

# REDDIT NLP

*Kade Higgins*

Web Scraping & NLP  
Text Classification

A black square and a yellow square are located in the top right corner of the slide.

# PROBLEM STATEMENT

Did this subreddit text come from `r/C_S_T` or  
`r/conspiracy`?



# LET'S PLAY A GAME!

*"Am I missing something else going on?"*

*"Did time speed up after 9/11?"*

*"Is Our Whole Government Corrupt Or Is It Certain Networks?"*

*"Washington DC was staged just like Gorge Floyd"*



# METHODOLOGY

*Scrape r/C\_S\_T & r/conspiracy reddit threads*

*Clean (remove HTML, Lemmatize)*

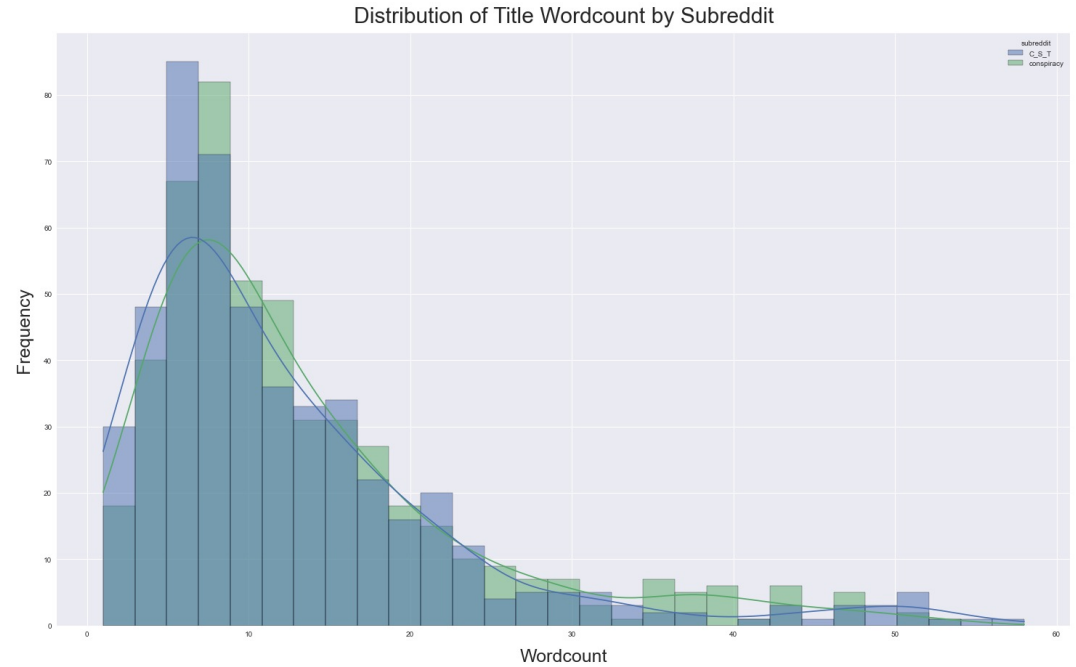
