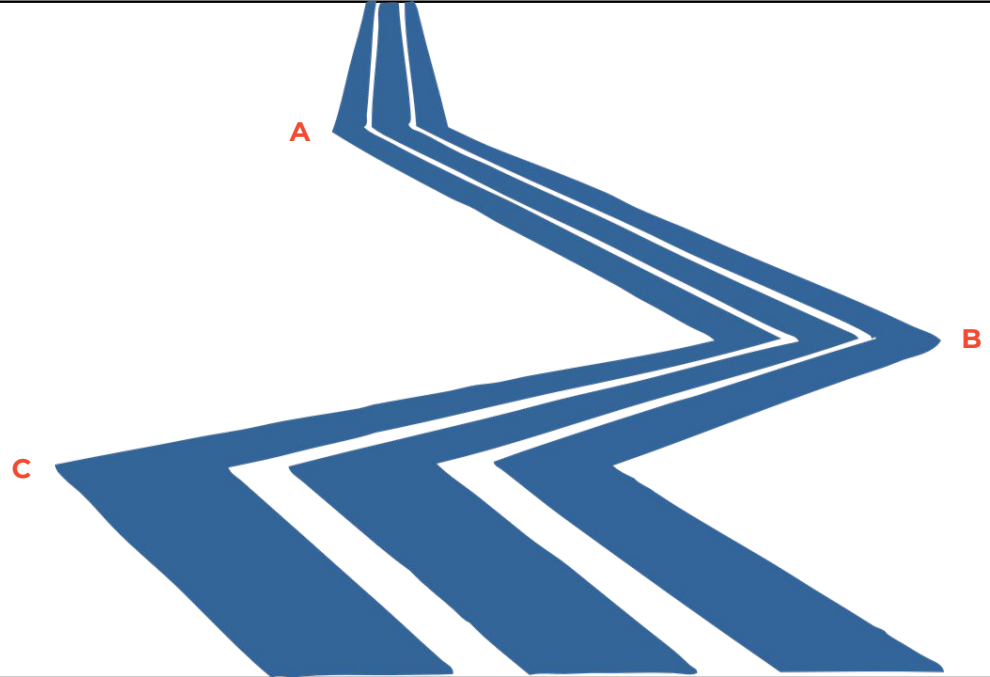# Reddit Scraping:
# Testing NLP for Targeting

Elliot Richardson, Junior Applied Data Scientist

targetsmart

# Outline

- ◎ Introduction
- ◎ Methodology
- ◎ Initial findings
- ◎ Modeling
- ◎ Conclusion
- ◎ Questions

targetsmart

# Introduction

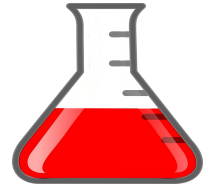◎ **Elliot Richardson,** Jr Applied Data Scientist

◎ **Problem Statement:**

- Digital outreach increasingly important for voter engagement

- Difficult to connect social media users to voter file & therefore scores

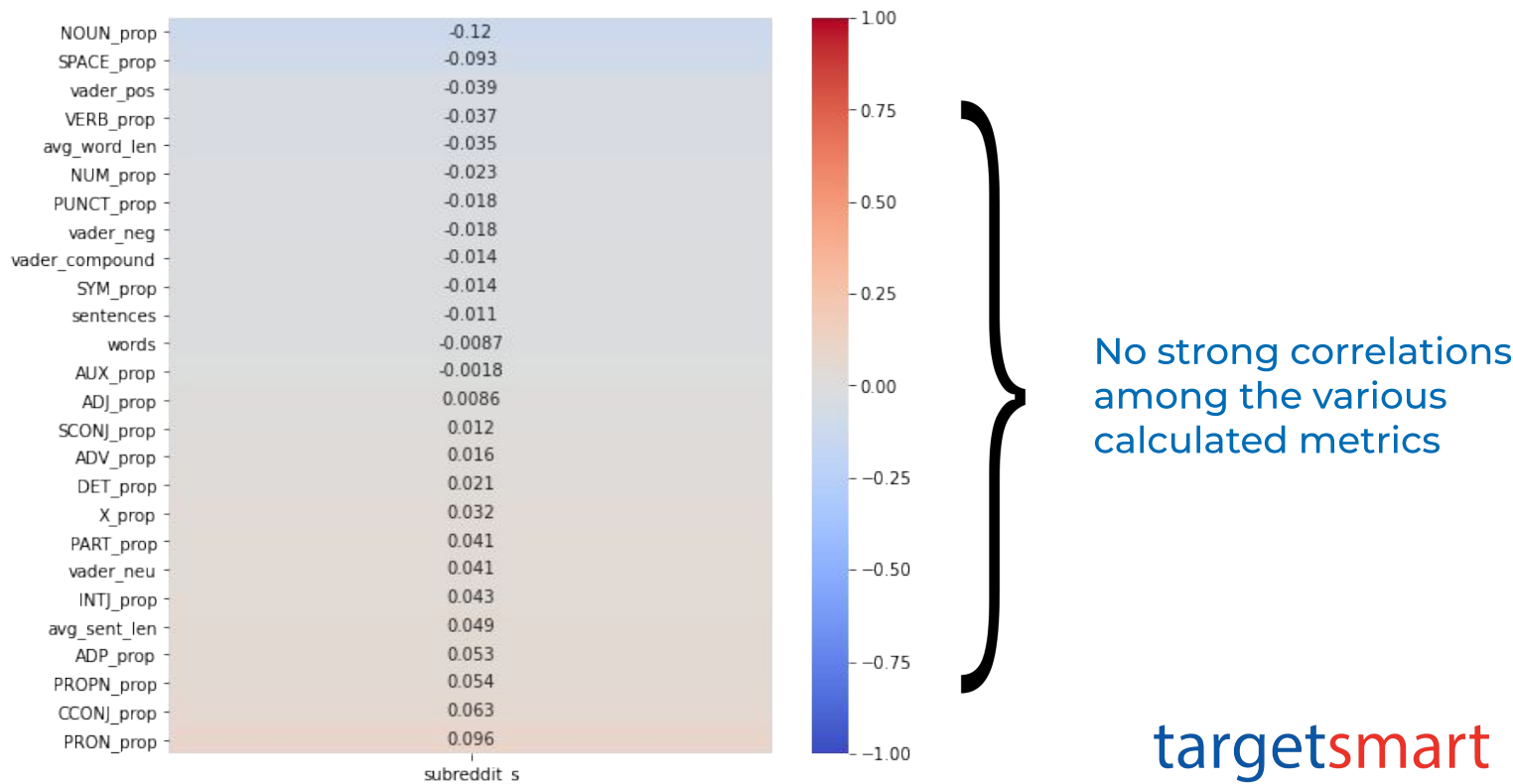- Identifying receptive targets manually is slow and laborious

