# Topic Modeling of Tweets: #ElectionNight #MidTerms 2018

HIRANYA KRISHNA KUMAR
*16th Nov 2018*

# PIPELINE

**Data Collection and Storage**

Tweepy :Twitter search API
MongoDB: 70,029 tweets

**Preprocessing**

Clean data, remove stop words , stemming, tokenization

**Sentiment Analysis**

Tweet Polarity, Subjectivity

**Count Vectorizer**

Term Frequencies , unigrams, bigrams, trigrams

**Topic Modeling**

LDA-10 topics
Dimensionality Reduction

**Gradient Boosting Regressor**

Predict retweet counts
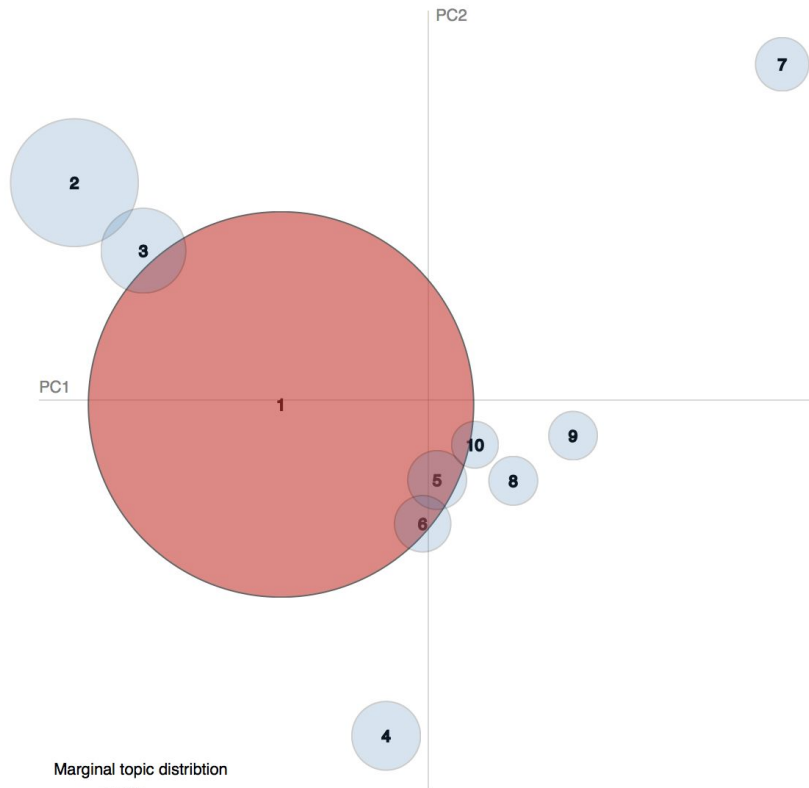Features:
LDA doc_topic probability, polarity,subjectivity
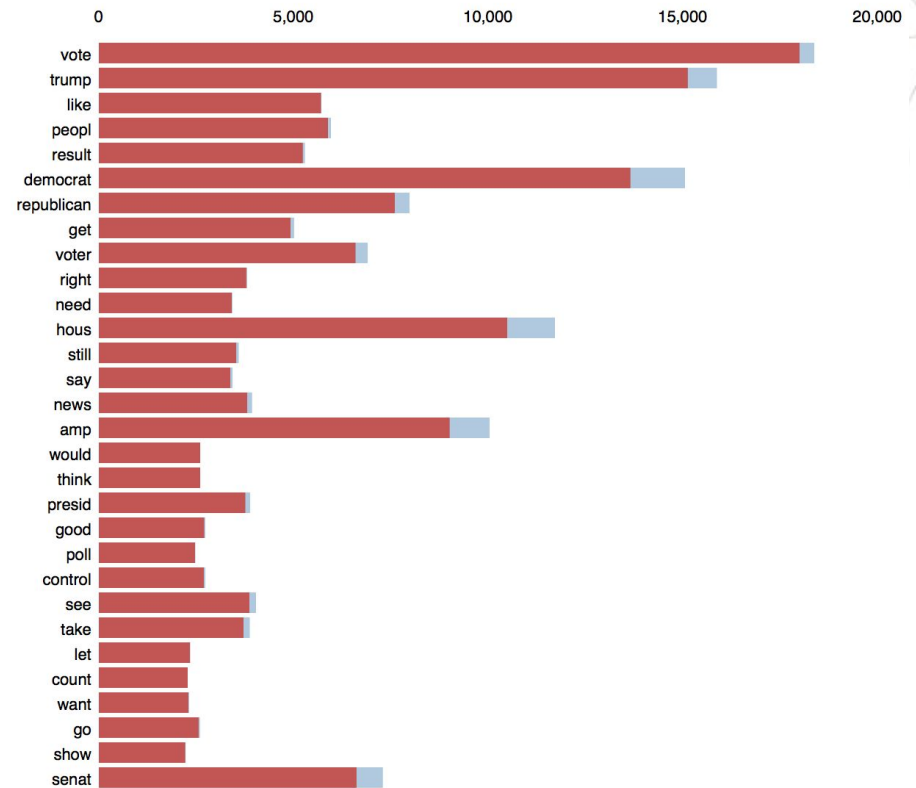
# GENERAL TOPIC -HOUSE, SENATE ELECTIONS

