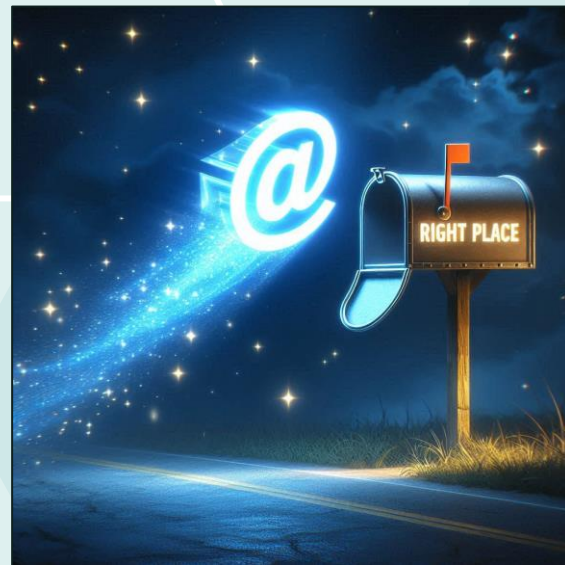# Spam Email Filter - Machine Learning Model

# Agenda

# 01

# Business Problem

Our telecommunications company wants to break into the email market and needs to implement an effective SPAM filter.

We are charged with:

- Create a model that reliably detects SPAM in a way that will maximize customer satisfaction

# 02

# Data Analysis

Actual emails in the form of individual HTML files were used in the analysis

- 5280 HTML files in total

- Each email was previously labeled as HAM or SPAM and sorted into folders

# Email Importing and Preprocessing

# Subject Line vs Email Body

Subject Line:
- We expect 'dense' information about the nature of the email
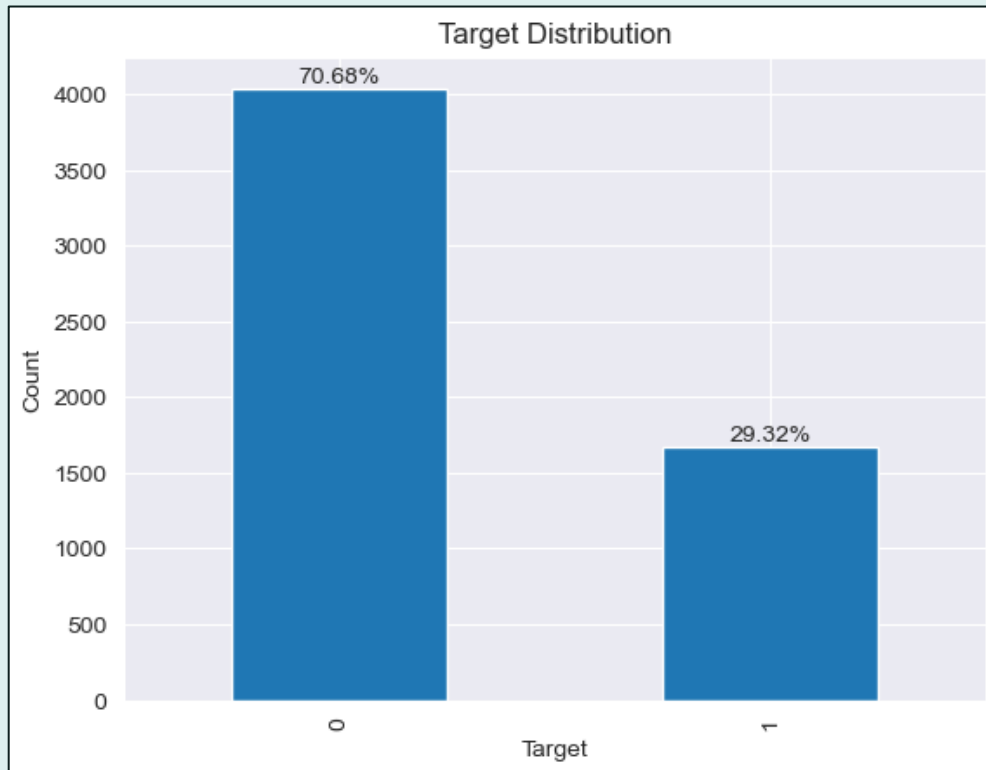
Email Body:
- More information about the nature of the email, more verbose

```
Subject: Species at risk of extinction growing


The latest "Red List" adds 124 to the 11,000 endangered species around the
globe - but also includes a stick insect revival
```