◎ **Proposed Solution:** Use Reddit to test NLP as online targeting method

targetsmart

# Methodology

◎ Scraped 2500+ posts from r/Socialism and r/Capitalism

◎ Boiled text down to important words

◎ Calculated some other metrics (i.e. sentence length, etc)

◎ Utilized sentiment analysis to assign positive/negative scores

◎ Testing and assembling various models to find reliable patterns

# Initial Finding: Subreddits similar on surface



| | subreddit_s |
|---|---|
| NOUN_prop | -0.12 |
| SPACE_prop | -0.093 |
| vader_pos | -0.039 |
| VERB_prop | -0.037 |
| avg_word_len | -0.035 |
| NUM_prop | -0.023 |
| PUNCT_prop | -0.018 |
| vader_neg | -0.018 |
| vader_compound | -0.014 |
| SYM_prop | -0.014 |
| sentences | -0.011 |
| words | -0.0087 |
| AUX_prop | -0.0018 |
| ADJ_prop | 0.0086 |
| SCONJ_prop | 0.012 |
| ADV_prop | 0.016 |
| DET_prop | 0.021 |
| X_prop | 0.032 |
| PART_prop | 0.041 |
| vader_neu | 0.041 |
| INTJ_prop | 0.043 |
| avg_sent_len | 0.049 |
| ADP_prop | 0.053 |
| PROPN_prop | 0.054 |
| CCONJ_prop | 0.063 |
| PRON_prop | 0.096 |

No strong correlations among the various calculated metrics

targetsmart

# Initial Finding: Subreddits similar on surface



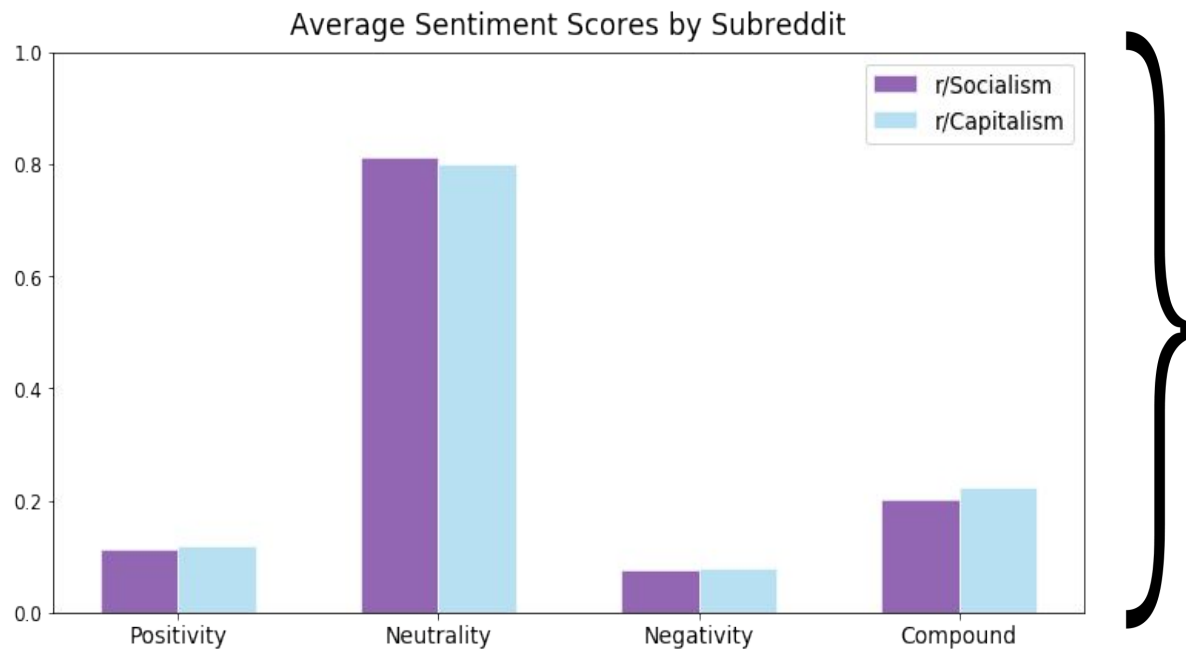Overlap Between Top 100 Most Common Words