# HISTORIC WINS

Selected Topic: 2 | Previous Topic | Next Topic | Clear Topic

Slide to adjust relevance metric:(2)
λ = 0.09

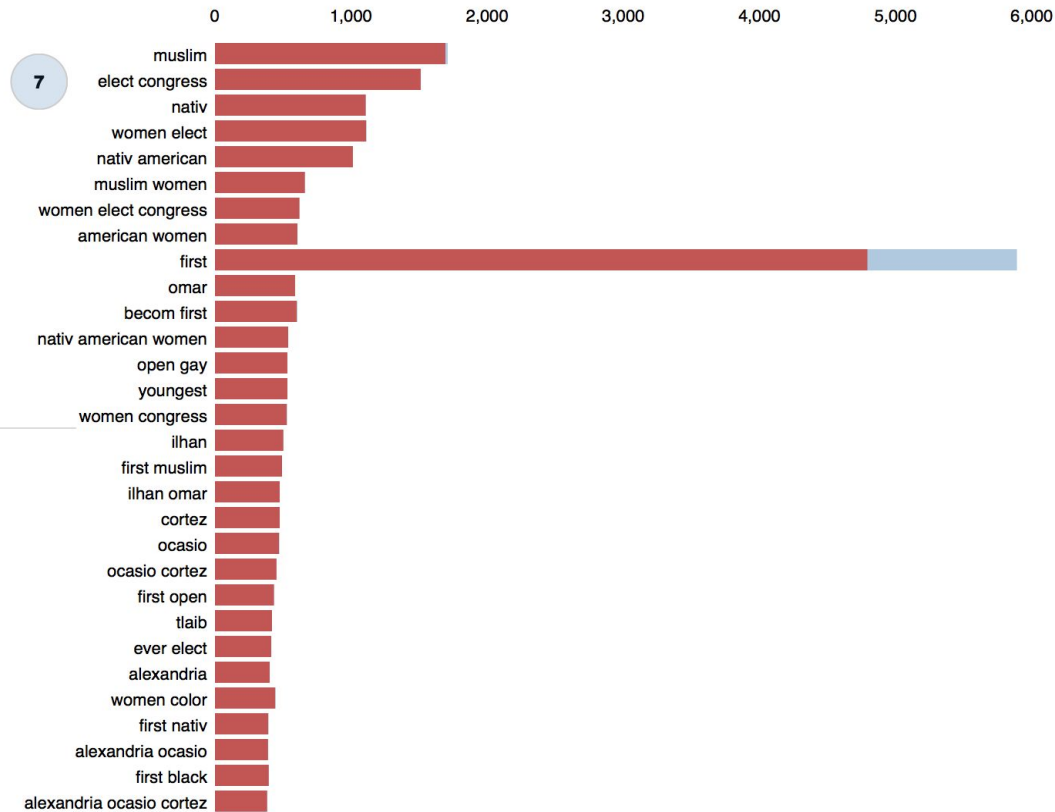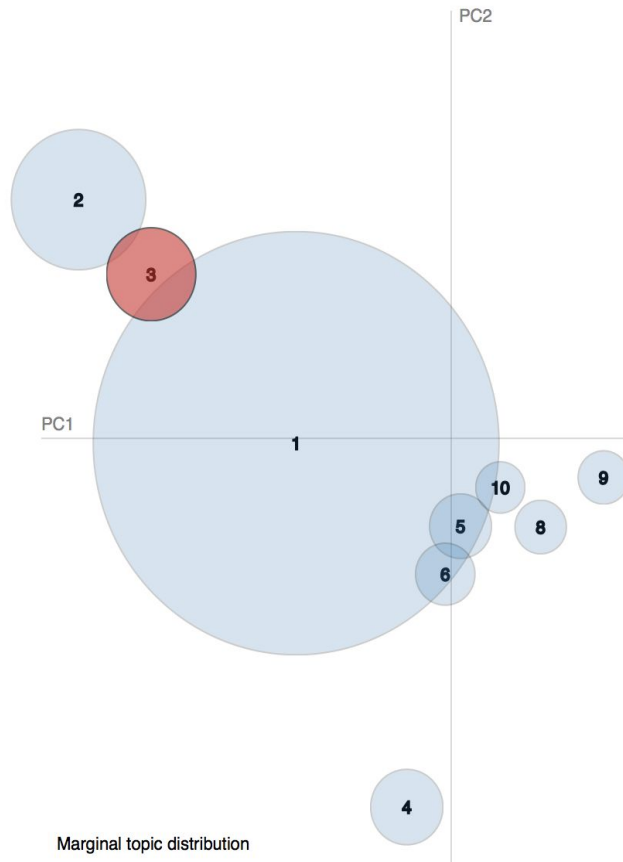0.0   0.2   0.4   0.6   0.8   1

