# Subreddit Classification

Shandeep
DSI14

# Table of Contents

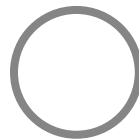# Define Problem



**Problem Statement:**

The aim of this project was to analyse the language used in these subreddits and to examine whether or not they were unique to both Singapore and the United Kingdom

# Gather Data

**Data Collection**: r/singapore & r/unitedkingdom

**Number of data:** 944 and 956 comments

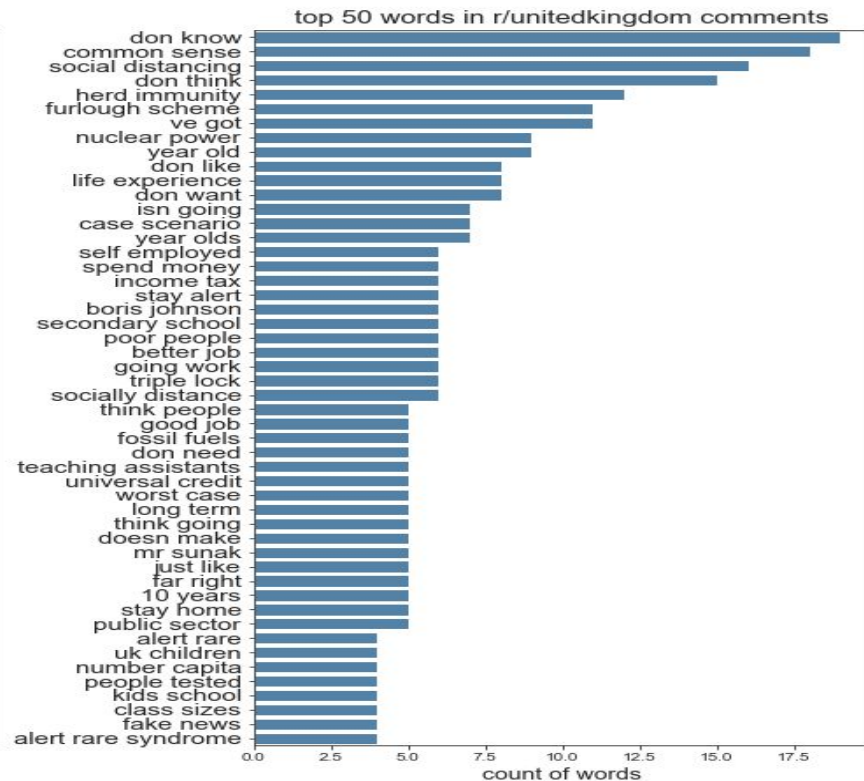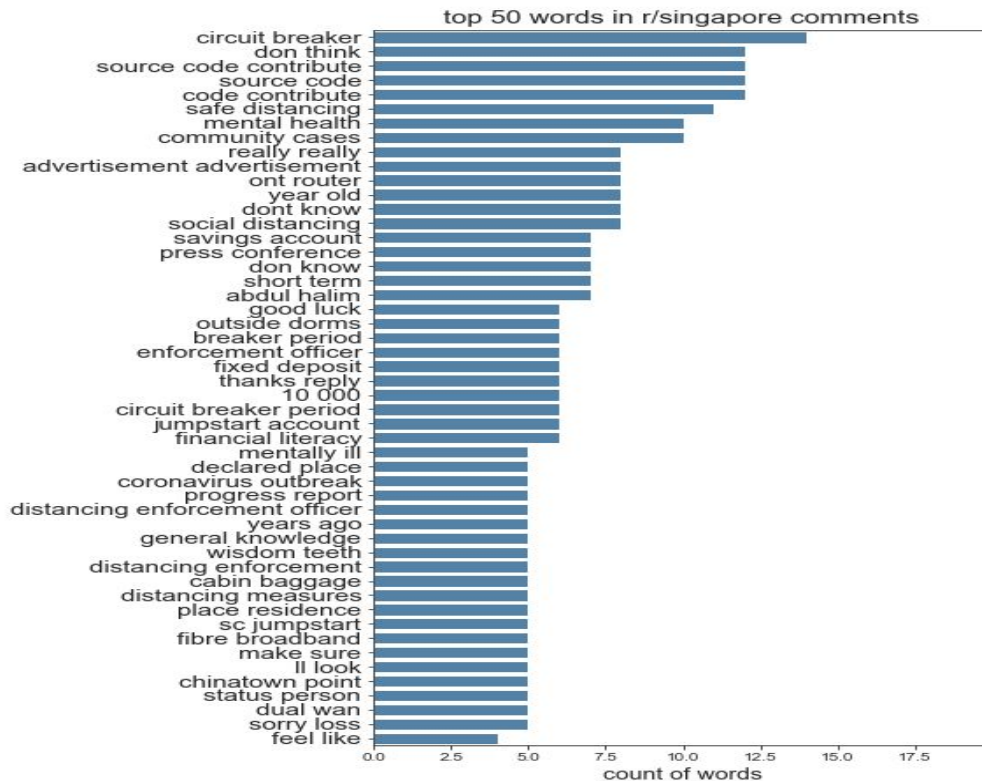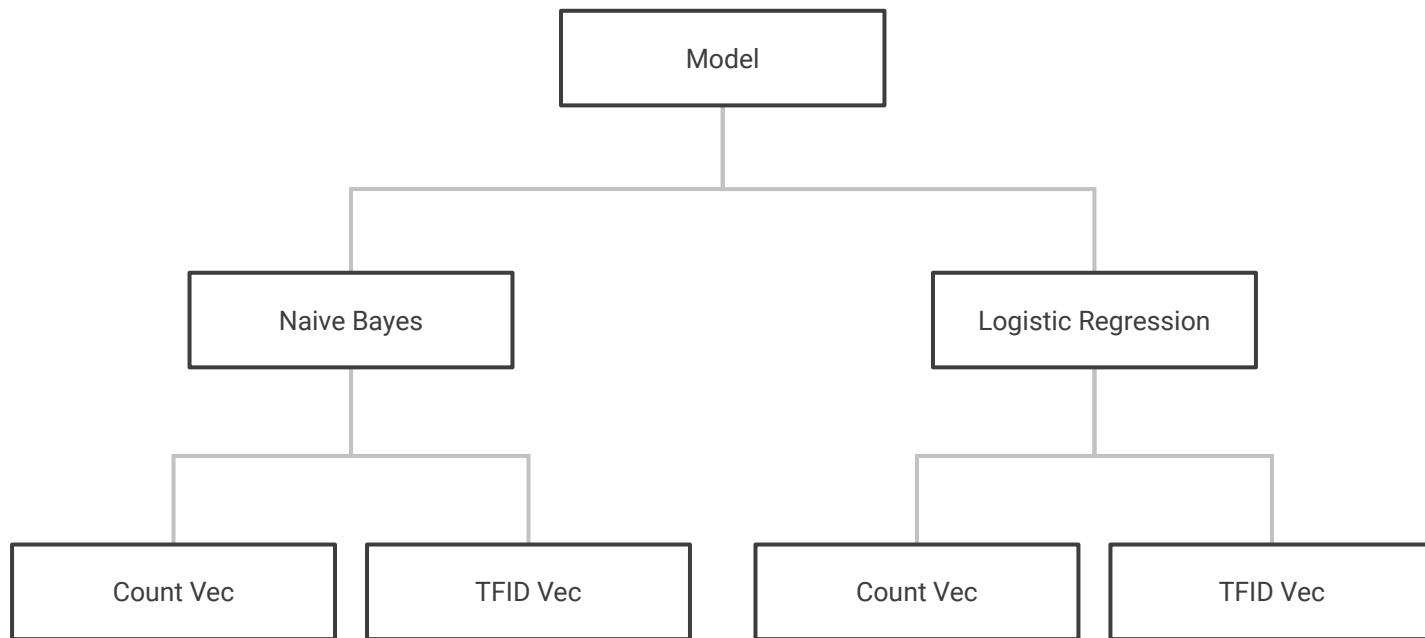X = comments, y = subreddit