27    73    27

Socialism    Capitalism

Lots of overlap between most common words

targetsmart

# Initial Finding: Subreddits similar on surface



Average Sentiment Scores by Subreddit

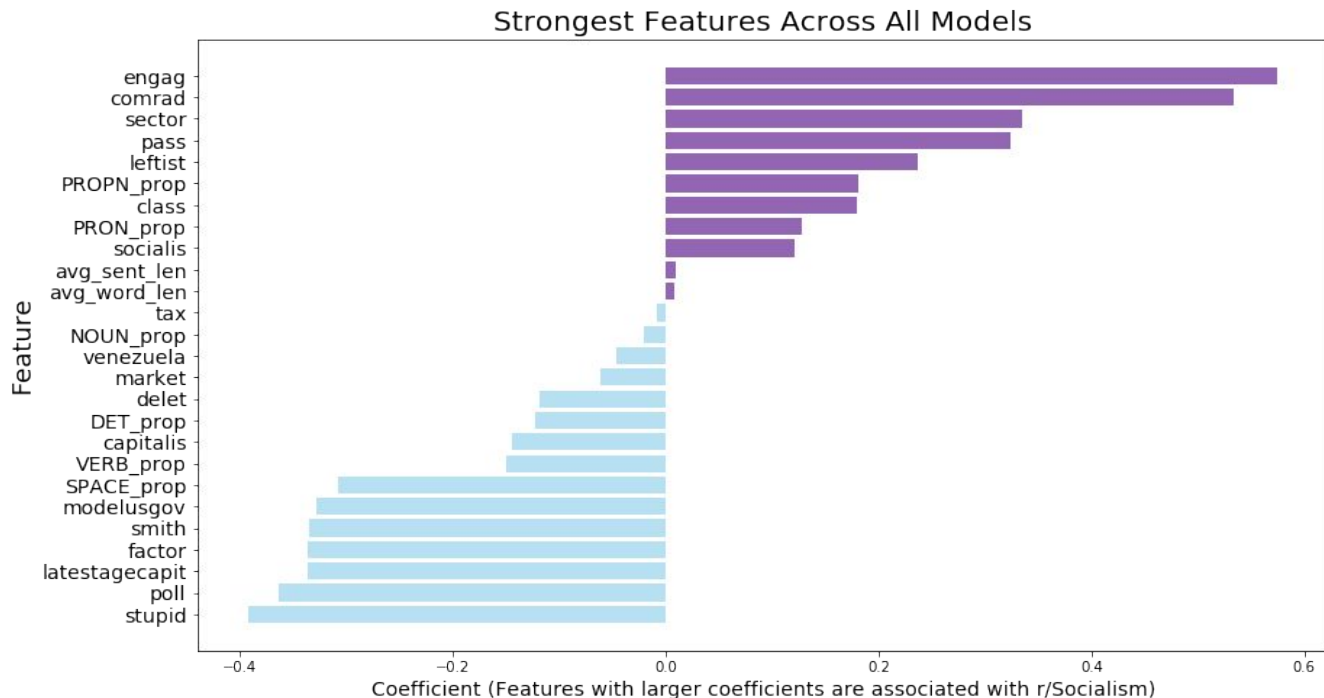Very little difference in sentiment scores

targetsmart

# Modeling

◎ **Ensembled models performed best! Used Voting classifier to ensemble:**

- **Logistic regression**

- **Random forest**

- **Extra trees**

- **Ada boost**

- **Gradient boost**

◎ **77.4% accuracy: outperformed baseline by 51.7%**

$2 > -3$

$0.999... = 1$

$\infty$

$+$ $-$

$\times$ $\div$

$\pi \approx 3.14$

$5^2$

$\sqrt{2}$

$1 + 2 \cdot 3$

$(1 - 2) + 3$

$5(2 + 2)$

$101_2 = 5_{10}$

targetsmart

# Conclusion: Promising but needs more work!



Strongest Features Across All Models

Coefficient (Features with larger coefficients are associated with r/Socialism)

targetsmart