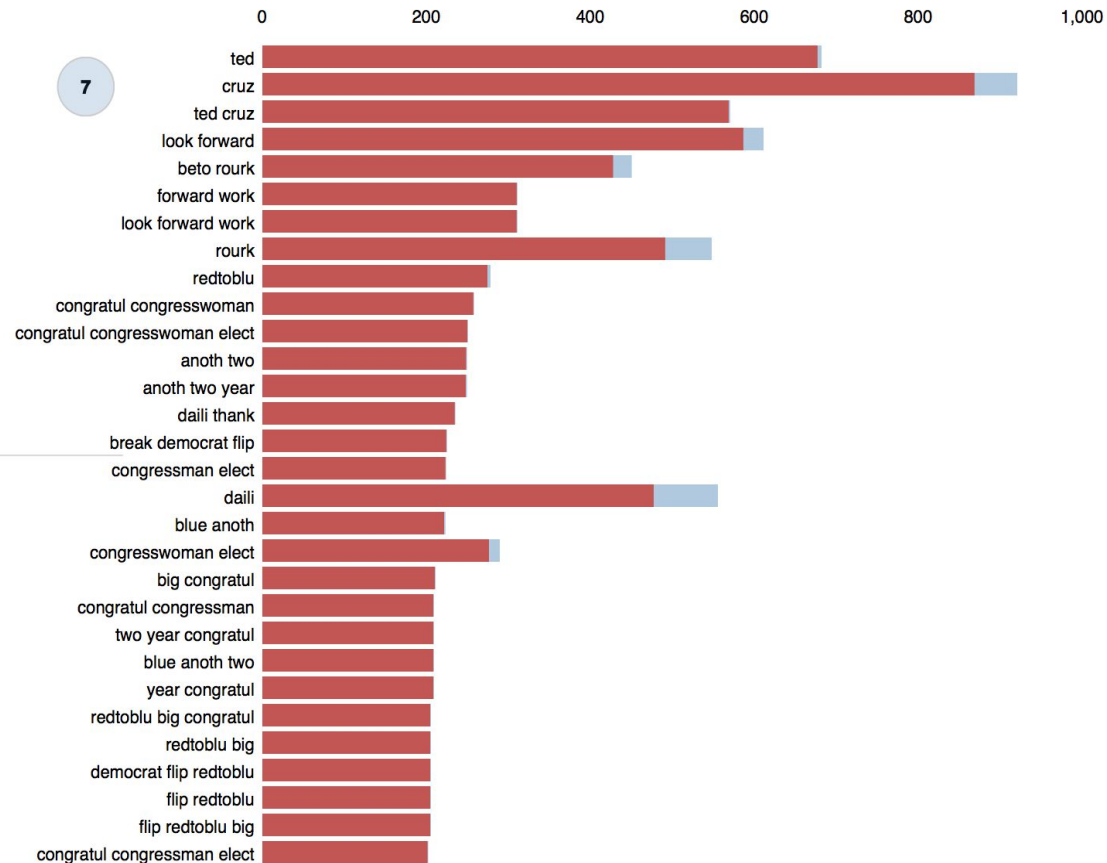## Intertopic Distance Map (via multidimensional scaling)