*EDA & Visualization*

*Model - Naïve Bayes*

*Model – Logistic Regression*

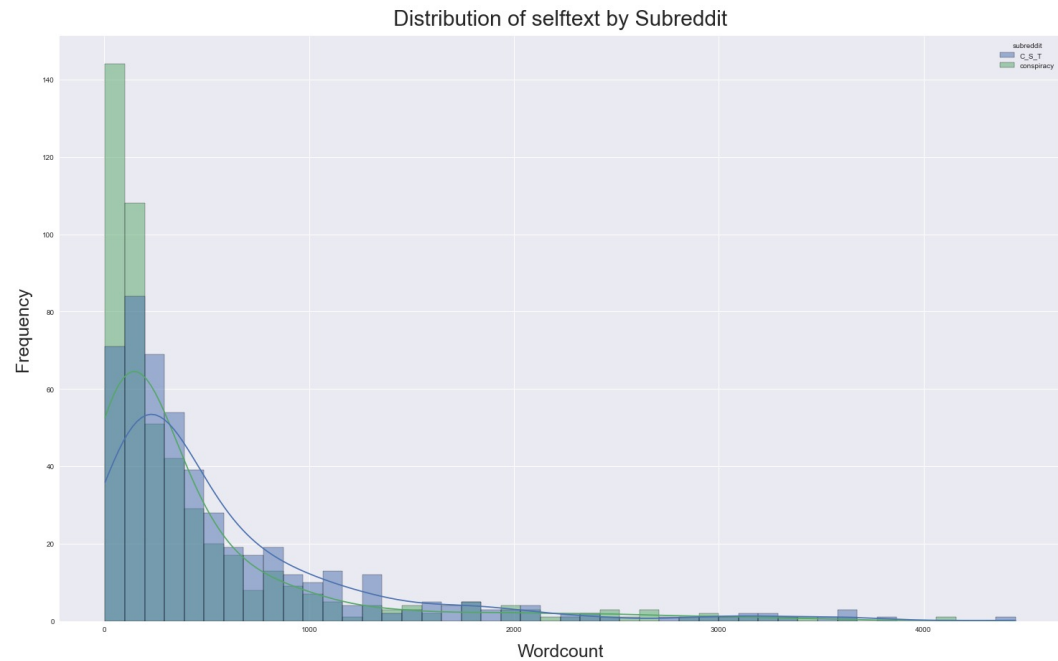
# EXPLORATORY DATA ANALYSIS

## *TITLE WORD COUNT*

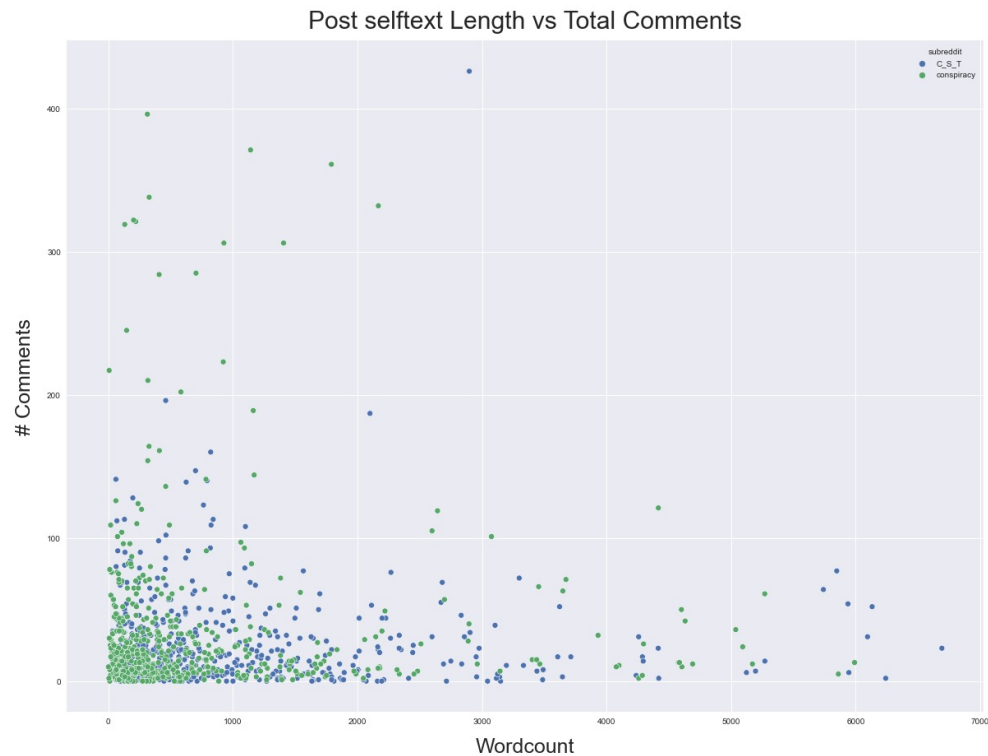


EXPLORATORY  
DATA  
ANALYSIS

*SELFTEXT  
WORD COUNT*

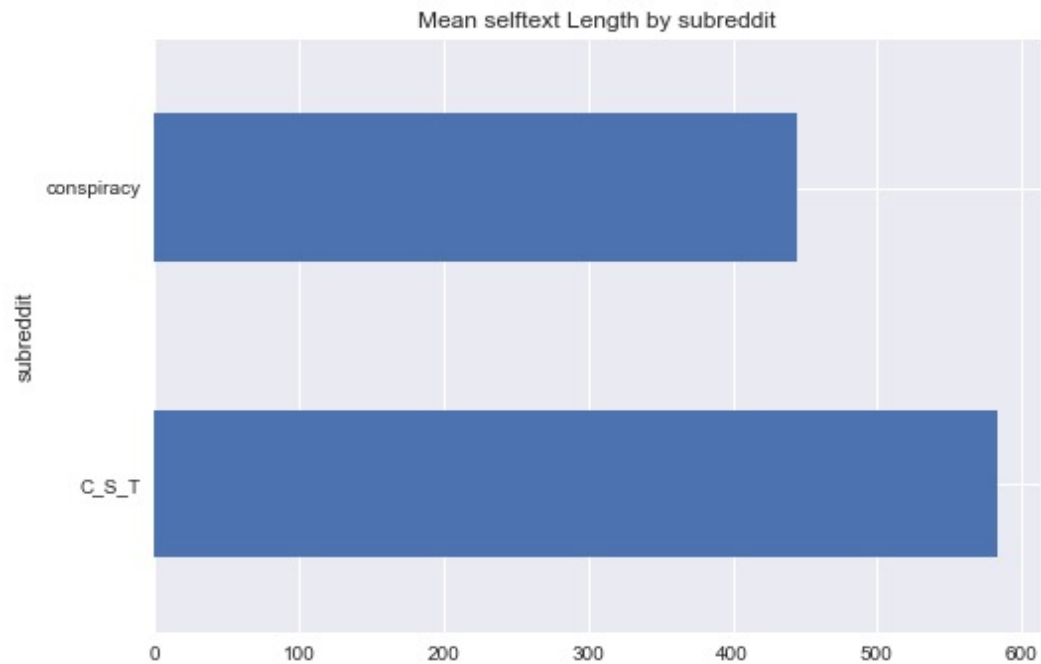


