

# Web APIs & NLP

---

Predicting between the coffee and tea subreddits



# Problem Statement

---

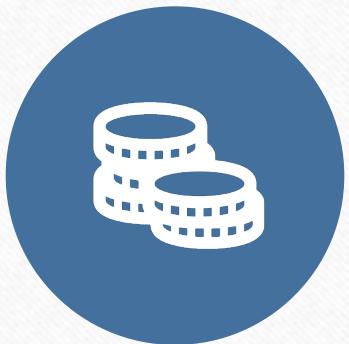
- Up and coming coffee & tea company
- Leverage Data to form sales strategy
- Aid of classification models

# Background

---



COFFEE AND TEA (2/3  
MOST POPULAR  
BEVERAGES WORLDWIDE)



TWO OF THE MOST  
PROFITABLE  
COMMODITIES.



PEOPLE LOVE CAFFEINE

# The Data

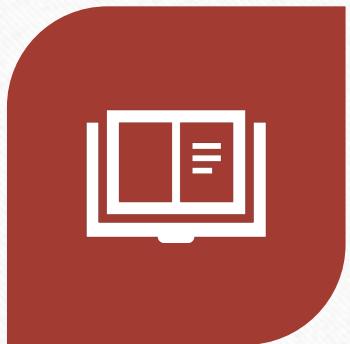
---



PUSHSHIFT API

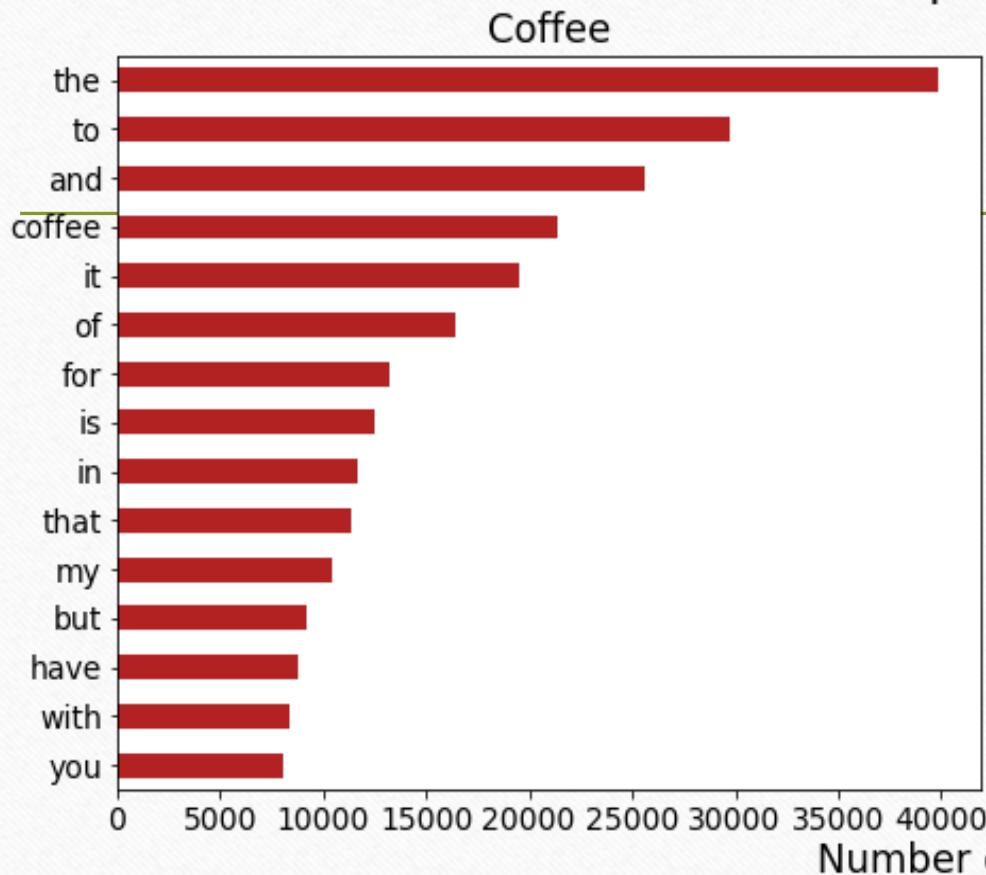


SCRAPED TEA &  
COFFEE SUBREDDITS



19592 POSTS

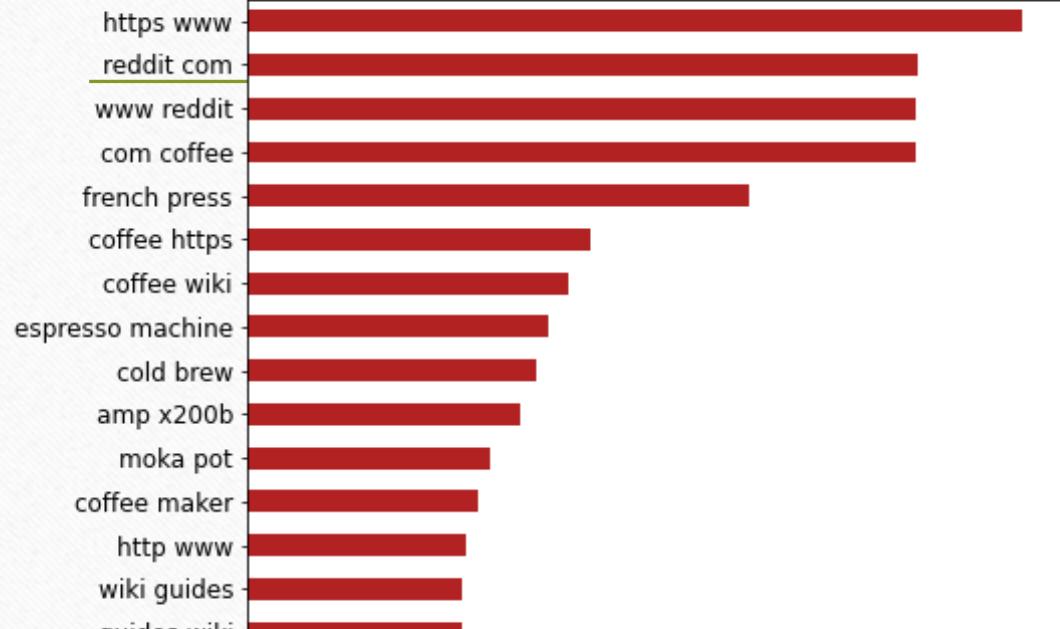
## Top 15 Words



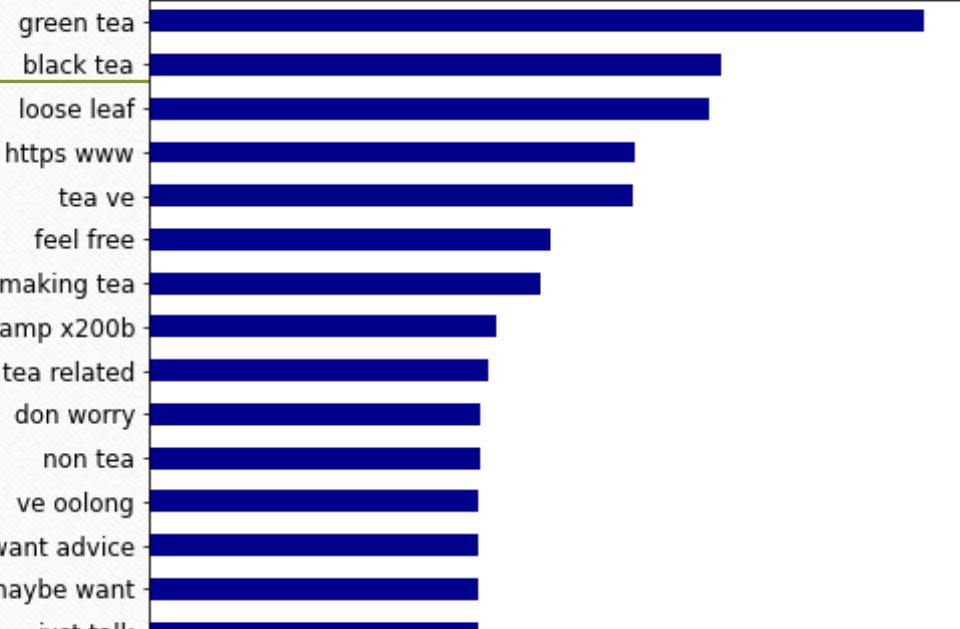
Number of Occurrences

### Top 15 Word Pairs(No stop words)

Coffee



Tea



Frequency

# Model Evaluation

---



ACCURACY



RECALL

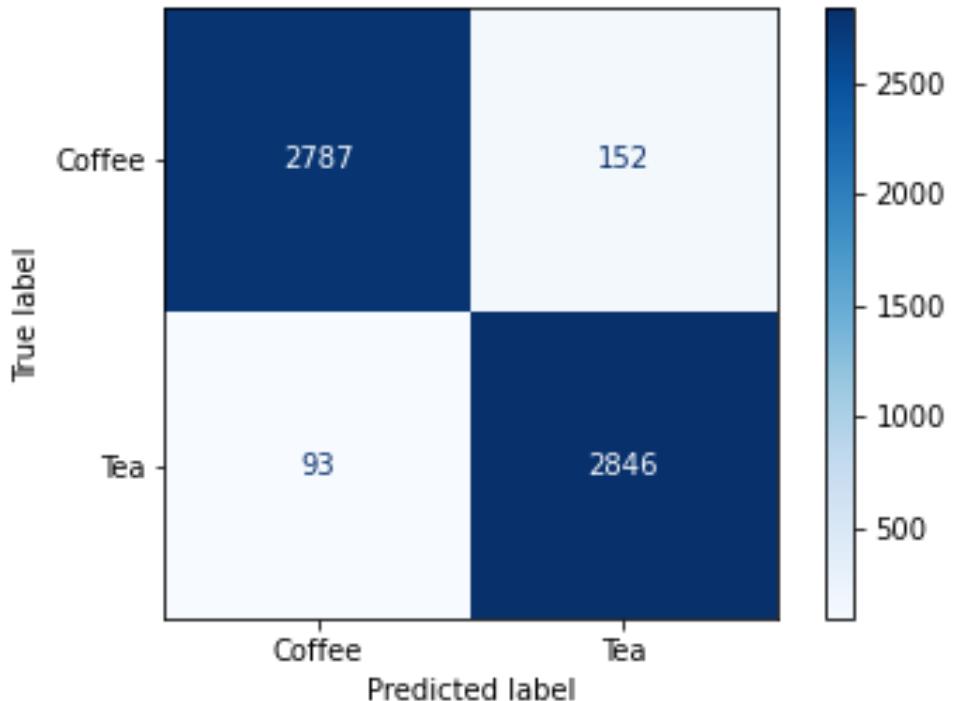