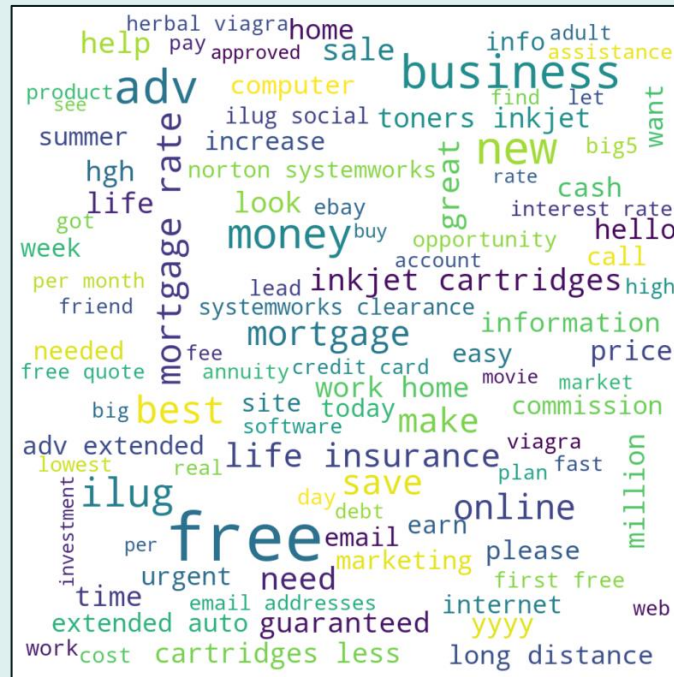➢ Modeling the Subject and Body separately...

# Target Distribution Imbalance
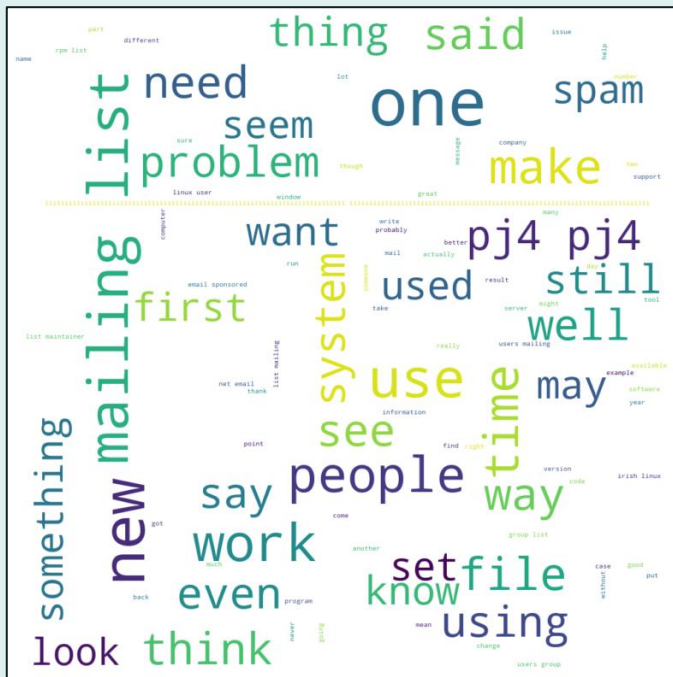
# Subject Line Word Clouds by Class



HAM

SPAM

Note: Models are built on company internal email data, will need to adjust models trained on anonymized customer emails data after deployment

