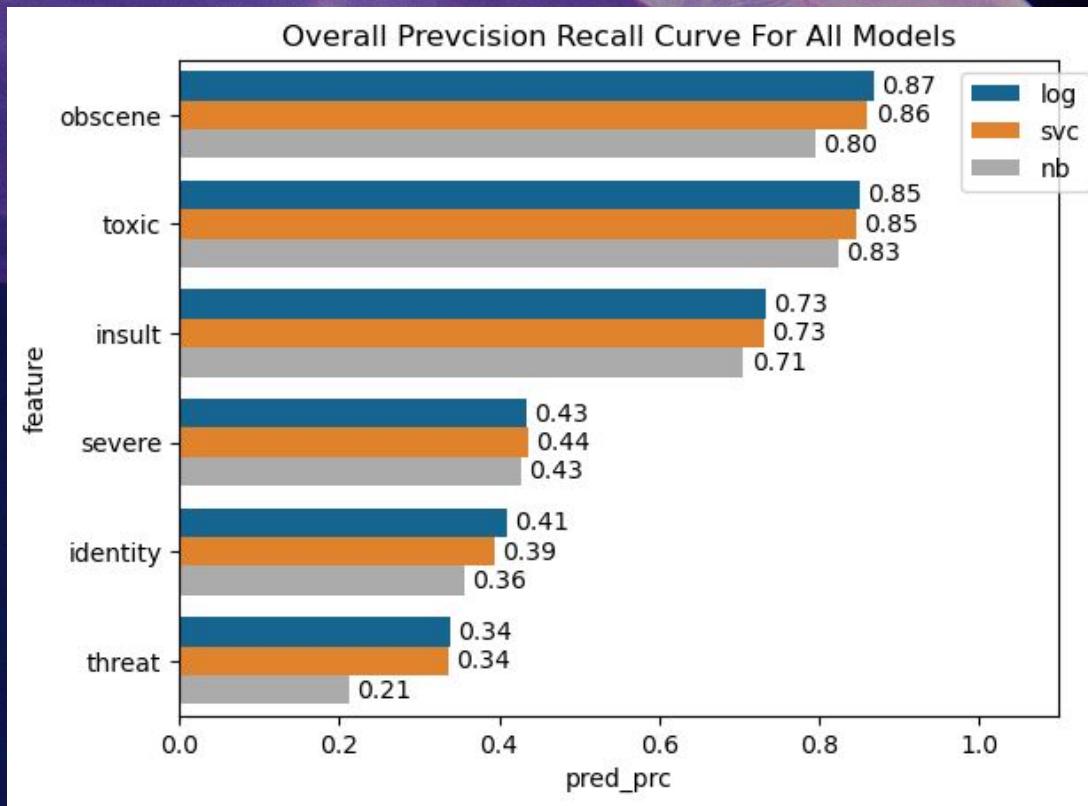
Model Flow



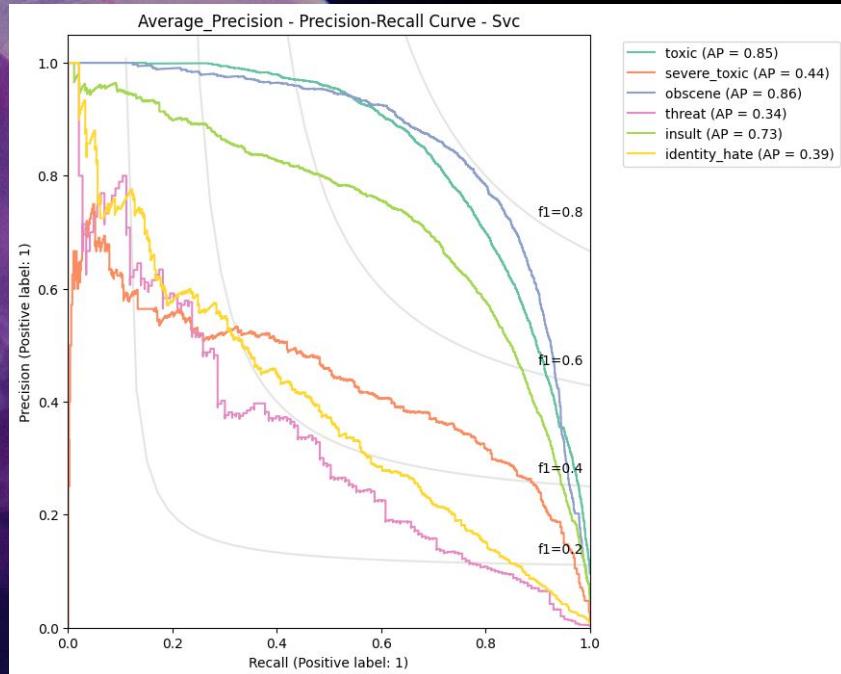
