

# Spam Classification

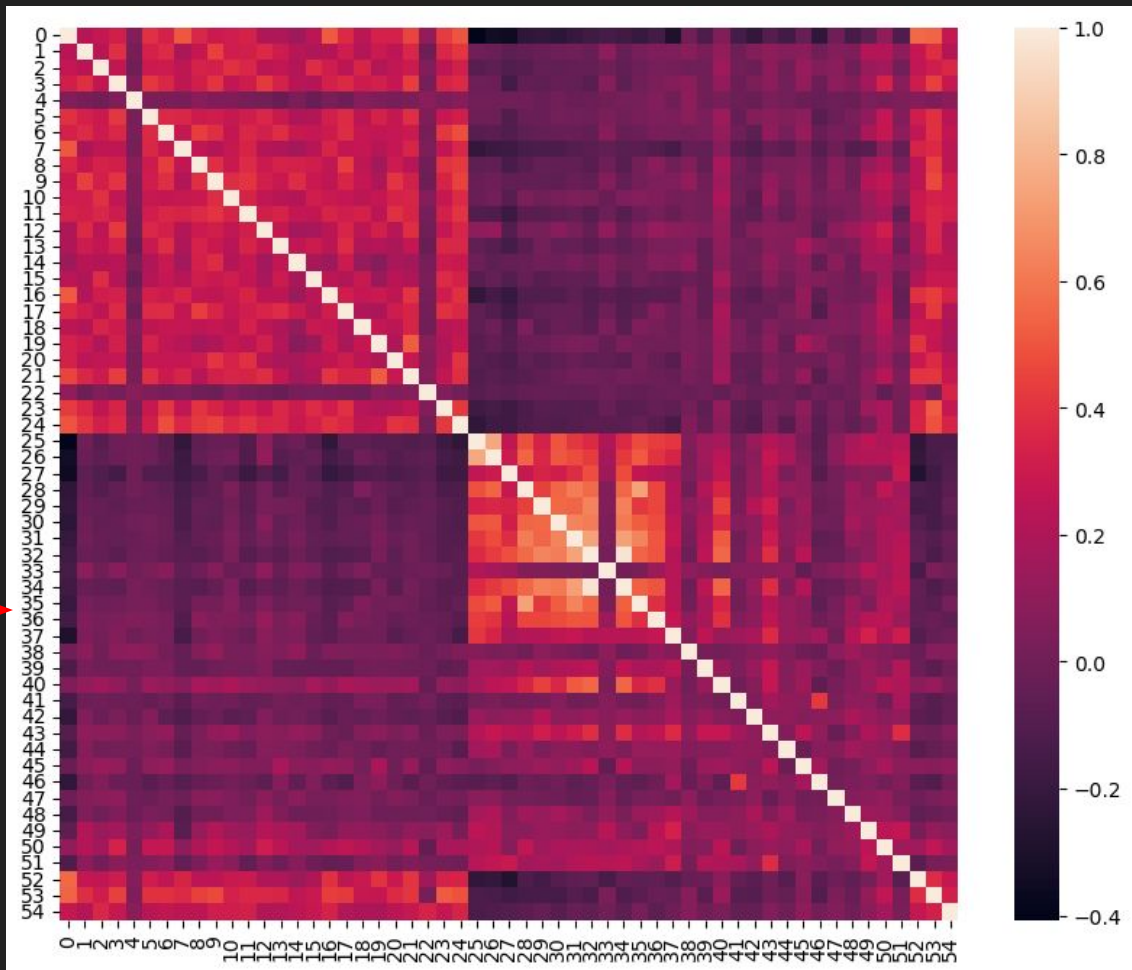
## - Naive Bayes

# What's in the data?



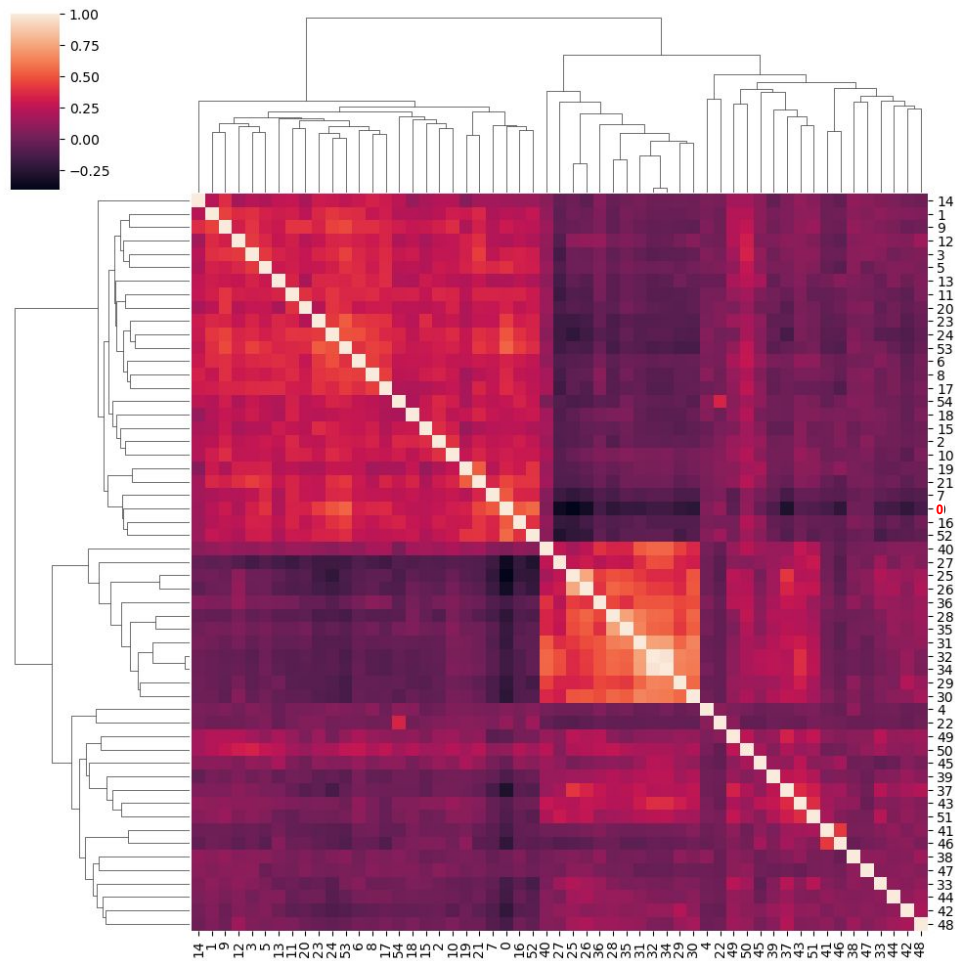