SPECIFICITY



TYPE I AND  
TYPE II

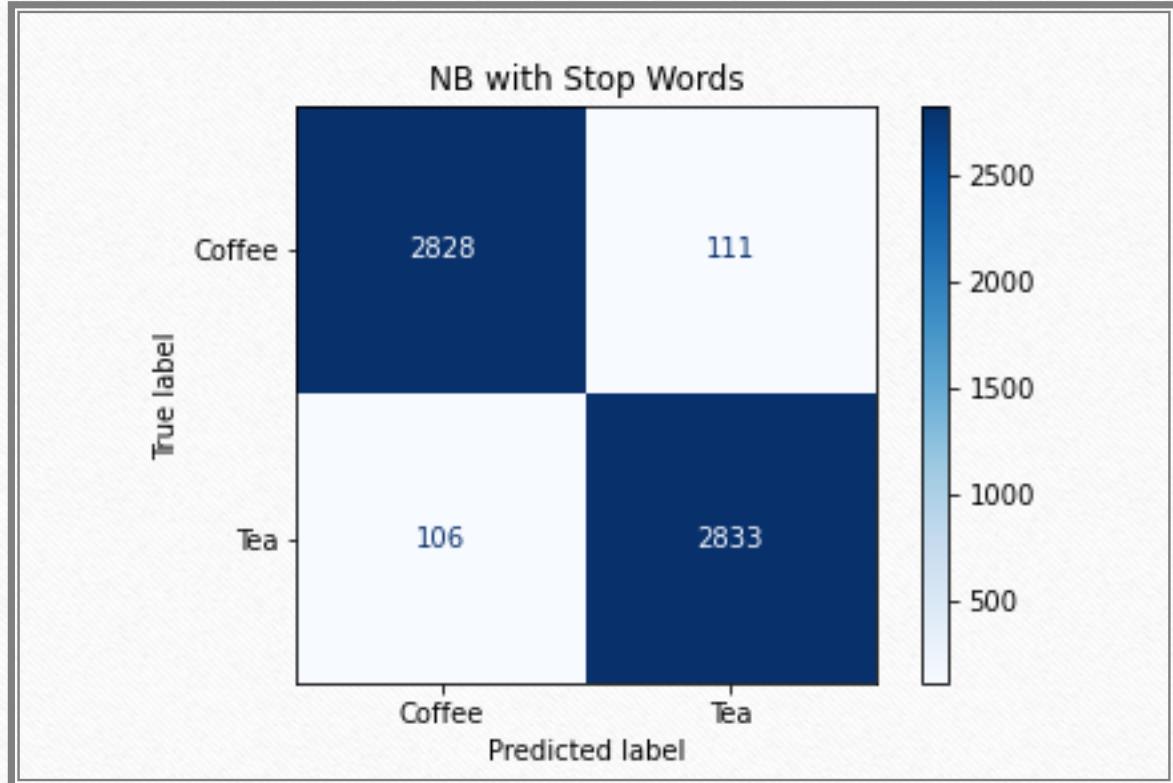
# Multinomial NB



|                | Multinomial NB Model |
|----------------|----------------------|
| Training Score | 96.1%                |
| Testing Score  | 95.8%                |
| Specificity    | 94.8%                |
| Recall         | 96.8%                |
| Type I error   | 5.2%                 |
| Type II error  | 3.2%                 |