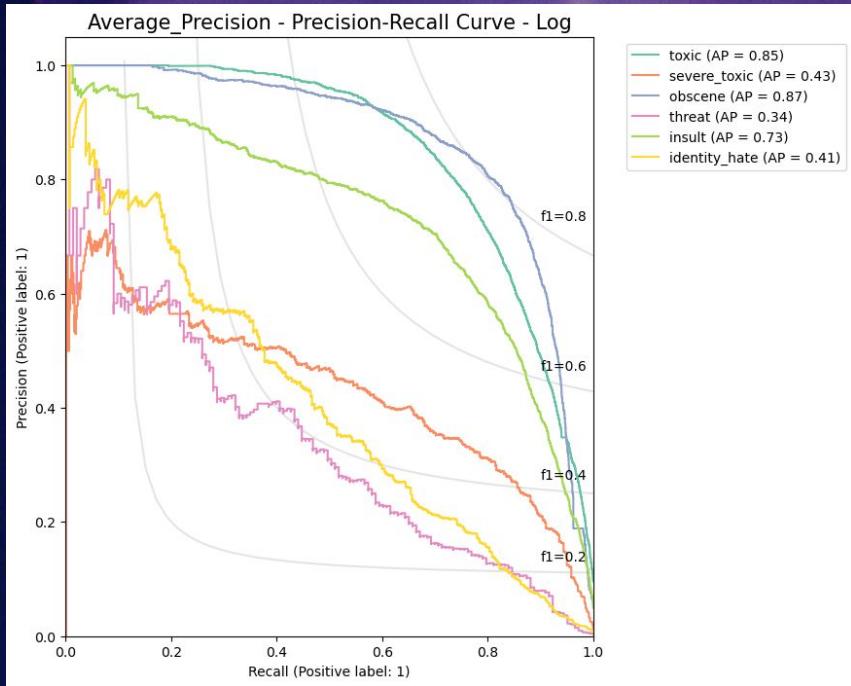
Evaluation