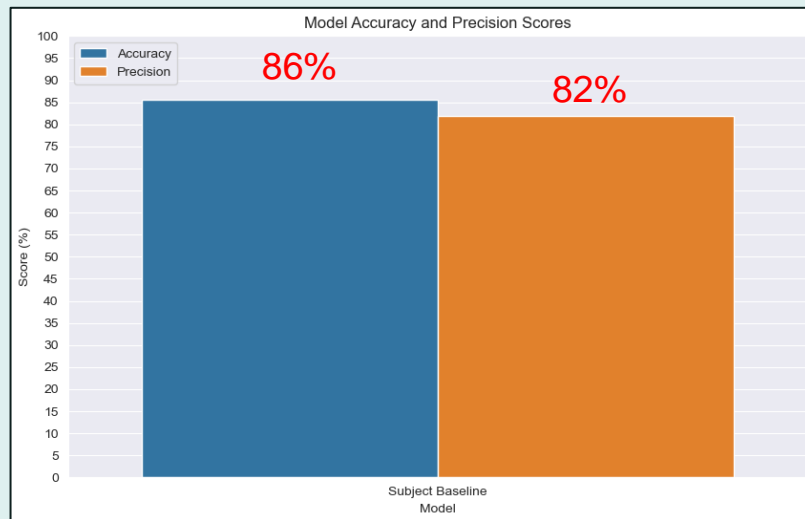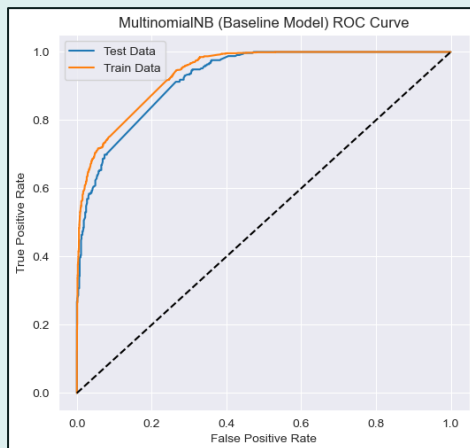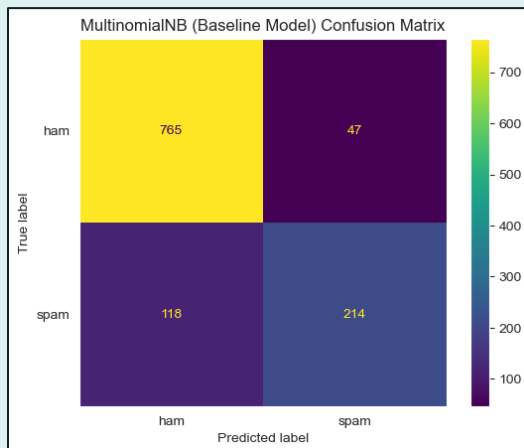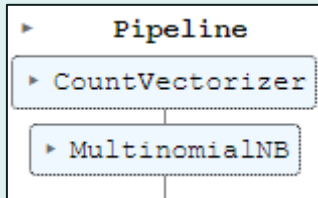# Email Body Word Clouds by Class

HAM



SPAM

# 03 Predictive Modeling



BUSINESS REQUIREMENT:

- Getting SPAM in the inbox irritates our customers, but missing an important message because we classified it as SPAM will really anger them!

