

# Project 5 - Naive Bayes Classifier

Jefferson Roylance  
Justin Fairbourn

# Dataset

Over 1 million reddit comments, with the following columns and labeled as sarcastic or not.

- label
- comment
- author
- subreddit
- score (# of upvotes - # of downvotes)
- ups (# of upvotes)
- downs (# of downvotes)
- date
- created\_utc
- parent\_comment

# Analysis

- Used multinomial bayes classifier
  - This was because the data didn't seem to be normally distributed, and tests confirmed that multinomial got better results than gaussian
- 2 sections - textual analysis of the comment and analysis of the rest of the data

Textual analysis of the comment examined the following characteristics using 4-fold cross-validation:

- Amounts of each letter
- Length of the comment (in characters)
- Presence of punctuation (boolean)
- Average word length
- Words used - this was found out by taking the top 500 words used in all comments and then finding the counts of each of those words in the comment
- Checking for predefined patterns (we only got around to checking for the presence of '...')
- Number of uppercase letters

# Results

Using just the comment textual analysis, we got an f-score of 0.6209

Use cases include automatic classification of sarcastic comments by social media as well as messaging apps.