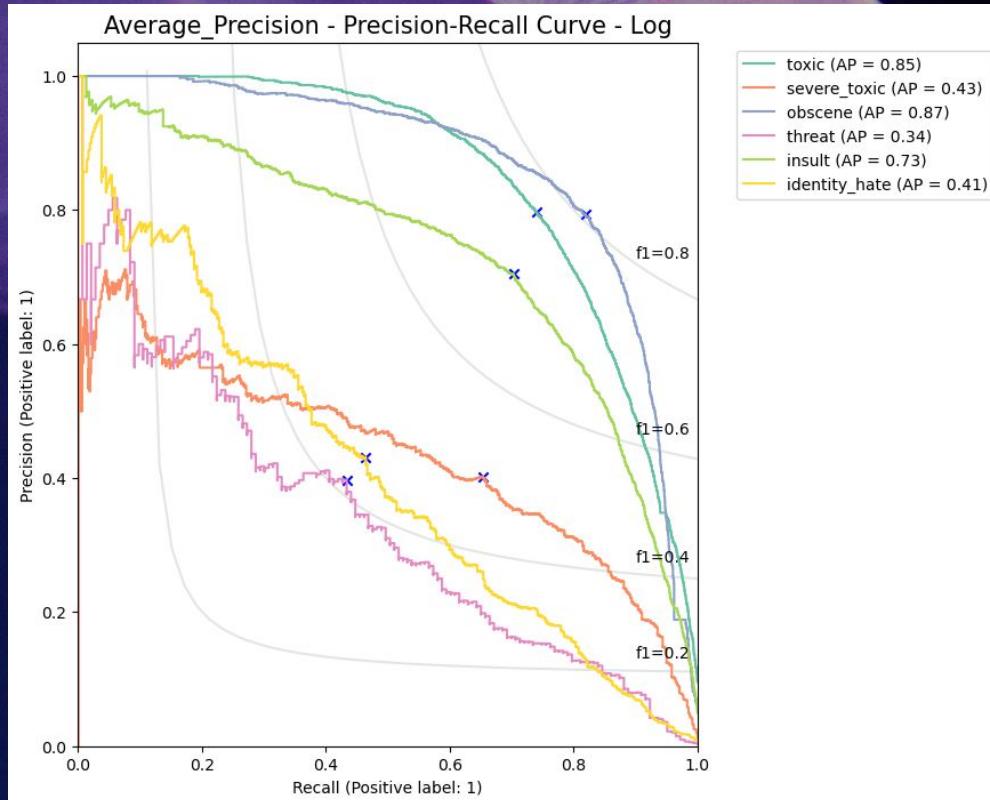
Evaluation



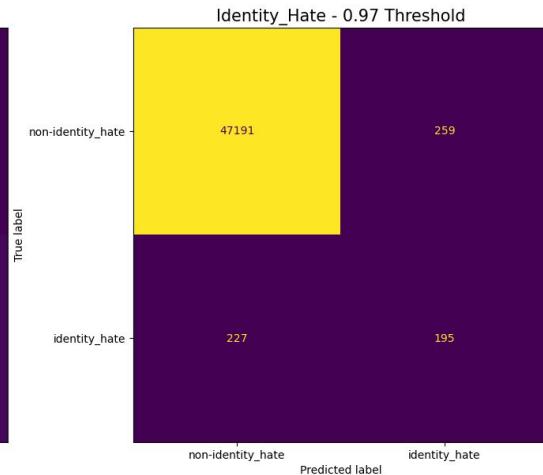
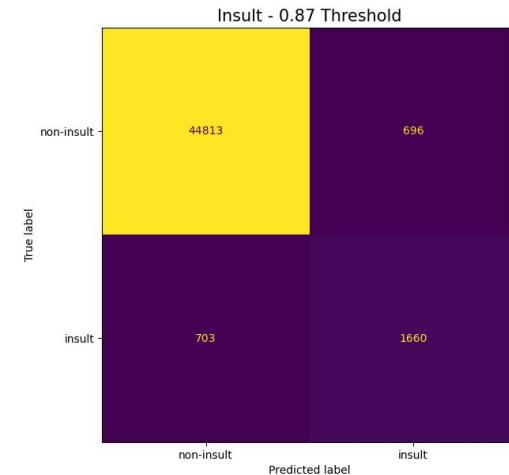