- Optimizing models on **Precision (minimize False-Pos)**, then **Accuracy**

- Less than 1% of customers HAM emails are allowed to go to the SPAM classification
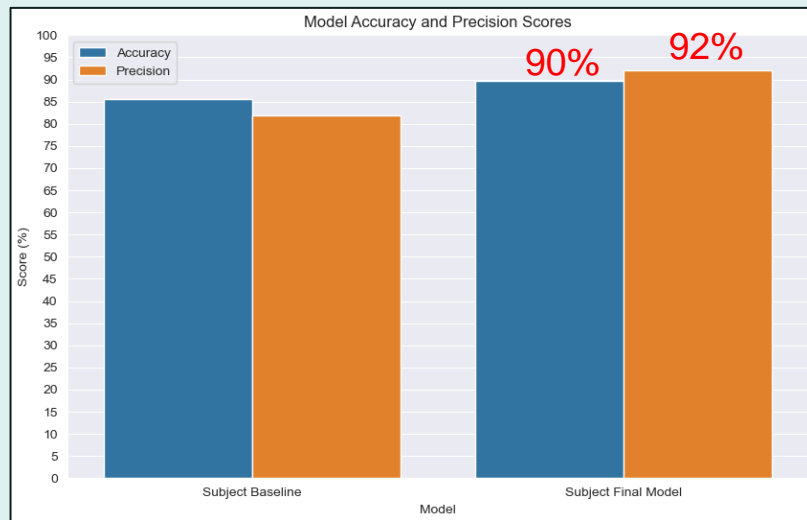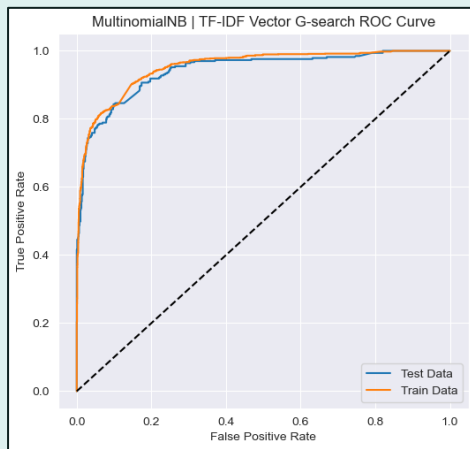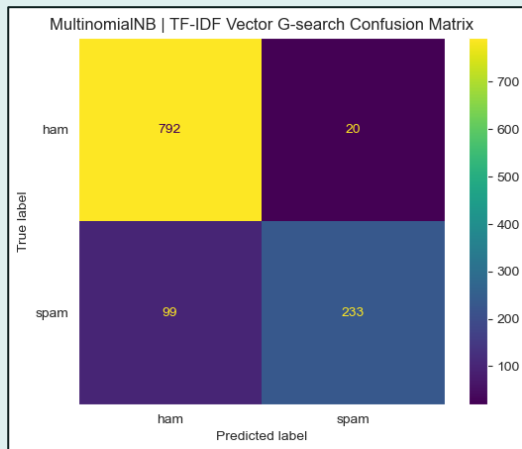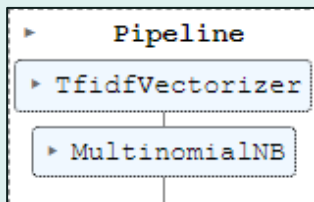
# Subject Line Baseline Model

- Regex tokenizer of >2 characters
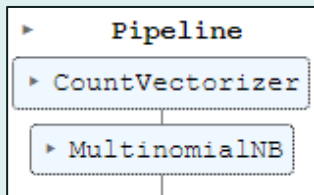- Binary vectorizer
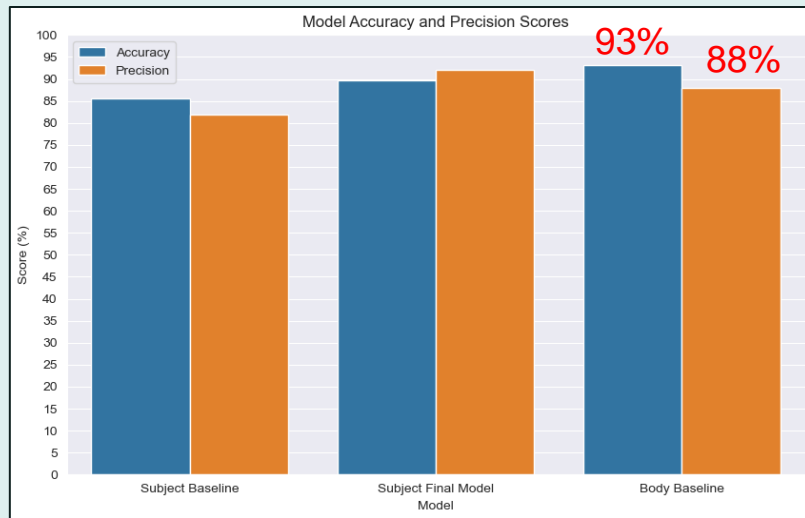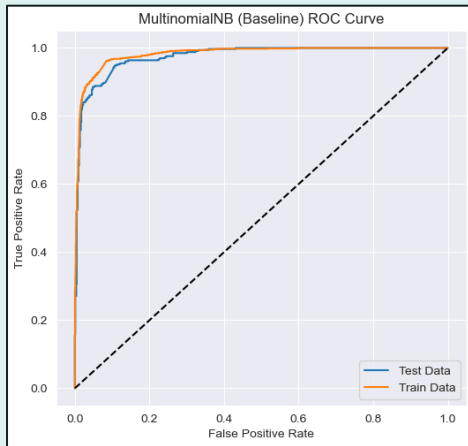- Multinomial Bayes Classifier
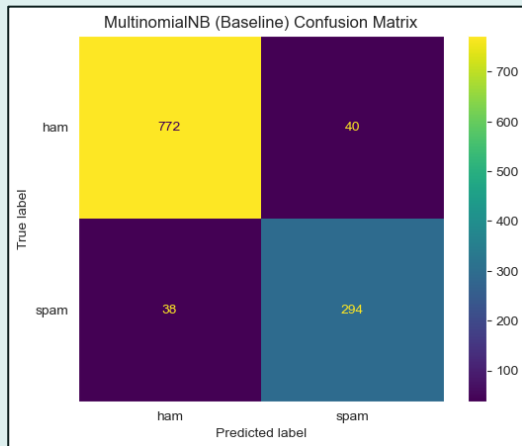
# Subject Line Final Model