## Top-30 Most Relevant Terms for Topic 2 (8.4% of tokens)

Marginal topic distribution

2%

5%

10%

Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

# DEMOCRATS FLIPPING HOUSE, TEXAS SENATE
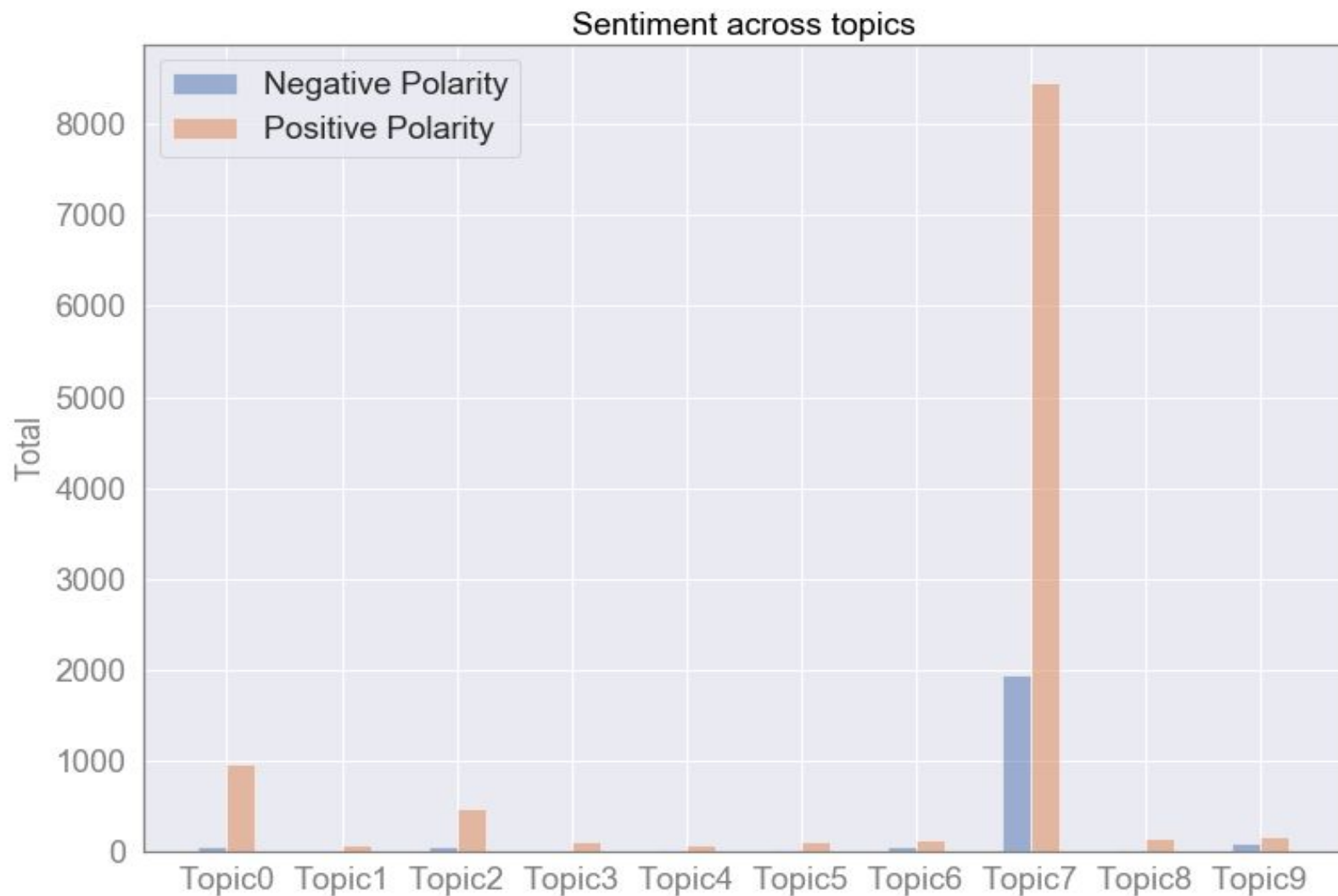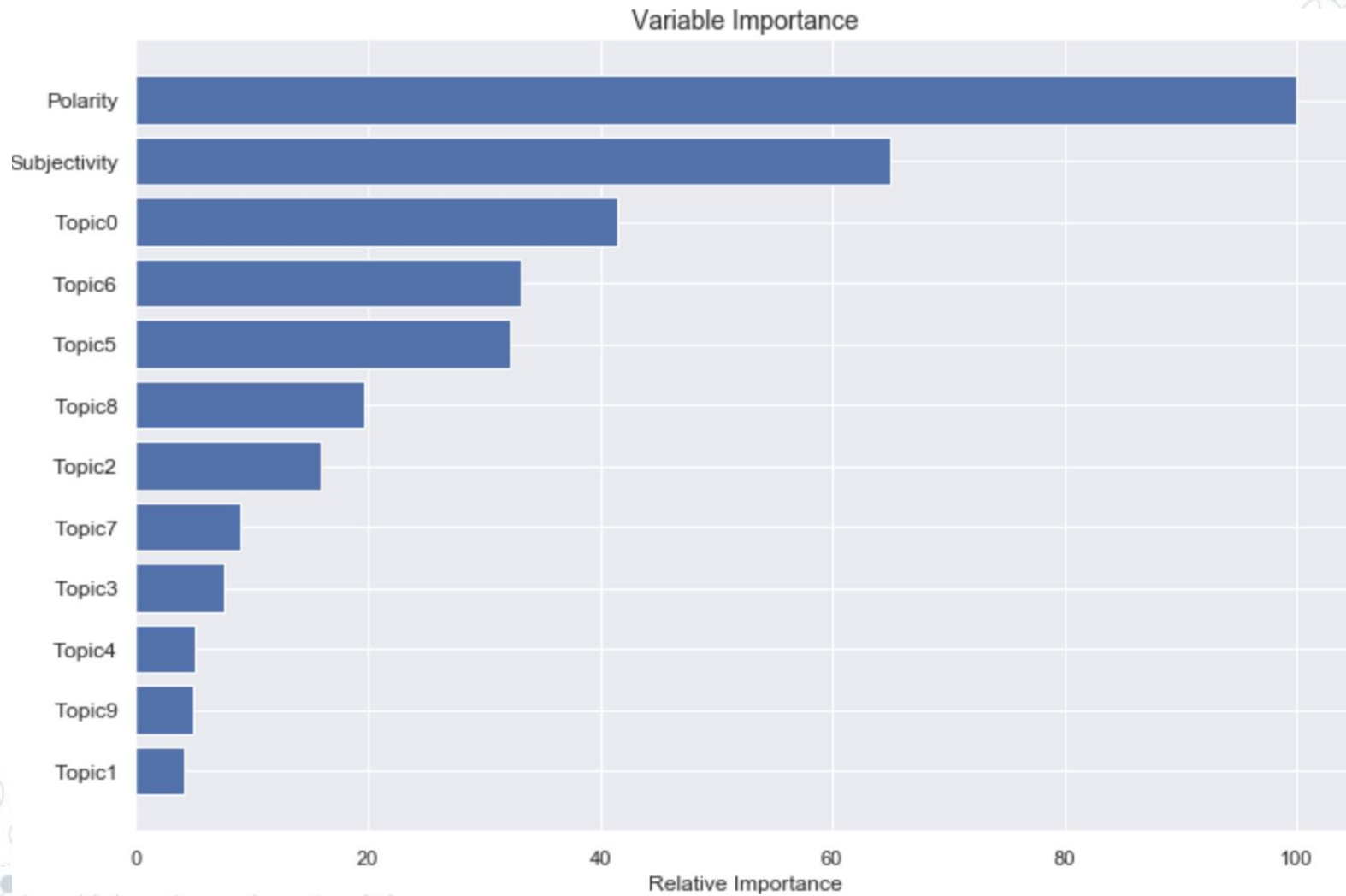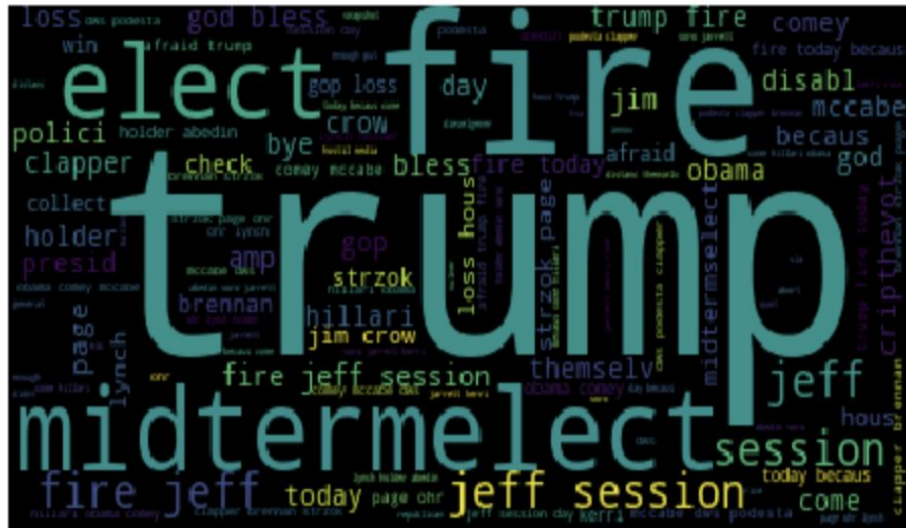
# STATES LEGALIZING MARIJUANA