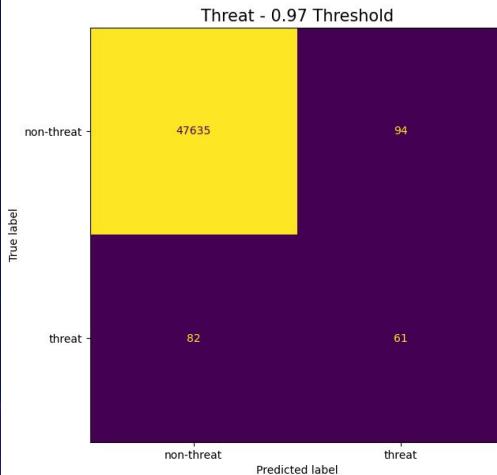
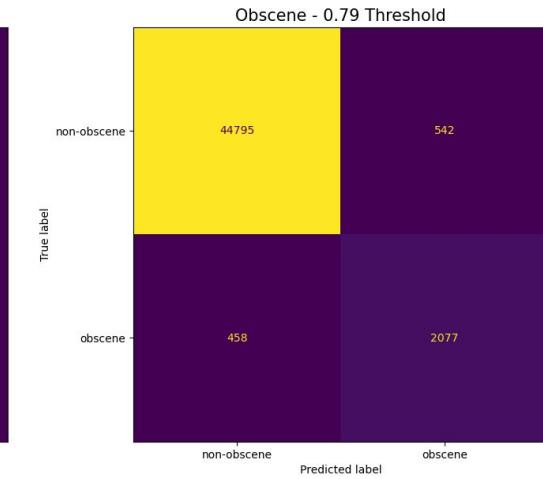
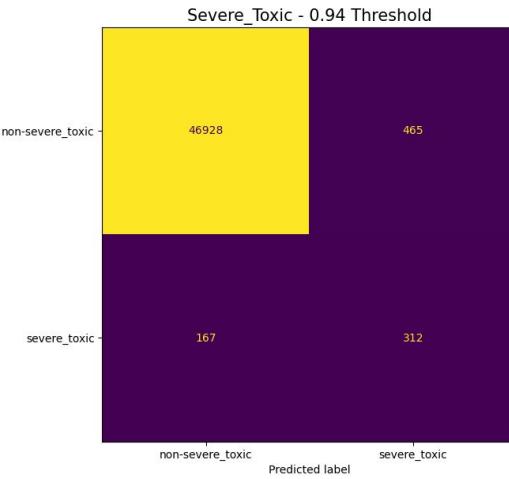
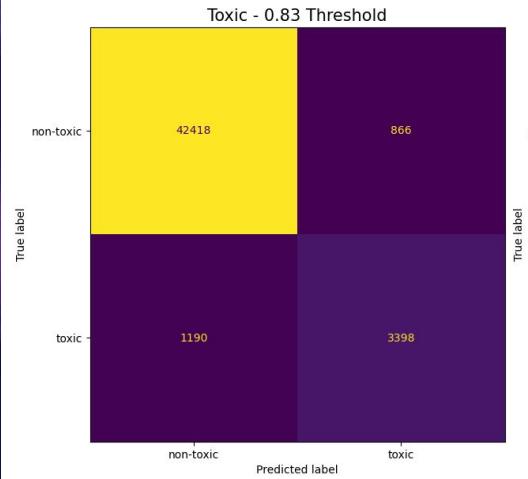
Evaluation