# Analysis

Heat Map of all features



# Analysis

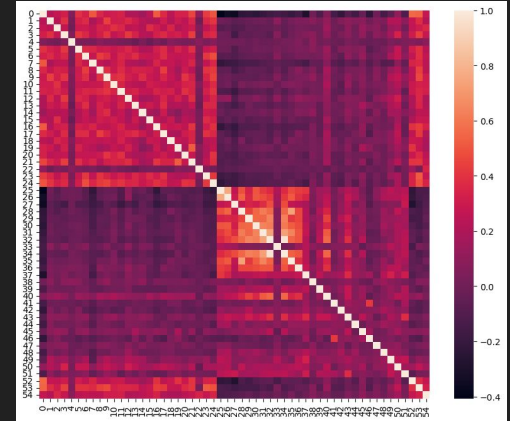
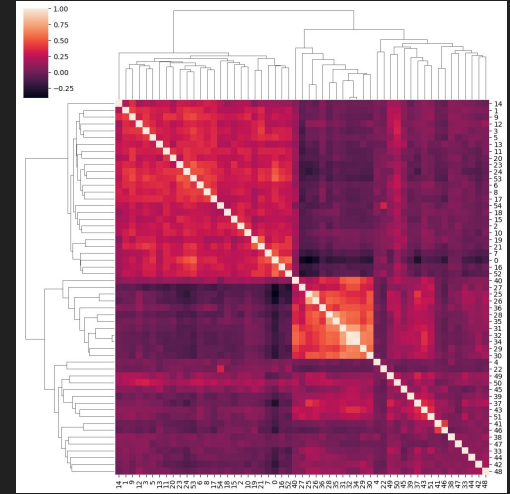
Cluster Map of all features