# Explore Data



**Exploratory Data Analysis:**

- "Covid 19" top word in /singapore, "Don know" top word in /unitedkingdom
- Covid related issues dominate both subreddits
- Only 5 out of top 50 phrases similar (2,3 ngram)
- **Average number of words in comments:** 39 for r/singapore and 45 for r/unitedkingdom

# Explore Data



top 50 words in r/singapore comments

| word | count |
|------|-------|
| circuit breaker | |
| don think | |
| source code contribute | |
| source code | |
| code contribute | |
| safe distancing | |
| mental health | |
| community cases | |
| really really | |
| advertisement advertisement | |
| ont router | |
| year old | |
| dont know | |
| social distancing | |
| savings account | |
| press conference | |
| don know | |
| short term | |
| abdul halim | |
| good luck | |
| outside dorms | |
| breaker period | |
| enforcement officer | |
| fixed deposit | |
| thanks reply | |
| 10 000 | |
| circuit breaker period | |
| jumpstart account | |
| financial literacy | |
| mentally ill | |
| declared place | |
| coronavirus outbreak | |
| progress report | |
| distancing enforcement officer | |
| years ago | |
| general knowledge | |
| wisdom teeth | |
| distancing enforcement | |
| cabin baggage | |
| distancing measures | |
| place residence | |
| sc jumpstart | |
| fibre broadband | |
| make sure | |
| ll look | |
| chinatown point | |
| status person | |
| dual wan | |
| sorry loss | |
| feel like | |

top 50 words in r/unitedkingdom comments

| word | count |
|------|-------|
| don know | |
| common sense | |
| social distancing | |
| don think | |
| herd immunity | |
| furlough scheme | |
| ve got | |
| nuclear power | |
| year old | |
| don like | |
| life experience | |
| don want | |
| isn going | |
| case scenario | |
| year olds | |
| self employed | |
| spend money | |
| income tax | |
| stay alert | |
| boris johnson | |
| secondary school | |
| poor people | |
| better job | |
| going work | |
| triple lock | |
| socially distance | |
| think people | |
| good job | |
| fossil fuels | |
| don need | |
| teaching assistants | |
| universal credit | |
| worst case | |
| long term | |
| think going | |
| doesn make | |
| mr sunak | |
| just like | |
| far right | |
| 10 years | |
| stay home | |
| public sector | |
| alert rare | |
| uk children | |
| number capita | |
| people tested | |
| kids school | |
| class sizes | |
| fake news | |
| alert rare syndrome | |

# Model with Data

# Modeling with Data

| Model | Training Score | Test Score |
|---|---|---|
| CVEC + Naive Bayes | 0.99 | 0.68 |
| TFID + Naive Bayes | 0.99 | 0.69 |
| CVEC + Logistic Regression | 0.99 | 0.76 |
| TFID + Logistic Regression | 1 | 0.80 |

# Model Evaluation

|  | Predicted UK | Predicted Singapore |
|---|---|---|
| **Actual UK** | 147 | 30 |
| **Actual Singapore** | 36 | 129 |

```
Accuracy: 0.8070175438596491
              precision    recall  f1-score   support

           0       0.80      0.83      0.82       177
           1       0.81      0.78      0.80       165

    accuracy                           0.81       342
   macro avg       0.81      0.81      0.81       342
weighted avg       0.81      0.81      0.81       342
```

**Conclusion:**

- Despite limitations, model has predictive value

- both subreddits have a lot of covid language even after stop words

- differences stem mainly from current affairs

- larger corpus which incorporates larger vocab could be useful