# Flatiron Capstone Project

## Analysis of former reddit the_donald

**Angelo Turri**

# Stakeholder

- A social media communications team for the Republican Party wants me to extract meaningful insights from Reddit data

  - Goal: have a better grasp of the sentiment within their voter base

- Reddit: r/the_donald (no longer exists)

# Data

- Origin: https://the-eye.eu/redarcs/

- Dataset properties

  - Date range: 2015–2020

  - 178,308 unique authors

  - 48 million posts

  - Reduced to 2.1 million due to excessively large dataset size

# Cleaning

- Removing:

  - Bot posts

  - Spam (any post containing numerous consecutive duplicate words or phrases)

  - Stop-words

    - Commonly used words with little meaning that make very small contributions at best to models

Spam example D:

```
In [36]:    1  # Example of removed_posts
            2  removed_posts[removed_posts.joined.str.contains('isis isis isis')].post.iloc[0]
```

```
Out[36]:  'ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISI
          S ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS IS
          IS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS I
          SIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
          ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS ISIS
```

# Feature Engineering (Bag of Words)

- Three levels of analysis

  - Unigrams (single words)

  - Bigrams (two-word phrases)

  - Trigrams (three-word phrases)

- Limited vocabulary to top-100 unigrams, bigrams, or trigrams

- Every sentence including one or more of these words/phrases was used for modeling

# How is a sentence represented?

Vocabulary: 5 words - ['Donald', 'Trump', 'became', 'President', 'apple']

"Donald Trump became President"

| is_donald | is_trump | is_became | is_president | is_apple |
|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 0 |

# All Target Variables

- Post score (upvotes - downvotes)

- VADER sentiment score (engineered)

  - "Valence Aware Dictionary sEntiment Reasoner" 🤦‍♂️

- These are different ways of measuring sentiment, we will choose whichever one allows for a better model and more interpretable results.

# Type of model used: Linear Regression

- Highly interpretable model: it provides the values of coefficients

  - Lets us know which words the model thought negatively impacted sentiment and which ones positively impacted it

- I will create a separate linear regression model for each set of features, so that unigrams, bigrams and trigrams can be analyzed separately

# Target Variable: Score

- Abysmal results across the board

- R-squared < 0.1 for unigrams, bigrams and trigrams models

  - R-squared ranges from 0 to 1

- Models categorized negative words as positive

- Models categorized positive words as negative

- Off by an average of 7-8 upvotes where the average score was 7-8 upvotes

# Target Variable: VADER sentiment

- <span style="color:red">Much better results</span>

- R-squared of 0.243, 0.169, and 0.333 for unigrams, bigrams and trigrams respectively

- The <span style="color:red">broad sentiment</span> of the top positive words and negative words was <span style="color:red">what you'd expect.</span>

- Off on average by 0.33, 0.433, and 0.363 for unigrams, bigrams and trigrams respectively

  - VADER sentiment ranges from -1 to 1

# Top unigrams after tuning

- The longer the bar, the more emotionally charged that word/phrase

- Positive column: Best, God, Trump

- Negative column: News, media, Hillary

# Top bigrams after tuning

- The longer the bar, the more emotionally charged that word/phrase

- Positive column: God bless, America great, thank God, god emperor, United States, free speech

- Negative column: civil war, fake news, Bill Clinton, anti-Trump, Middle East, Deep State



Top positive words

_god_bless
_good_luck
_america_great
_would_love
_pretty_sure
_supreme_court
_thank_god
_high_energy
_pretty_much
_god_emperor
_good_thing
_would_like
_seth_rich
_make_sure
_seems_like
_looks_like
_free_speech
_sounds_like
_united_states
_something_like

0.0  0.1  0.2  0.3  0.4  0.5  0.6

Top negative words

_piece_shit
_civil_war
_holy_shit
_fake_news
_bill_clinton
_shit_like
_anti_trump
_white_people
_black_people
_9_11
_middle_east
_every_single
_going_get
_people_get
_right_wing
_year_old
_people_need
_deep_state
_climate_change
_even_know

−0.6  −0.5  −0.4  −0.3  −0.2  −0.1  0.0

# Top trigrams after tuning

- The longer the bar, the more emotionally charged that word/phrase

- Positive column: God bless America, make America great, God emperor Trump

- Negative column: Bill Clinton rapist, orange man bad, CNN fake news, liberalism mental disorder



Top positive words

| | |
|---|---|
| _super_male_vitality | |
| _ha_ha_ha | |
| _nobel_peace_prize | |
| _god_bless_america | |
| _supreme_court_justice | |
| _make_america_great | |
| _making_america_great | |
| _keep_good_work | |
| _would_love_see | |
| _feels_good_man | |
| _name_seth_rich | |
| 🔲🔲🔲 | |
| _god_emperor_trump | |
| _could_care_less | |
| _clinton_rapist_infowars | |
| _would_like_know | |
| _take_high_energy | |
| _would_like_see | |
| _energy_🔲🔲 | |
| _makes_perfect_sense | |

Top negative words

| | |
|---|---|
| _bill_clinton_rapist | |
| _holy_fucking_shit | |
| _fuck_u_spez | |
| _feel_like_going | |
| _orange_man_bad | |
| _cnn_fake_news | |
| _fake_news_media | |
| _get_shit_together | |
| _us_look_bad | |
| _liberalism_mental_disorder | |
| _blah_blah_blah | |
| _rapist_infowars_com | |
| _law_abiding_citizens | |
| _high_energy_🔲 | |
| _moral_high_ground | |
| _12_year_old | |
| _voter_id_laws | |
| _3_2_1 | |
| _year_old_girl | |
| _black_lives_matter | |

# Insight #1

- The base is passionate about their representative

  - "Best" is listed in the most positive words

  - The word "Trump" is listed in the most positive words

  - The phrase "god emperor Trump" is listed in the most positive trigrams

  - Phrases insulting their representative such as "orange man bad" and "anti trump" are listed in the most negative trigrams

# Insight #2

- The base is passionately against opposing politicians and organizations

  - "News", "media", "fake news", "fake news media", "CNN fake news" and "deep state" appear in the most negative phrases

  - "Hillary", "Bill Clinton," "Bill Clinton rapist", "left", and "liberalism mental disorder" appear in the most negative phrases

# Insight #3

- The base is highly religious

  - "God", "God bless," "thank God", and "God bless America" appear in the most positive phrases.

# Insight #4

- The base is very <span style="color:red">passionate about the relevant issues</span>

    - Swear words such as the f-word and "shit" commonly appear in both the most positive and negative phrases column.

# Insight #5

- The base is very fond of their country

  - Country-oriented phrases such as "America great", "United States", "God bless America," "make America great," and "making America great" appear in the positive column.

  - Phrases such as "civil war" and "middle east" appear in the negative column.

# Recommendations

- Try to match the base's enthusiasm with your own

- Be aware that they are angry against certain people and organizations and want immediate change

# Thank You!

- All questions are welcome.

- You may contact me here:

  - angelo.turri@gmail.com