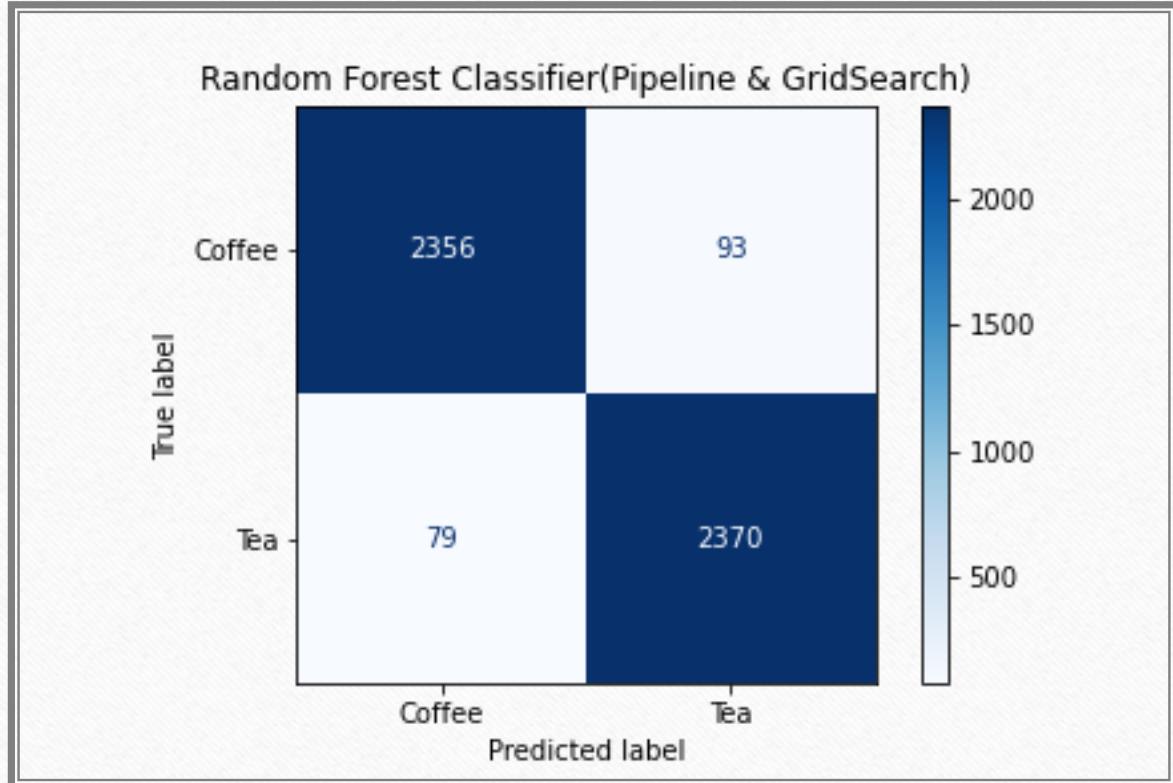
## Multinomial NB(w/stop words)

|                | Multinomial NB (w/Stop words) |
|----------------|-------------------------------|
| Training Score | 96.2%                         |
| Testing Score  | 96.3%                         |
| Specificity    | 95.4%                         |
| Recall         | 96.4%                         |
| Type I error   | 3.8%                          |
| Type II error  | 3.6%                          |