# Analysis of data set - why Naive Bayes?

- K-Nearest Neighbours 💀
- Naive Bayes ?
- Neural Net 💕💕

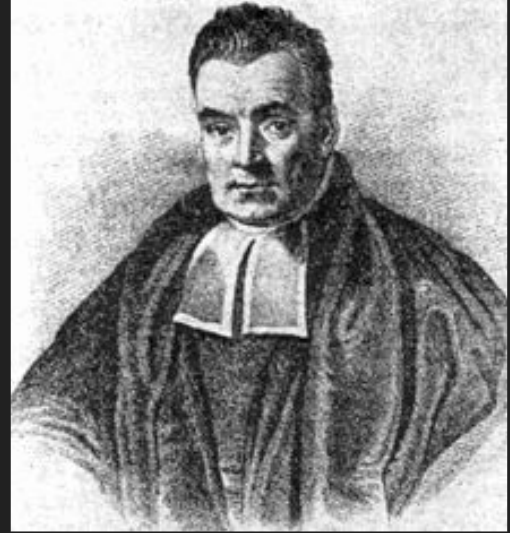
Whilst clustermap shows high dependency... I was curious to see NB's efficiency regardless



# Naive bayes - How does it work?

## Core Assumptions

- All features are independent
- Features are all equally important



Mr Thomas Bayes

At the complete basics :

Naive Bayes classifier works by calculating and comparing the probability of classes based on probabilities of features.

Most likely class for a new message is  $\hat{y} = \max_{class} P(class|message)$

Which simplifies down

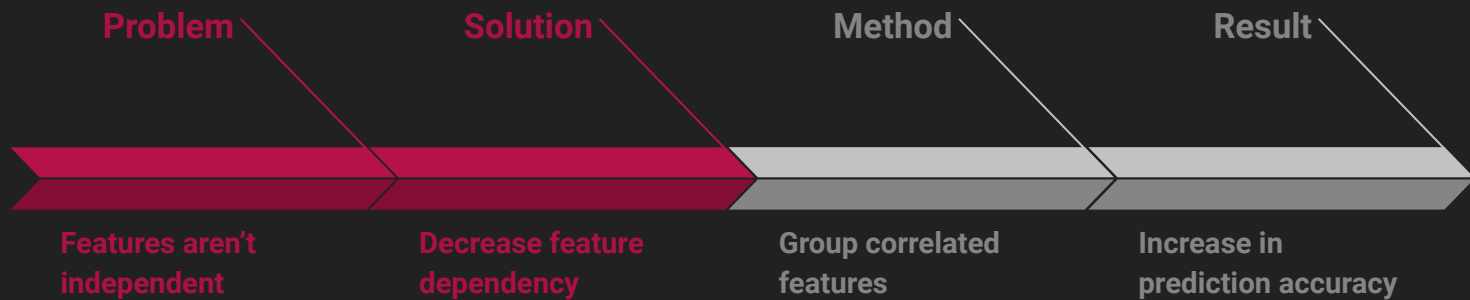
$$\hat{y} = \max_{class} P(class) \cdot P(message|class)$$

# Optimisations - Logarithmic

Logarithmic - reduce information loss from underflow of multiplying probabilities