- Word_punct_tokenizer
- TFIDF vectorizer
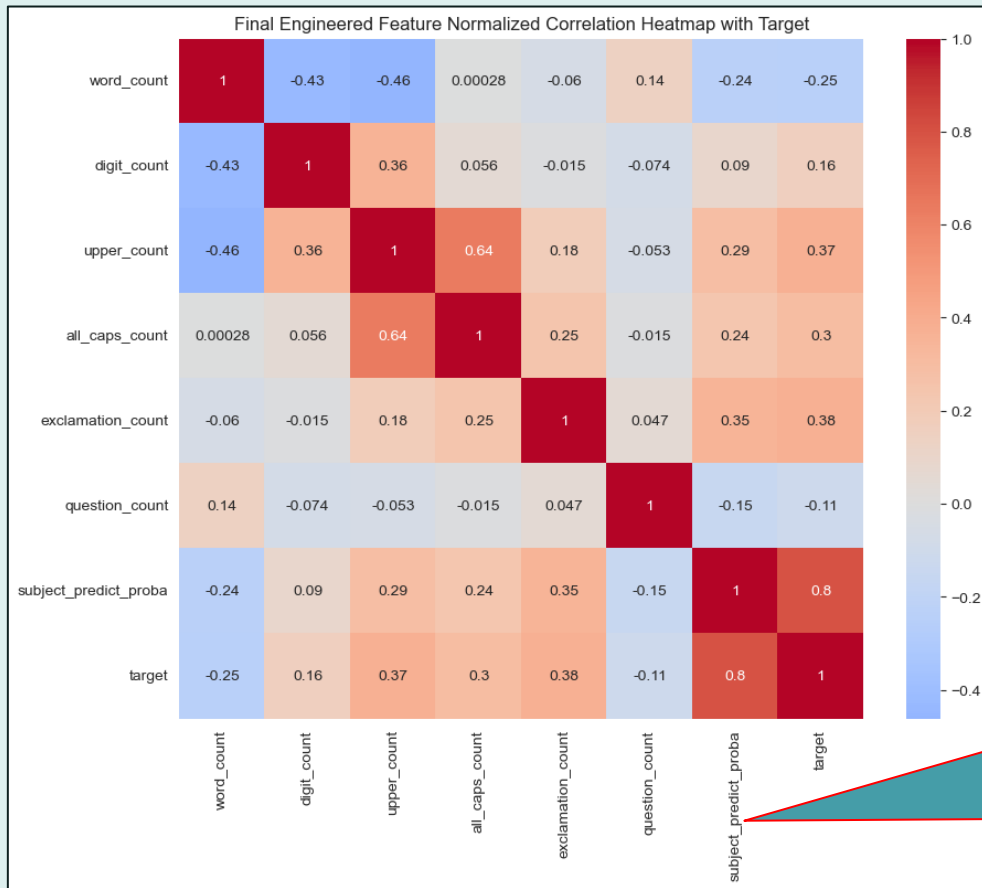- Multinomial Bayes Classifier

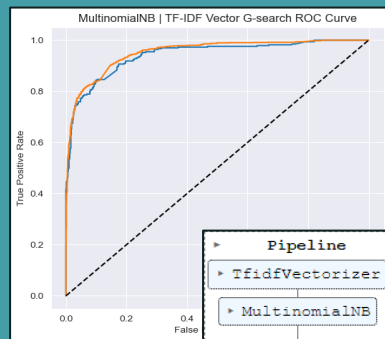# Email Body Baseline Model



- Regex tokenizer of >2 characters
- Binary vectorizer
- Multinomial Bayes Classifier