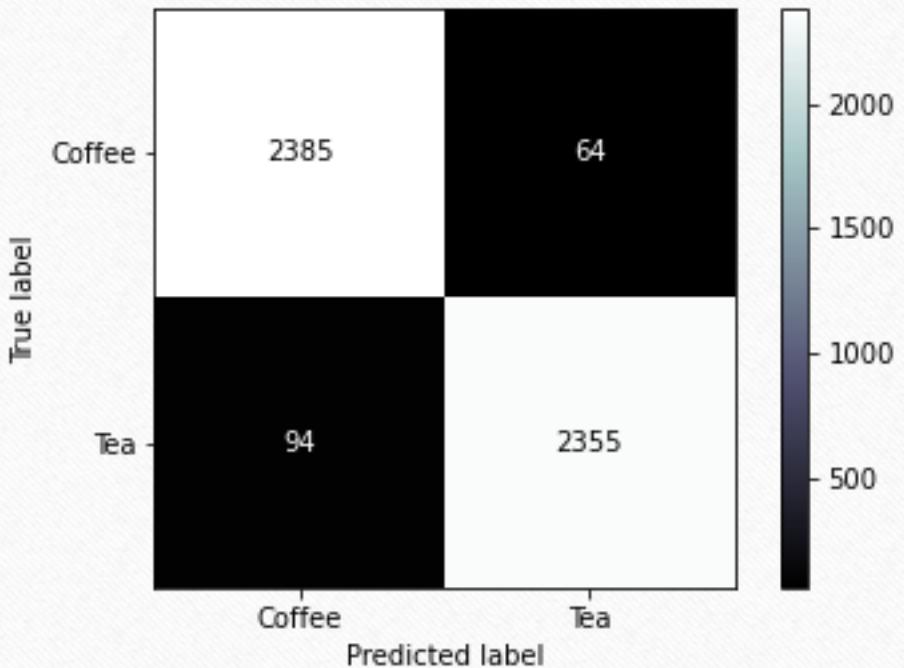
# Random Forest Classification Model

|                | Random Forest Classifier<br>(Pipeline &<br>Grid search) |
|----------------|---------------------------------------------------------|
| Training Score | 99.9%                                                   |
| Testing Score  | 96.5%                                                   |
| Specificity    | 96.1%                                                   |
| Recall         | 98.3%                                                   |
| Type I error   | 3.8%                                                    |
| Type II error  | 3.2%                                                    |



# Logistic Regression

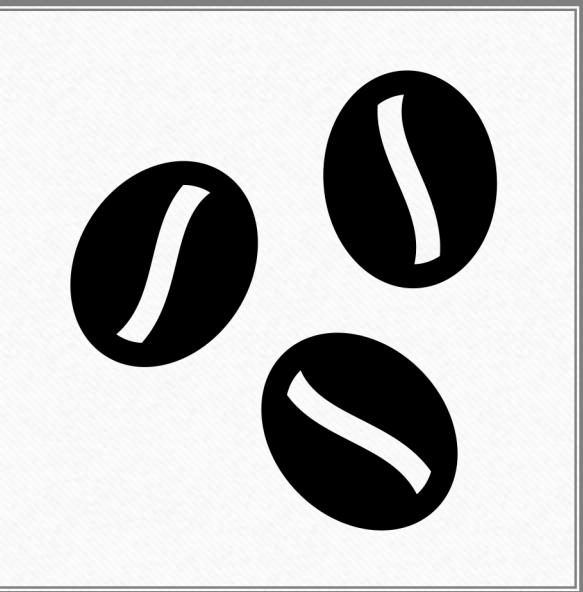
|                | Logistic Regression(TFID/GS) |
|----------------|------------------------------|
| Training Score | 99.9%                        |
| Testing Score  | 96.8%                        |
| Specificity    | 97.4%                        |
| Recall         | 96.2%                        |
| Type I error   | 2.6%                         |
| Type 2 error   | 3.8%                         |



|                       | Multinomial NB Model | Multinomial NB (w/Stop words) | Random Forest(TFID/GS) | Random Forest Classifier (Pipeline& Grid search) | Logistic Regression(TFID) |
|-----------------------|----------------------|-------------------------------|------------------------|--------------------------------------------------|---------------------------|
| <b>Training Score</b> | 96.1%                | 96.2%                         | 99.9%                  | 99.9%                                            | 98.9%                     |
| <b>Testing Score</b>  | 95.8%                | 96.3%                         | 96.3%                  | 96.5%                                            | 96.8%                     |
| <b>Specificity</b>    | 94.8%                | 95.4%                         | 96.1%                  | 96.1%                                            | 97.4%                     |
| <b>Recall</b>         | 96.8%                | 96.4%                         | 96.5%                  | 98.3%                                            | 96.2%                     |
| <b>Type I error</b>   | 5.2%                 | 3.8%                          | 3.9%                   | 3.8%                                             | 2.6%                      |
| <b>Type II error</b>  | 3.2%                 | 3.6%                          | 3.5%                   | 3.2%                                             | 3.8%                      |

# Conclusion

---



Word Pairs tell a story

Logistic Regression performed the best

Looking at a company's sales data to complete the picture