

Using NLP to Predict Posts from Reddit

Dennis Tran



Problem Statement

A Multimedia Company who owns several Youtube channels specializing on reposting narrated Reddit posts seeks to automate its post collection.

Of interest are two subreddits, r/NoSleep and r/IDontWorkHereLady, known for user submitted stories.

Our team's task is to create a machine learning model that can determine where a post came from.

The Process

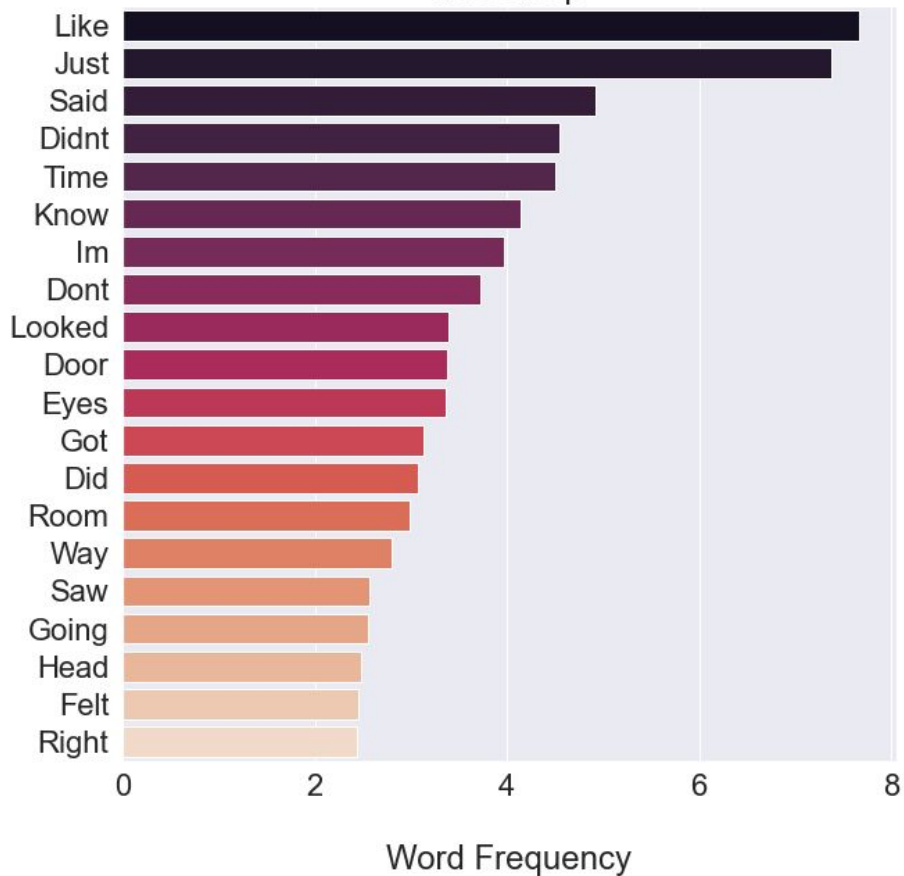
We scraped our data with Pushshift

Cleaned as we scraped and after

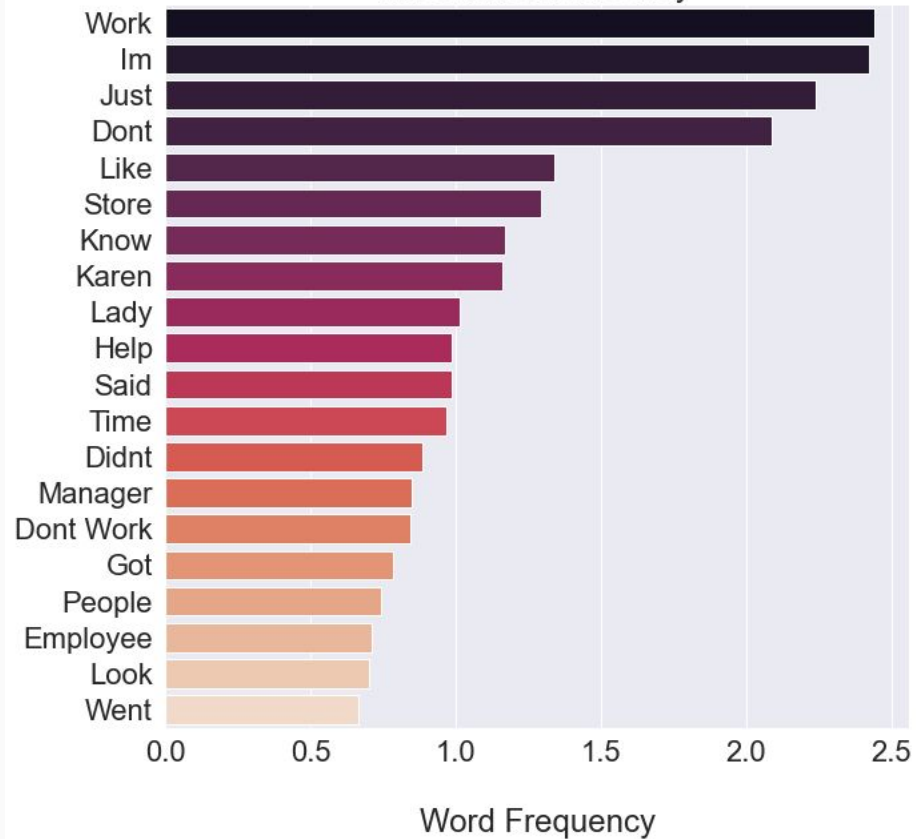
Tokenized then lemmatized

Modeled

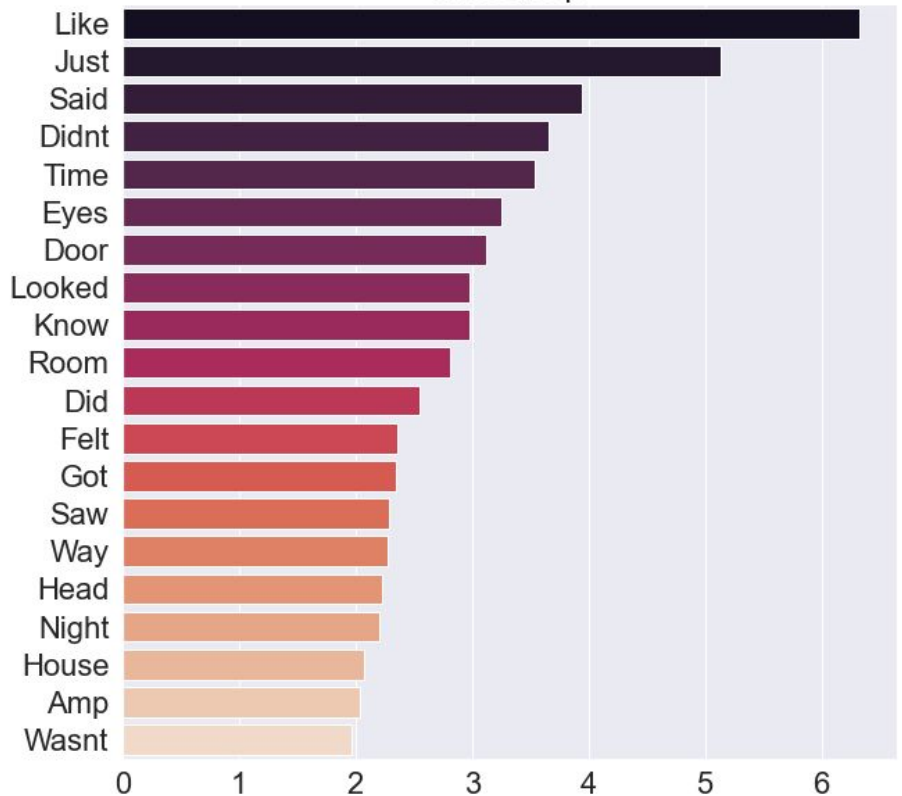
r/NoSleep



r/IDontWorkHereLady

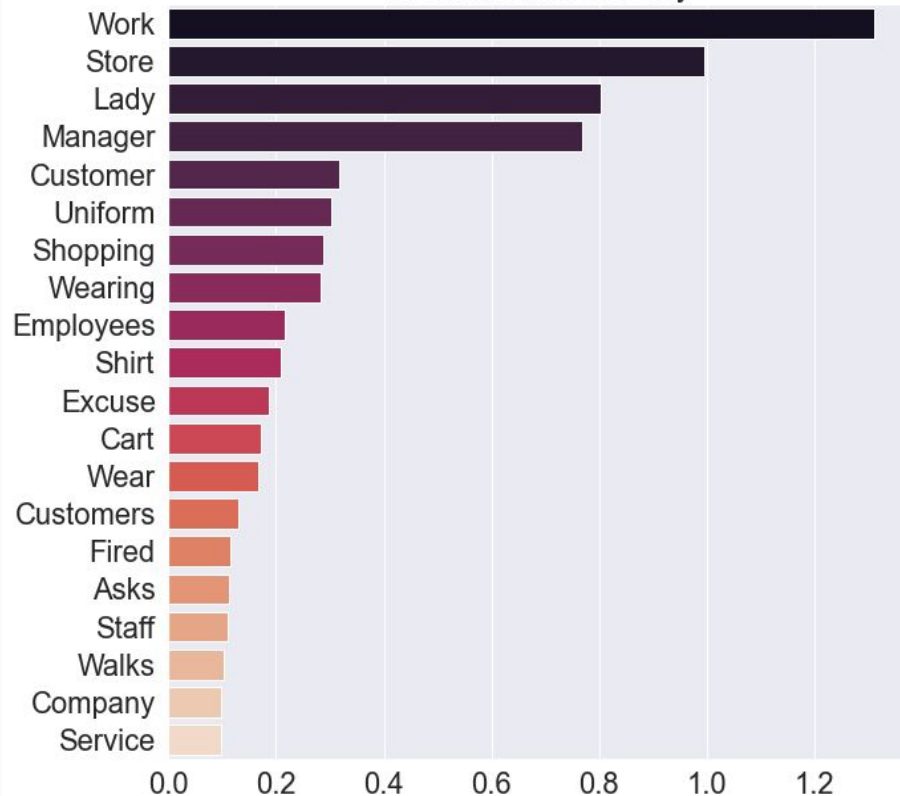


r/NoSleep



Word Frequency Compared to r/IDontWorkHereLady

r/IDontWorkHereLady



Word Frequency Compared to r/NoSleep

Modeling

Tested for accuracy

Baseline 50%

Used Ridge and Lasso Regression, Multinomial Naive Bayes, and Random Forest

All models except Naive Bayes overfit

However, we expected much more severe overfitting than what resulted

	Training Score	Testing Score
Logistic_Regression	1.000000	0.976471
Rand_Naive_Bayes	0.993464	0.988235
Grid_Naive_Bayes	0.993464	0.988235
Lasso	1.000000	0.976471
Modified_LogReg	1.000000	0.968627
TFIDF_Random_Forest	1.000000	0.988235
CVec_Random_Forest	1.000000	0.956863

Conclusion

The Multinomial Naive Bayes performed the best, with an accuracy score of 98.8%.

May be in part due to the different nature of the subreddits despite both containing written stories.

Further testing and tuning could be done on the Bayes model, but we must consider if the cost is worth it.