Evaluation



Confusion Matrix - Threshold Tuning

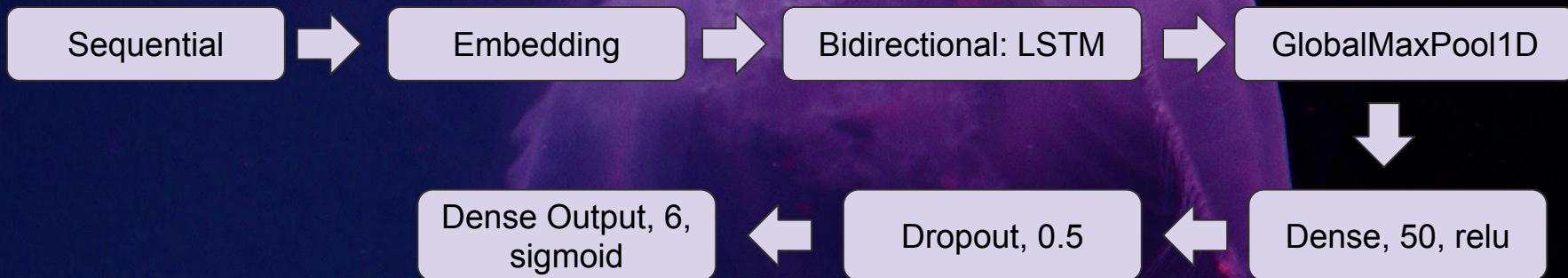


Neural Network

Data Pre-processing



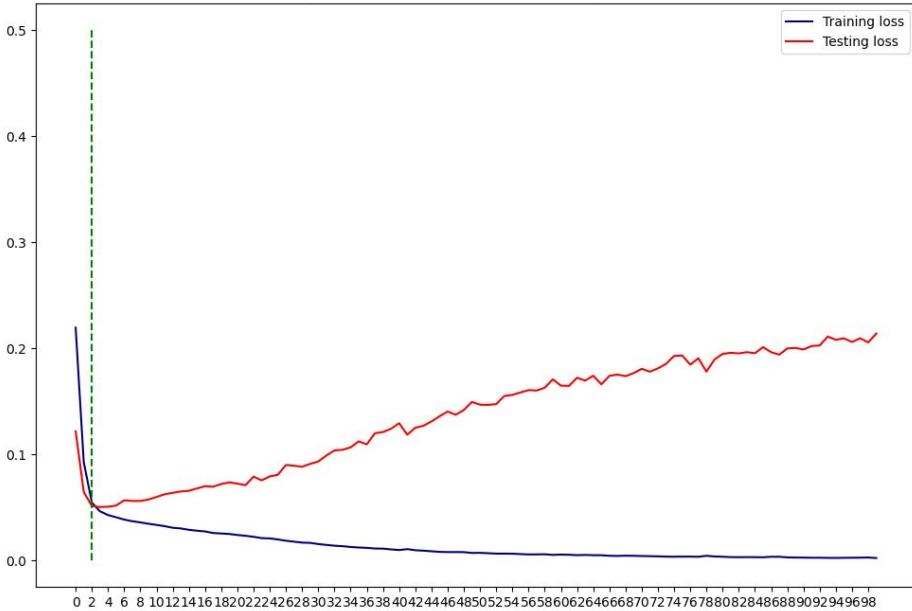
Neural Network Model