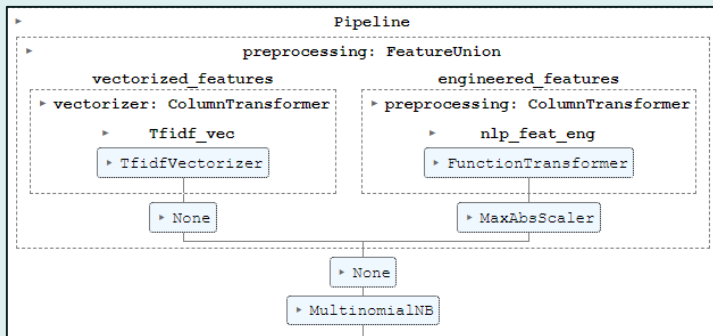# Email Body Feature Engineering
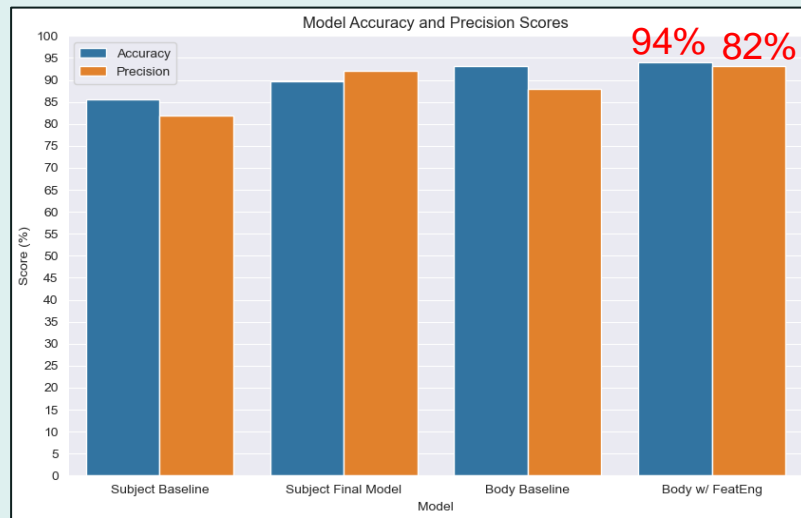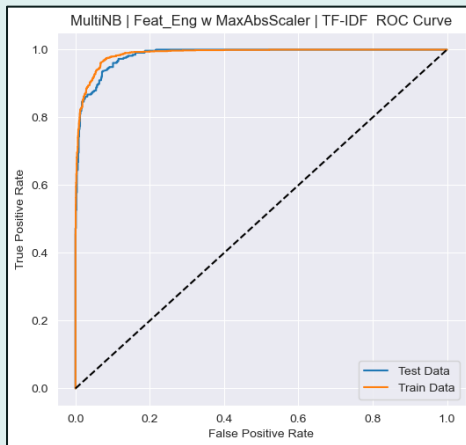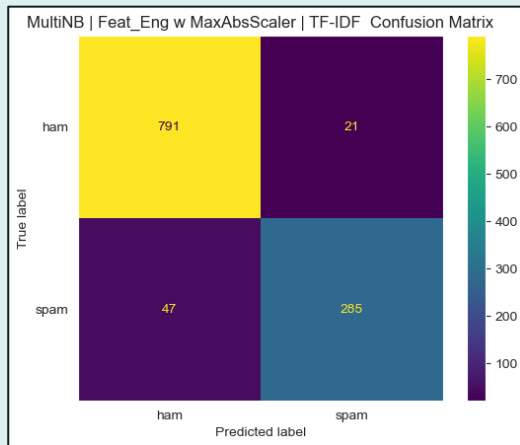