# GOVERNOR ELECTIONS IN GEORGIA

# SENTIMENT ANALYSIS



Sentiment across topics

# GRADIENT BOOSTING REGRESSOR FEATURE IMPORTANCE TO PREDICT RETWEET COUNTS
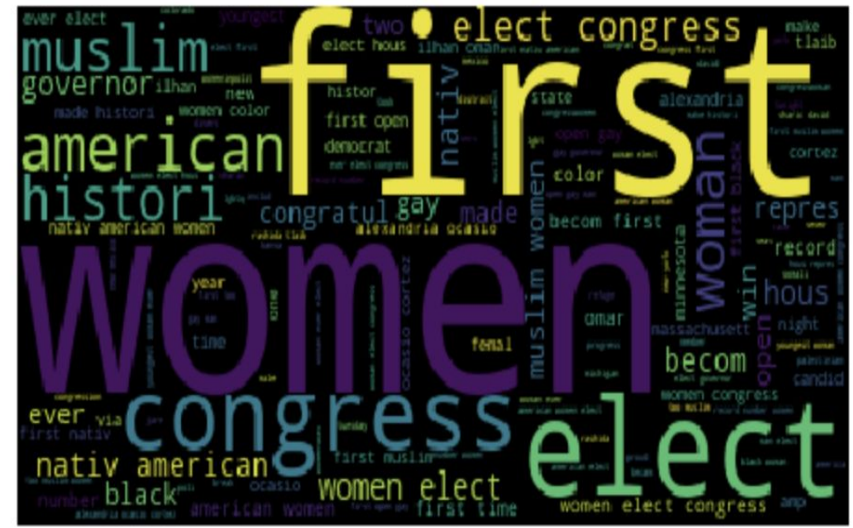


Variable Importance
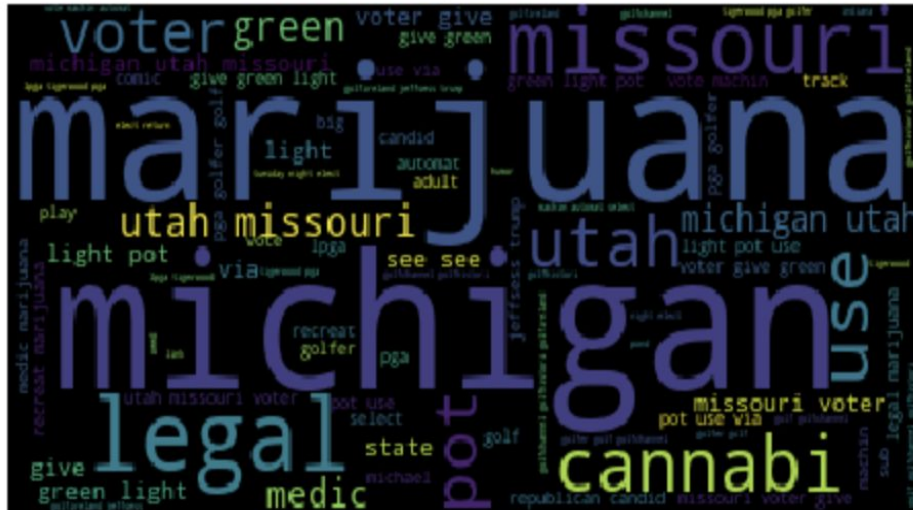
# WORD CLOUD OF TOPICS