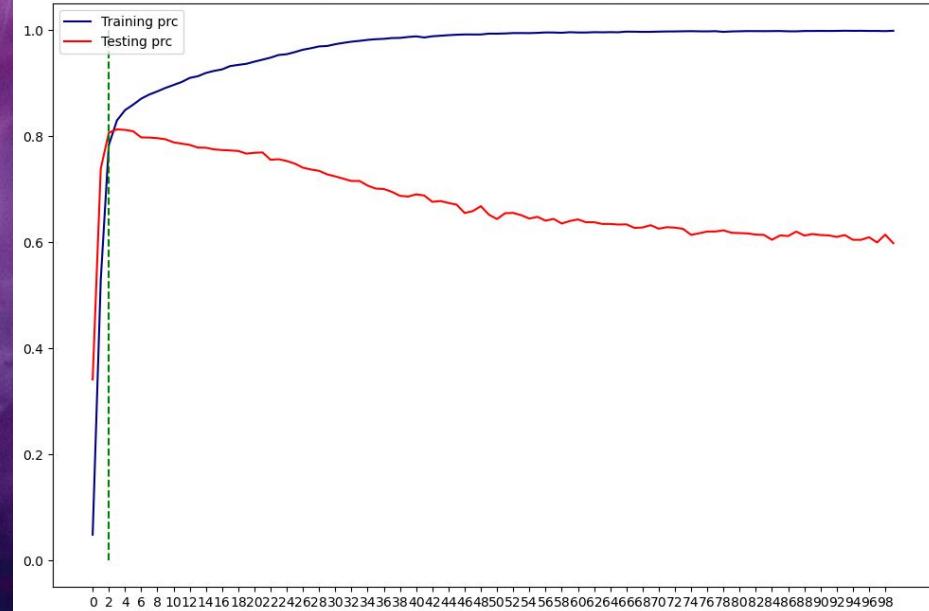
Evaluation

Batch Size = 1024

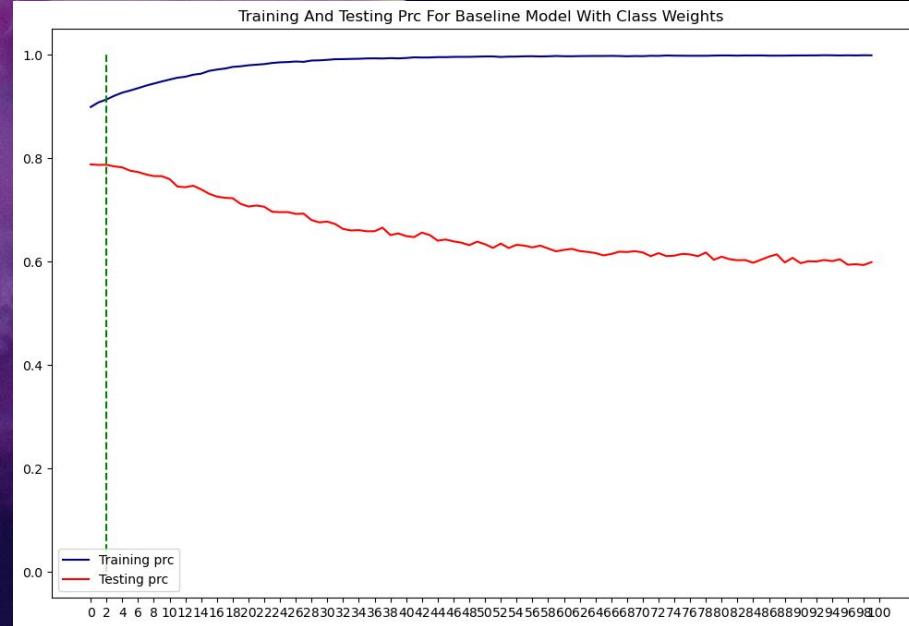
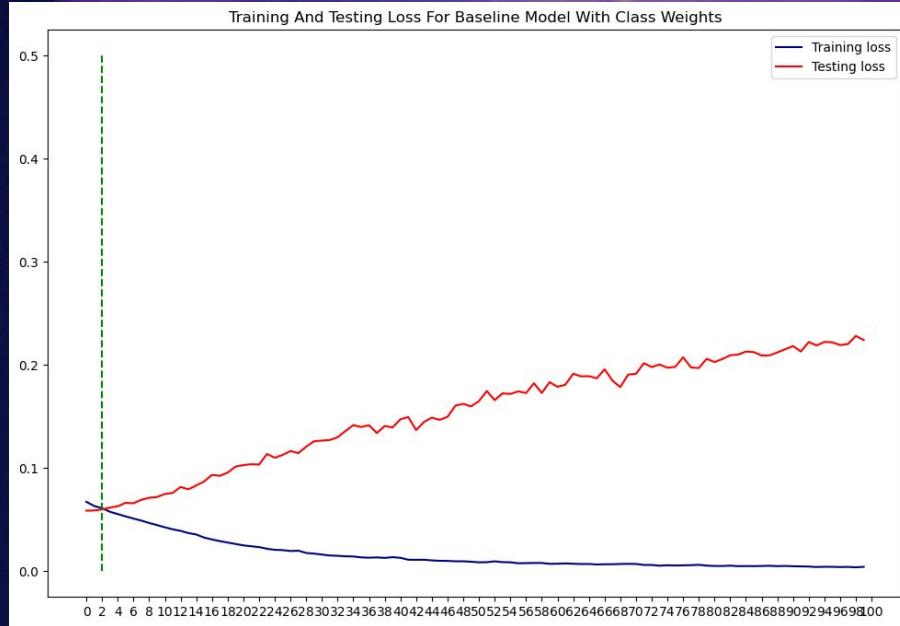
Training And Testing Loss For Model With Large Batch Size