EDA  
POST  
LENGTH VS.  
COMMENTS



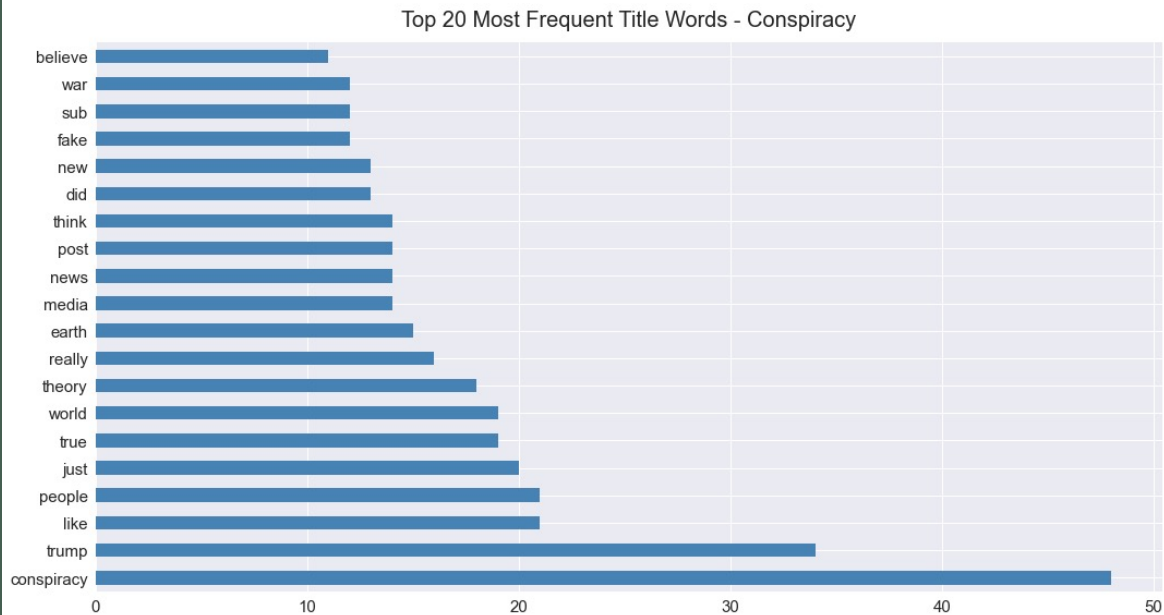
EDA

AVG  
POST  
LENGTH



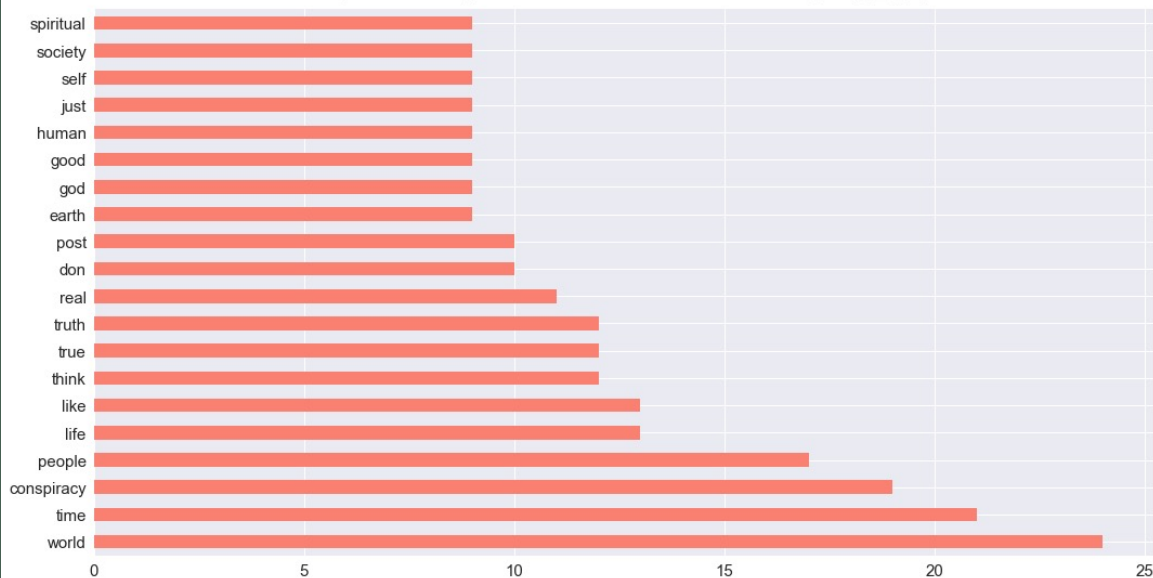


# CONSPIRACY WORDS



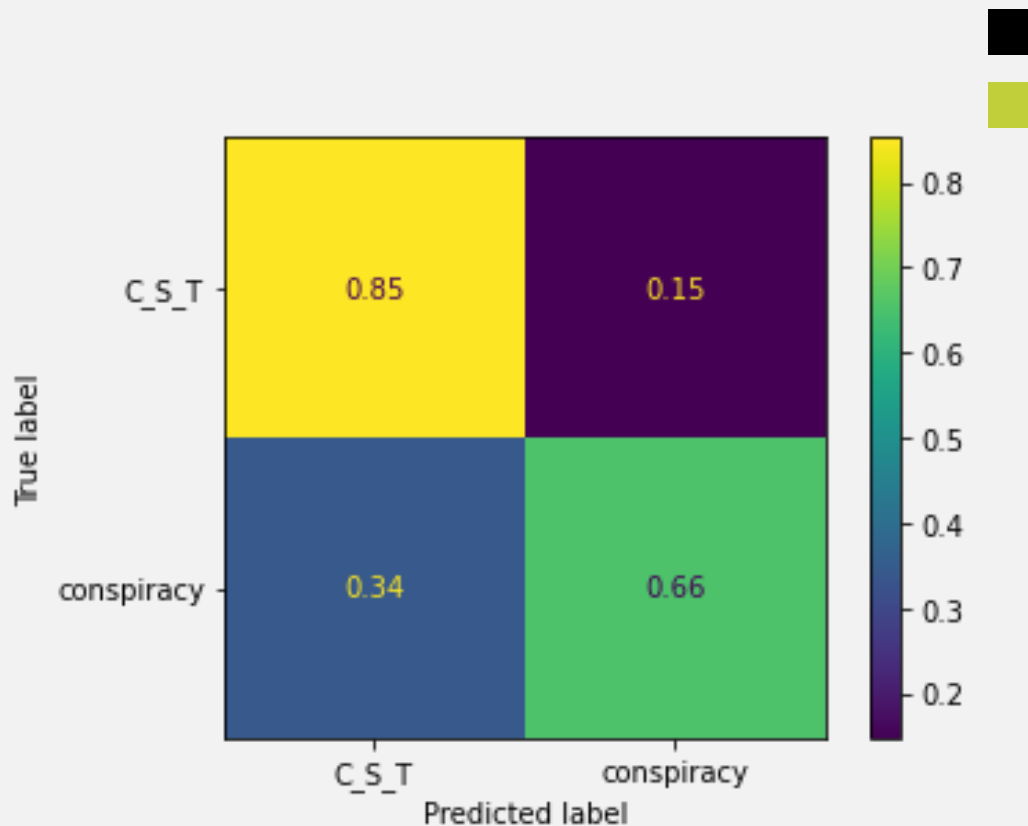
# C\_S\_T WORDS

Top 20 Most Frequent Title Words - Critical Shower Thoughts (C\_S\_T)



