# Classification

## of the political Reddit-verse

""Politicians are the same all over.  They promise to build a bridge even where there is no river."

—Nikita Khruschev

Google

#PickPete

**Bernie**
Politician since 1991.
Senator from Vermont since '07
Longest serving Independent
Running for president in 2020

**Pete**
Politician.
Mayor of South Bend, In
Former military officer
Running for president in 2020

**Liz**
Politician & Academic.
Senator from Massachussets
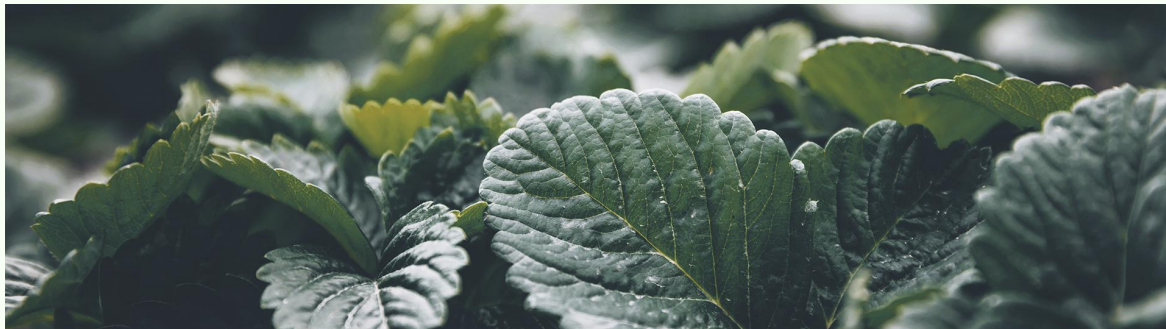Former law school professor
Running for president in 2020

**Kamala**
Lawyer & Politician.
Senator from California
Former district attorney
Running for president in 2020

Says

# Determine the Baseline

## 01

Slightly Unbalanced Class →

| PETE | 26% |
|---|---|
| KAMALA | 26% |
| BERNIE | 26% |
| LIZ | 21% |

# Three Classification Models were chosen to test

## RANDOM FOREST

An ensemble model made up of many boosted decision trees, with randomly selected features. Reduces the variance of a single decision tree.

## LOGISTIC REGRESSION

Aims to find the relationship between features, and the probability of a particular outcome

## GRADIENT BOOSTING

Another decision tree based ensemble model, using regression trees instead of decision trees.

**01**

In the end, the simplest model proved to be the best! (but still not so great)

## ... and now for some performance stats!

| | RANDOM FOREST | NAIVE BAYES | LOGISTIC REGRESSION |
|---|---|---|---|
| **ENGLISH STOP WORDS** | Train: 99.1%<br>Test: **85.2%** | Train: 96.1%<br>Test: **79.1%** | Train: 97.5%<br>Test: **87.5%** |
| **SPEED TO RUN** | MEDIUM | FAST | FAST |

## WOOT!!

# TOP 5 WORDS BY CANDIDATE SUBREDDIT

## BUTTIGIEG

he/him
pete
buttigieg
petebuttigieg
mayor
chasten

## HARRIS

her/she
kamala
harris
kamalaharris
barr
women

## SANDERS

he/him
bernie
sanders
bernie sanders
medicare
deadline

## WARREN

her/she
warren
elizabeth
elizabeth warren
plan
liz

# ... and now for some <u>USEFUL</u> stats!

| | RANDOM FOREST | NAIVE BAYES | LOGISTIC REGRESSION |
|---|---|---|---|
| **ENGLISH STOP WORDS** | Train: 99.1%<br>Test: 85.2% | Train: 96.1%<br>Test: 79.1% | Train: 97.5%<br>Test: 87.4% |
| **REMOVED CANDIDATE NAMES** | Train: 97.7%<br>Test: 52.3% | Train: 86.4%<br>Test: 54.7% | Train: 90.6%<br>Test: 54.9% |
| **SPEED TO RUN** | MEDIUM | FAST | FAST |

**Heatmap** of the **Confusion Matrix**

was there enough data?

# Classification Report

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Sanders | 0.57 | 0.64 | 0.60 | 240 |
| Buttigieg | 0.60 | 0.56 | 0.58 | 249 |
| Harris | 0.50 | 0.57 | 0.53 | 257 |
| Warren | 0.54 | 0.39 | 0.45 | 197 |
|  |  |  |  |  |
| accuracy |  |  | 0.55 | 943 |
| macro avg | 0.55 | 0.54 | 0.54 | 943 |
| weighted avg | 0.55 | 0.55 | 0.55 | 943 |

# TOP 5 WORDS BY CANDIDATE SUBREDDIT
## after name removals

### BUTTIGIEG

chasten **(12.2x)**
south bend
police
service
tomorrow
florida
shooting
team
hire

### HARRIS

barr **(7.31x)**
california
women
busing
sen
iowa
2020
senator
trump

### SANDERS

photo **(10x)**
insurance **(5.7x)**
medicare
deadline
donated
health
donation
today
revolution

### WARREN

plan **(14.58x)**
policy **(7.7x)**
theory
rise
leads
capitalism
debt
ask
winning

Google

#PickPete

**Bernie**
Politician since 1991.
Senator from Vermont since '07
Longest serving Independent
Running for president in 2020

**Pete**
Politician.
Mayor of South Bend, In
Former military officer
Running for president in 2020

**Liz**
Politician & Academic.
Senator from Massachussets
Former law school professor
Running for president in 2020

**Kamala**
Lawyer & Politician.
Senator from California
Former district attorney
Running for president in 2020
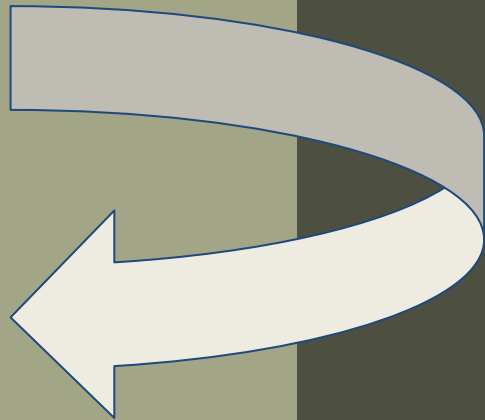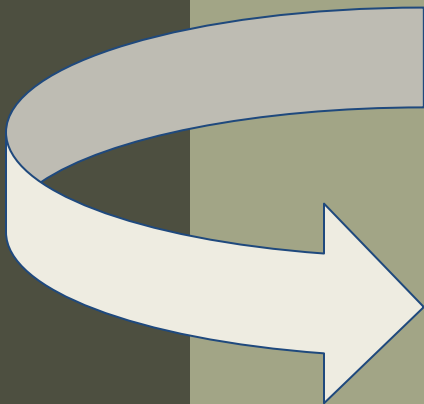
Says

THIS IS JUST THE START!

# Get More Data

and categorized by time posted

# Deeper Analysis

How does each candidate's subreddit correlate to policy subject threads

# Advanced Modeling

Try to generate better predictions by using advance NLP modeling techniques like Spacy, Amazon Comprehend, Google

# CAN YOU SEE THE POSSIBILITIES?!

Does anyone have any questions?

git.generalassmeb.ly/DSTrichter