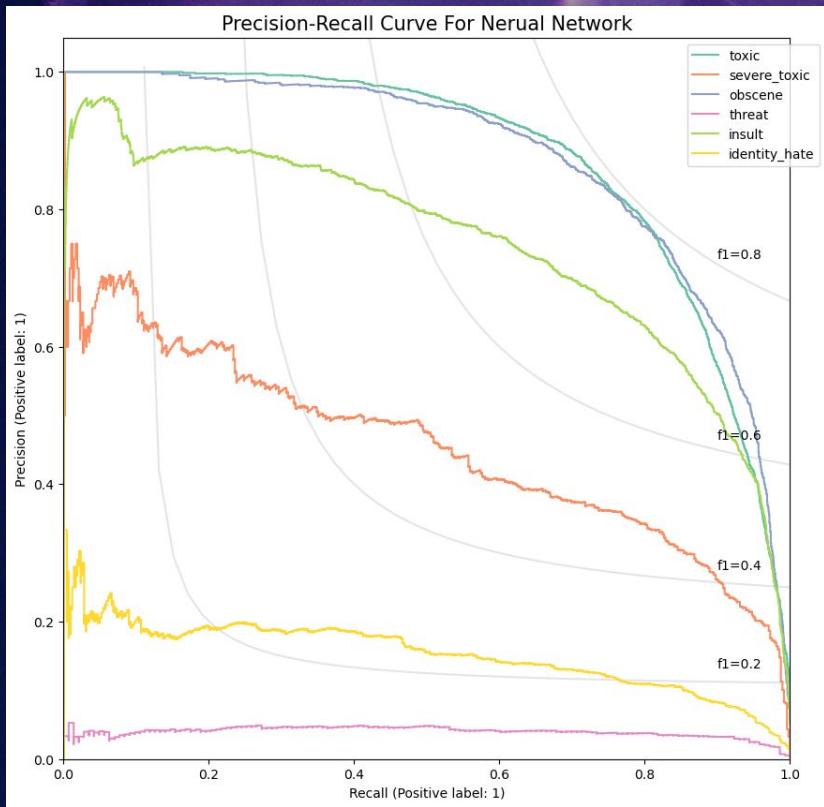
Training And Testing Prc For Model With Large Batch Size



Evaluation With Class Weights

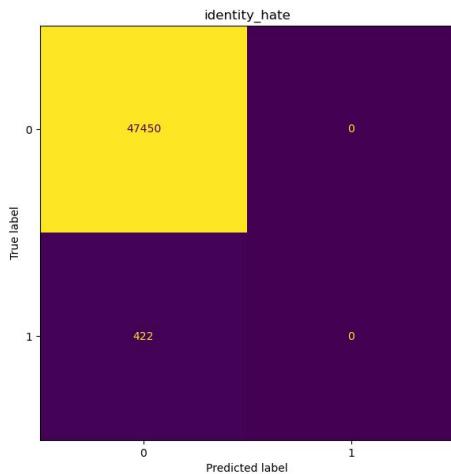
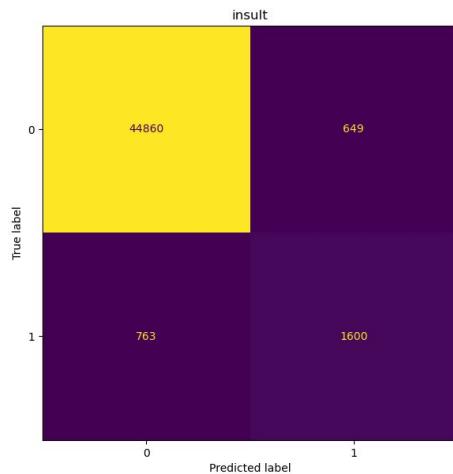
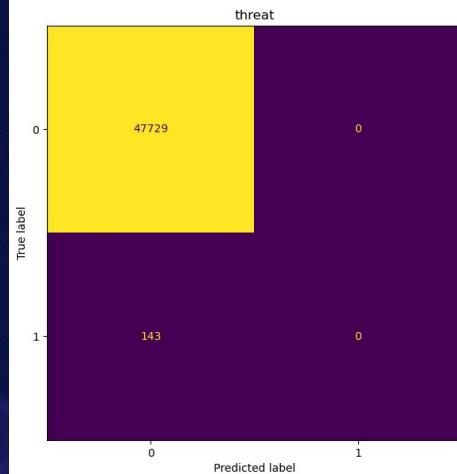
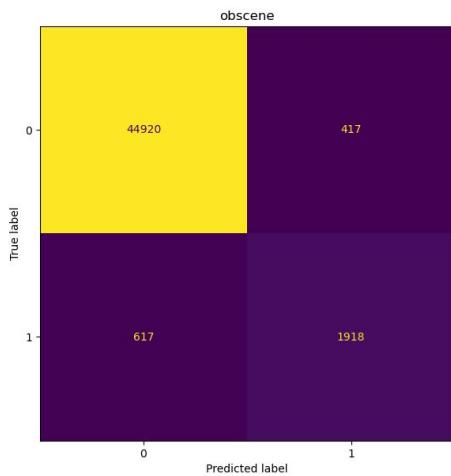
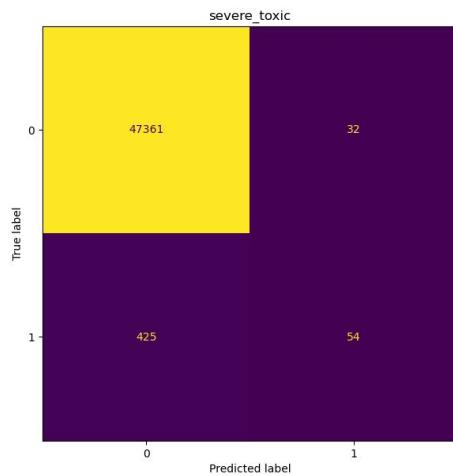
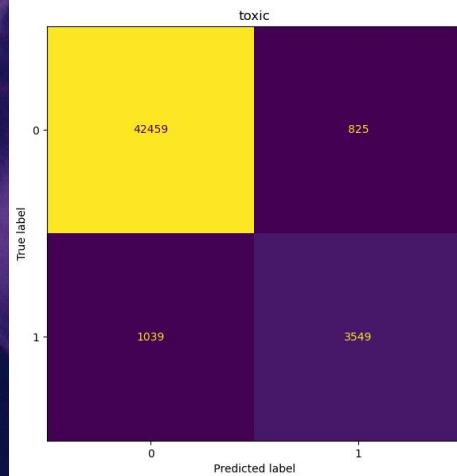