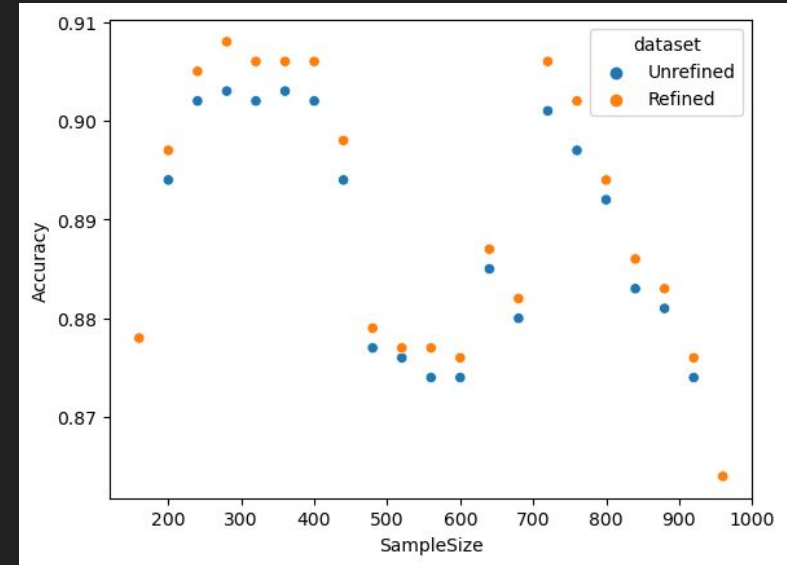
Grouping correlated columns - increase feature independency, good because naive works on assumption of all features independent






# Optimisations - Grouping correlated columns

Removing High Feature Correlation	
PMCC threshold <b>0.7</b>	PMCC threshold <b>0.6</b>
Grouping lead to higher feature independency <b>LEAD TO</b> Higher independence and balanced weighting <b>Result</b> Minor <b>Gain</b> in accuracy	Grouping lead to higher feature independency <b>BUT</b> Too much information lost from removal column <b>Result</b> Minor <b>Loss</b> in accuracy



Columns removed to provide Refined data:  
[25, 28, 31]  
21 Samples Sizes  
Refined accuracy total = 18.693 and  
Unrefined total = 18.636



Is it worth  
anything?



IMPROVMENT  
IS  
IMPROVEMENT





IMPROVEMENT  
IS  
IMPROVEMENT



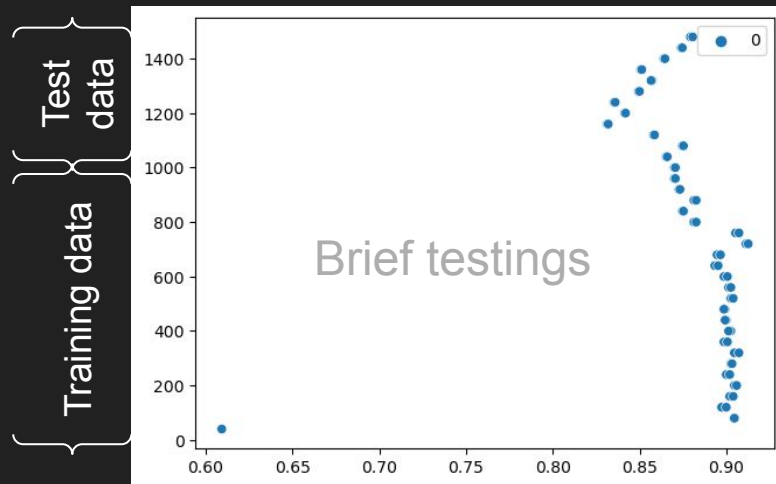
IMPROVMENT  
IS  
IMPROVEMENT

# Accuracy belief

If test data is representative, then accuracy of roughly 0.85

if test data contains anomalies, more akin to 0.8 or 0.75

Hence final bet is **0.8**



# If i had more time

Implement a Neural Network

Incorporate Machine Learning into Naive Bayes

- Detect trends alongside single feature probabilities





### **Difficulty of approach (10%)**

How ambitious and complicated? How well is it suited for task? Clearly state the approach you have taken to solve the task and why you have chosen it.

### **Description of algorithm (10%)**

How well understood and explained is algorithm? Describe algorithm you have implemented in terminology supporting mathematical notation

### **Implementation and Optimisation (10%)**

optimisation of your chosen algorithm? Why you made these choices

### **Explanation and Contextualisation of Results (10%)**

How well have you provided a contextualised assessment of your predicted accuracy? Both strengths and limitations. This could include proposing any future work that could address the current limitations of your model.

### **Presentation and Polish (10%)**

How well presented is your video? your video must deliver information clearly.