Final Engineered Feature Normalized Correlation Heatmap with Target
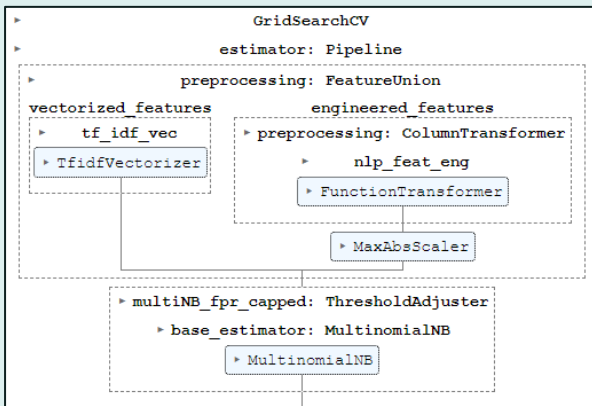
SUBJECT LINE PREDICTIONS
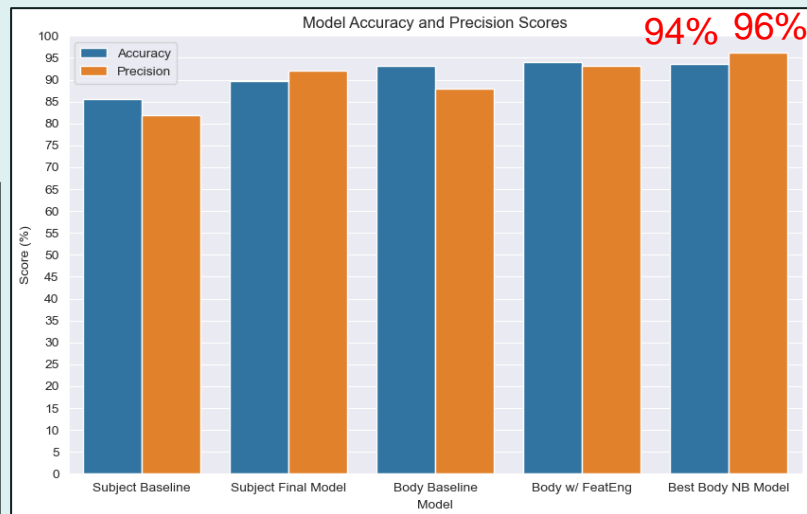
# Email Body With Feature Engineering



- Regex tokenizer of >2 characters
- TF-IDF vectorizer
- Multinomial Bayes Classifier

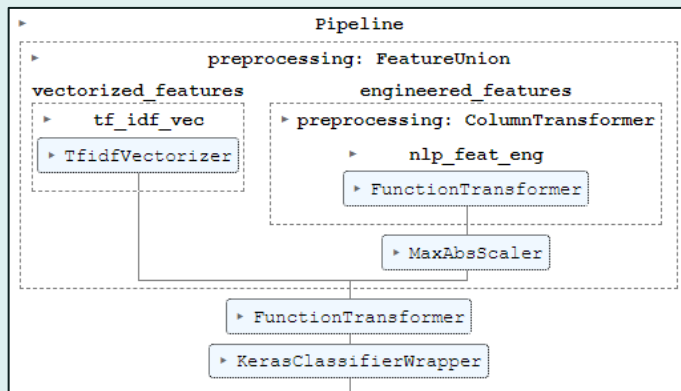# Email Body – Best Naive Bayes Model



- Hyperparameter Gridsearch
- NB Model wrapped in Prob Threshold Adjuster  (<1% FPR)

# Email Body – Neural Network Baseline



- One hidden layer
- Threshold Adjusted (<1% FPR)



97%  96%

# Email Body – Neural Network Final



- 5 hidden layers
- 2 dropout layers
- Gridsearch using Hyperband()
- Threshold Adjusted (<1% FPR)

# 04 Conclusion



Model Accuracy and Precision Scores

**Final Model Test Metrics:**

➢ Precision:  **97.7%**
➢ Accuracy:  **97.8%**

**Final Business Metrics:**

➢ False Pos Rate
   (HAM going to SPAM):  **0.9%**

➢ False Neg Rate
   (SPAM going to HAM):  **6%**

# 05 Recommendation

1. Implement the best Neural Network model during initial Email service rollout

2. Continue to retrain and retune the models on new anonymized customer data

3. Investigate other potential model improvements (see next steps)

# 06 Next Steps

1. Continue to improve the model

   1. Rerun and tune the models with the 'Cheat words' removed

   2. Investigate tokenizing based on wordcloud algorithm

   3. Check other model architectures (RandomForest, XGBoost)

   4. Tune NN for subject-line only, perhaps effective and more efficient

   5. Investigate verbiage and nature of SPAM, are there features we can capture

   6. Investigate the False cases; are there features we can capture

   7. Investigate keeping URLs instead of scrubbing them preprocess

2. Acquire larger data sets of up-to-date emails

# Thanks!

**Questions? Please contact me at:**

Dale DeFord
daledeford@gmail.com

 