Evaluation



	precision	recall	f1-score	support
toxic	0.81	0.77	0.79	4588
severe_toxic	0.63	0.11	0.19	479
obscene	0.82	0.76	0.79	2535
threat	0.00	0.00	0.00	143
insult	0.71	0.68	0.69	2363
identity_hate	0.00	0.00	0.00	422
micro avg	0.79	0.68	0.73	10530
macro avg	0.50	0.39	0.41	10530
weighted avg	0.74	0.68	0.70	10530
samples avg	0.07	0.06	0.06	10530

Overall Confusion Matrices



A close-up photograph of a glowing jellyfish against a dark background. The jellyfish has a translucent, glowing body with a bright yellow/orange center and a darker purple/pink outer edge. Its long, flowing tentacles are visible, with some small, glowing red spots along them.

Deployment

Comparisons

Method	Split Modeling	Combined Modeling
Model	Logistic Regression	Neural Network
Type	Binary Output	Multi-Label Output
Metrics	Higher F1-score	Lower F1-score
Run Time	Longer run time (~6 hours)	Short run time (~1h for 100 epoch) (Usage of GPU)
	Model only looks at the count of words	Model is able to learn from the word position
	Model cannot be improved further	Model can be improved by transfer learning (Word2Vec/Bert)

Conclusions:

- Split modeling shows that binary classification is more effective than multi-label modeling.
- Split modeling has a longer run time than combined modeling.
- Even logistic regression gives the best Precision-Recall ratio, due to the nature of neural networks, being able to do a bi-directional layer, shows that neural network is the better model.

Recommendations:

1. Usage of pre-trained models, (Word2Vec, Bert) to fit to the data.
2. Changing the dataset to a binary or multi-class problem.
3. Usage of larger datasets, especially for each feature.
4. The model is unable to be used now, as the data is from 2015.

Limitations

1. The data is very dependent on human intuition → Cannot expect machine to predict better than us
2. Very little data in certain feature labels

Questions?