TOPIC 6-TRUMP

TOPIC 0-HISTORIC WINS



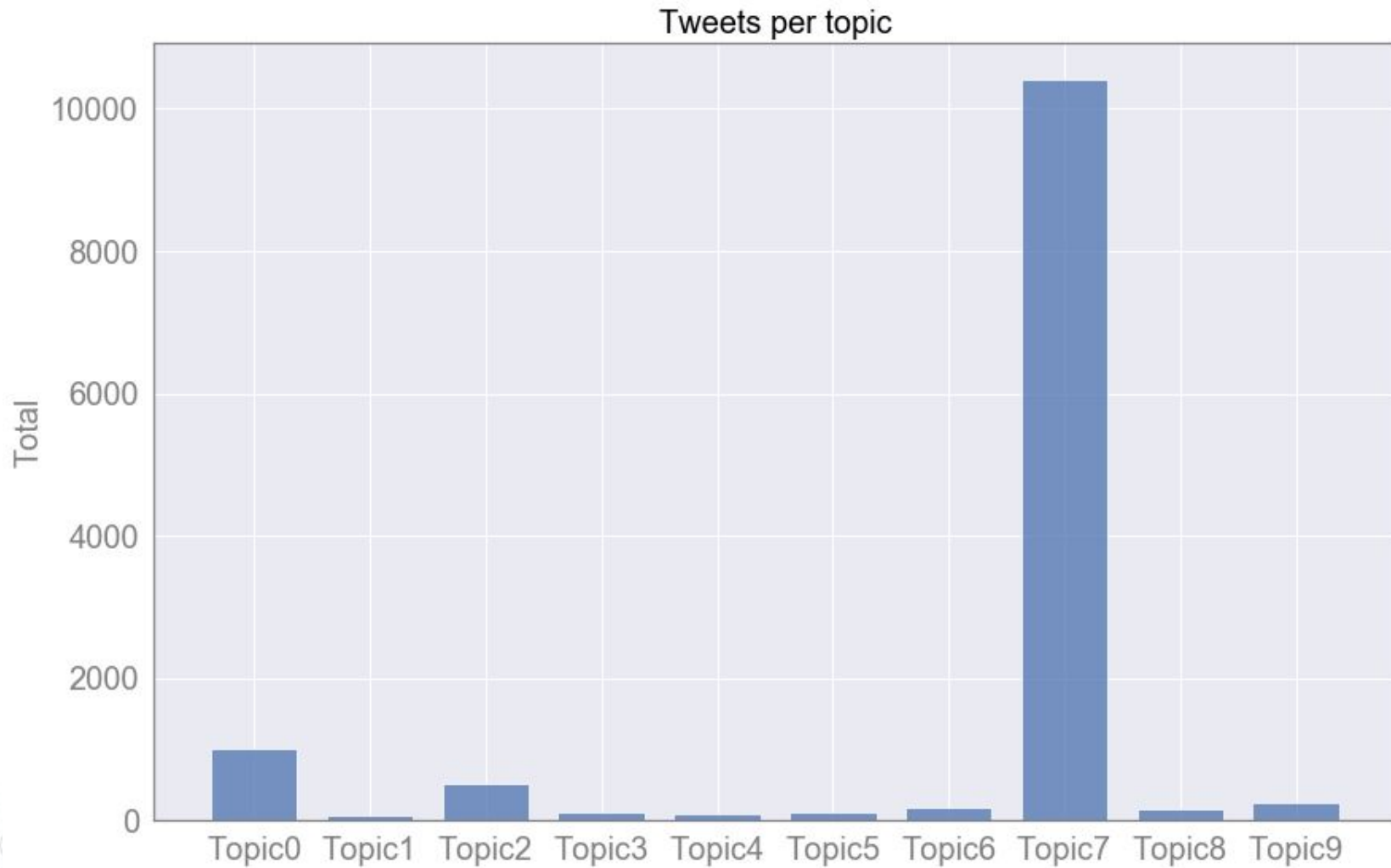TOPIC 5-STATES VOTING TO LEGALIZE MARIJUANA