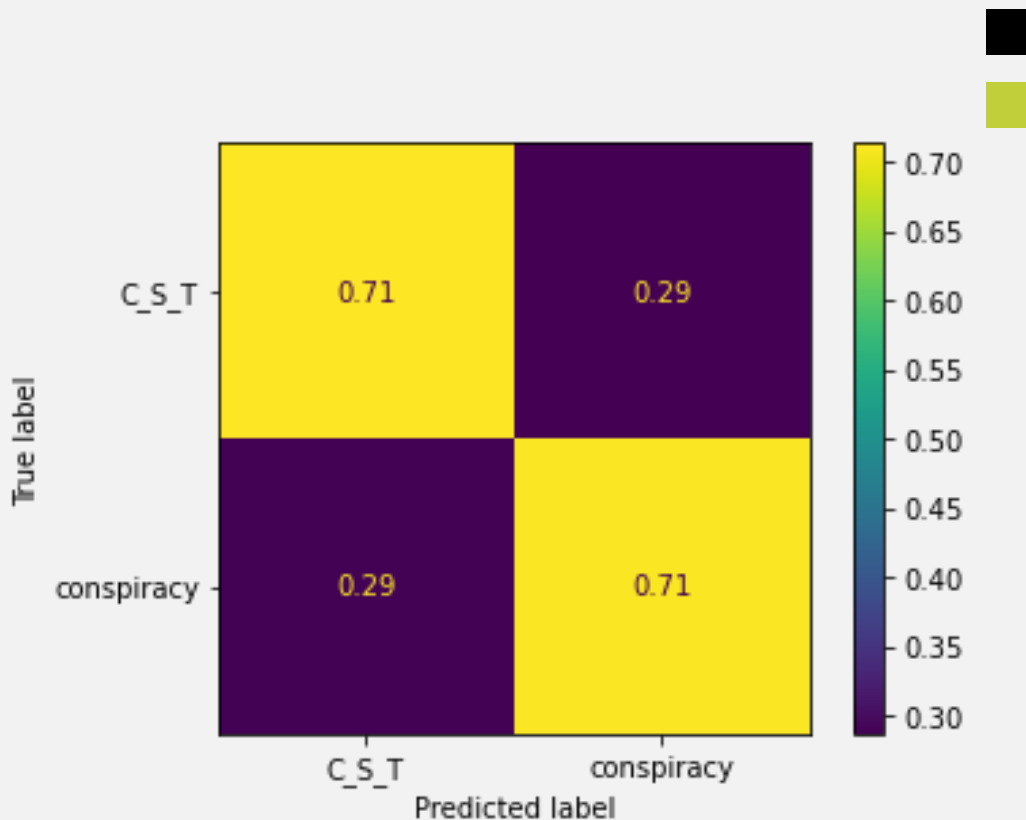
# MODEL 1: COMPLEMENT NB

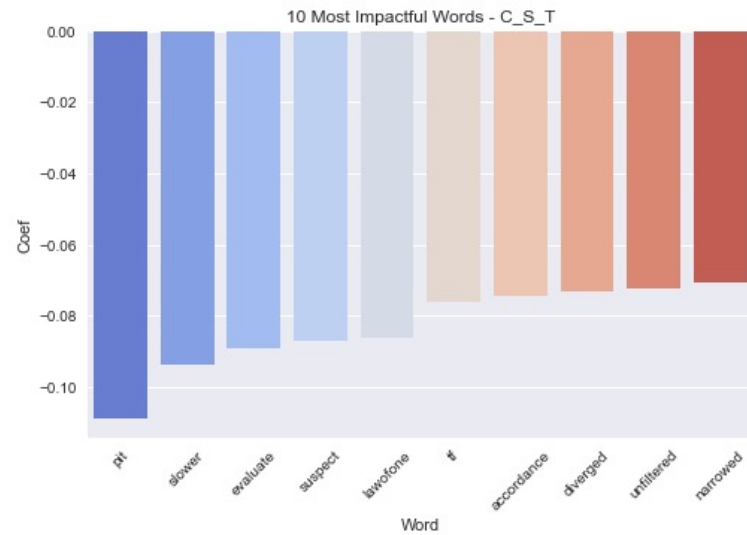
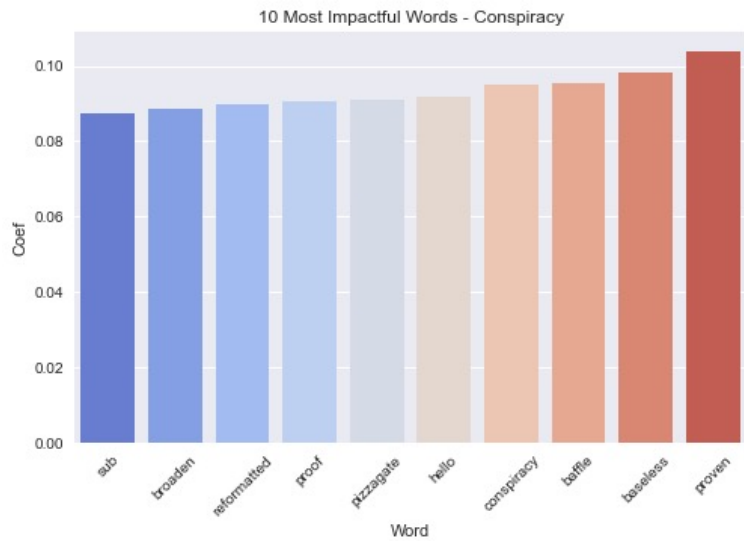
- CountVectorizer & Complement Naïve Bayes on selftext
- **0.94 Train / 0.76 Test**
- *Balanced Accuracy Score: 0.76*
- \*\*our baseline = 0.50



# LOGISTIC REGRESSION

- Using TfidfVectorizer & Logistic Regression, the results were:
- **Best Params** (*max\_iter=5\_000, penalty='elasticnet', C=1.0, solver='saga', l1\_ratio=0.65*)
- **0.77 Train / 0.72 Test**
- **Balanced Accuracy Score: 0.70**





# CONCLUSIONS

## Naïve Bayes Complement Wins!

With a slightly higher train/test, F1 & Baseline Accuracy scores

However, Logistic Regression might win out on an easier dataset. LR does typically have much more interpretable coefs

## Next Time...

Combine variables like num\_comments & score to see if they can increase the overall accuracy of the model.

THANK YOU FOR LISTENING

I'm happy to answer any questions

*Kade Higgins*