# FUTURE WORK

1. Use a different pipeline like word2vec and k-means clustering to identify topics and most frequent words in tweets
2. Gather more tweets and use guided LDA to get subtopics within a topic
3. Tune hyperparameters like doc_topic_prior,evaluate_topics frequency,word_topic_prior
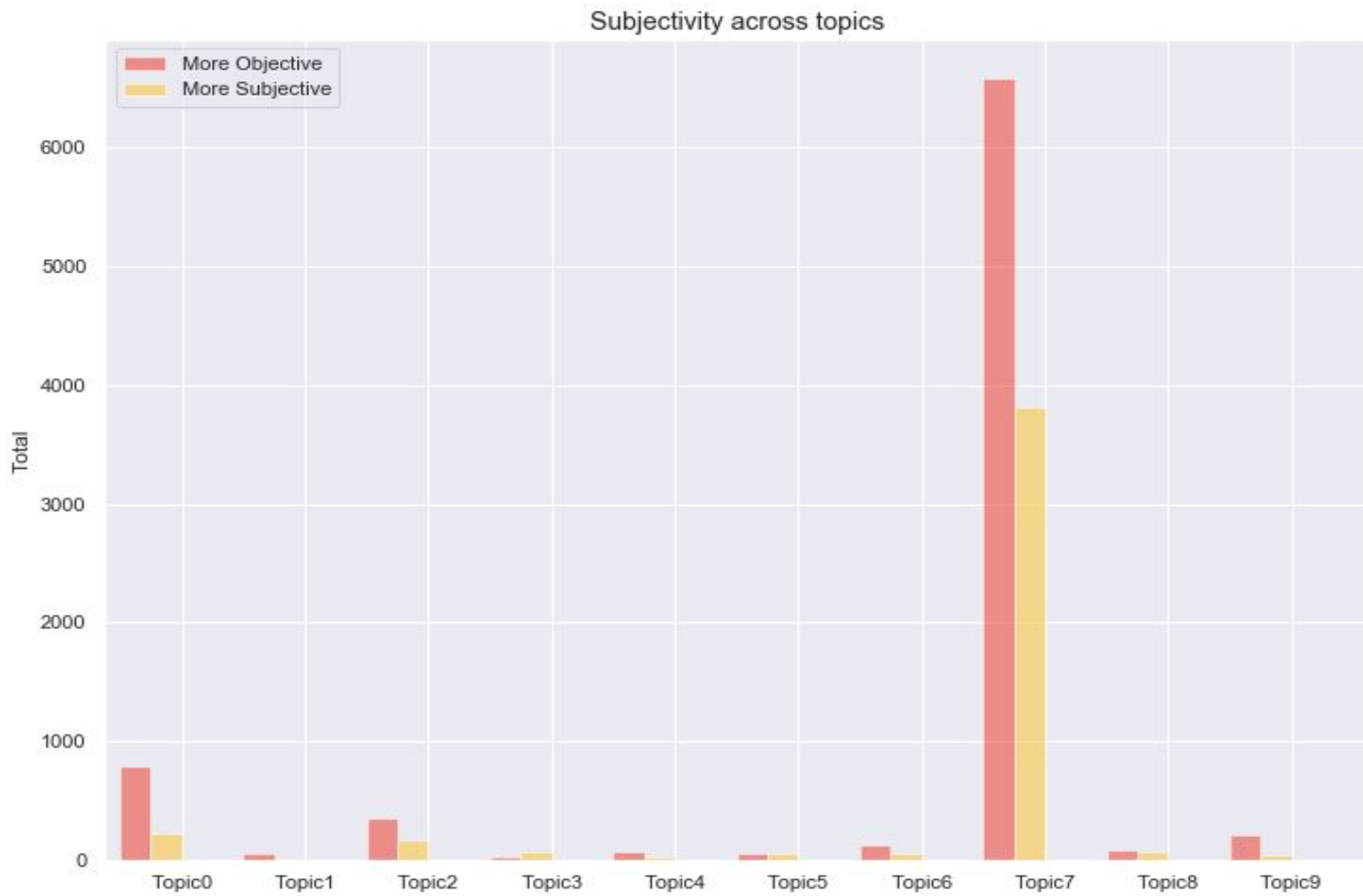4. Build a flask app to predict a topic for a tweet on midterm elections

# Thanks!

## Any questions?

# FUTURE WORK



Tweets per topic

# WORD CLOUD


Topic #8


Topic #2

# SENTIMENT ANALYSIS



Subjectivity across topics

# TOPIC ABOUT ONE PARTY RULE ENDING

# TOPIC ABOUT TRUMP FIRING JEFF SESSIONS