# Trump Tweets Analysis
# Project Report

# Group 13

Chaitanya Gokhale

Vijaya Vasavi Seenivasan

Vamsinadha Reddy Mallavaram

**Problem Statement and Background:**

Donald J. Trump became the 45$^{th}$ President of the United States on January 20, 2017.Since the inauguration, President Trump is actively initiating new policies and conversations, which generate active conversations in the social media. As Data Scientists, we want to take advantage of this opportunity by using text-mining approaches to conduct social media analytics

**Project Overview:**

This project is to analyze people's sentiment and topics about the new administration.We will use Twitter API to collect tweets about Trump. Then we will conduct sentiment analysis to measure how positive or negative the collected tweets are, which can be an indirect measure of Trump's approval. Next, to see what kinds of topics are discussed related to the new president, we will create word clouds and conduct topic modeling on the collected tweets. To see the geographic variation in opinions, we will collect tweets from 5different states and conduct the aforementioned three analyses. Finally, please conclude the project by describing the insights you gained based on the conducted analyses.

**Data Collection:**

As our analysis needs to be done around Donald J. Trump, the initial analysis need tweets that consists/relate anything to *Trump*, which are collected by using Twitter Streaming API's. The package we used to handle the data collection in python is Tweepy and the target keyword as "trump".

Since our experiment design is to also test the geographic variation in the opinions, we have collected tweets from other states i.e., Texas, California, New York, New Jersey and Florida. To achieve this, the tweets were filtered on user location & classified when desired State is found. These results were stored into 5 different files for further analysis.
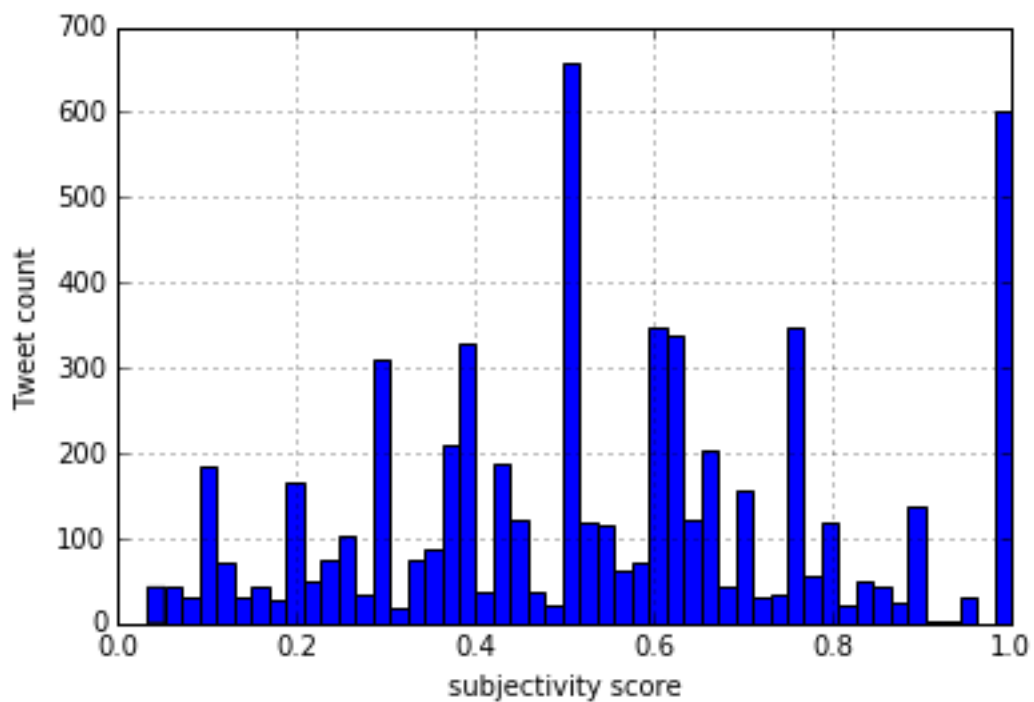
**Data Cleaning:**

One of the main challenges we faced during this project is data cleaning. While twitter allows user to input several special characters along with usual punctuations, fewother characters/Letters refer to certain actionson twitter i.e., @userName to refer a person in the tweet, RT for Re-Tweet. Such characters and strings that add no sense to the analysis needs to be removed from our corpus. Also, we should throw out any web links in our data which are not appropriate for text analysis.

**Note:** Though (') is a punctuation, it is often used as diacritical mark in the text than it is used as a separate special character. And replacing this with a space will make the word incomprehensible. So we have replaced it with nothing to preserve the meaning of the word. i.e., let's => lets, couldn't => couldn't, who's => whos etc.,.
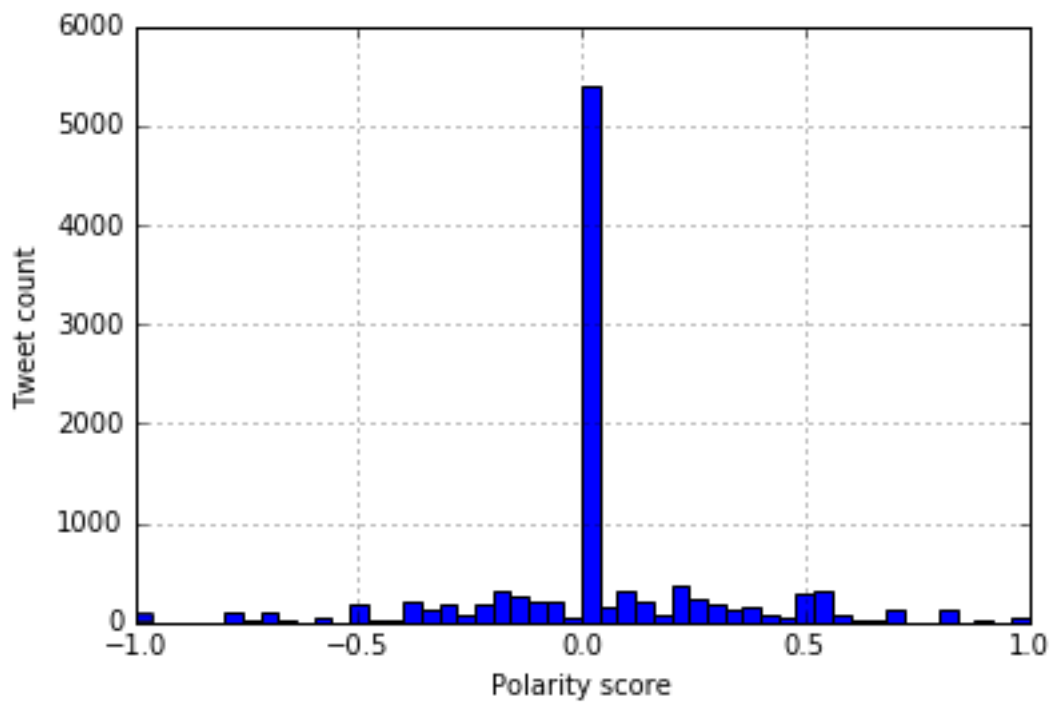
**Sentiment Analysis:**

To measure the sentiment of the people who tweeted about trump, we calculated the subjectivity and polarity scores of the whole corpus & the summarized scores are as displayed below:
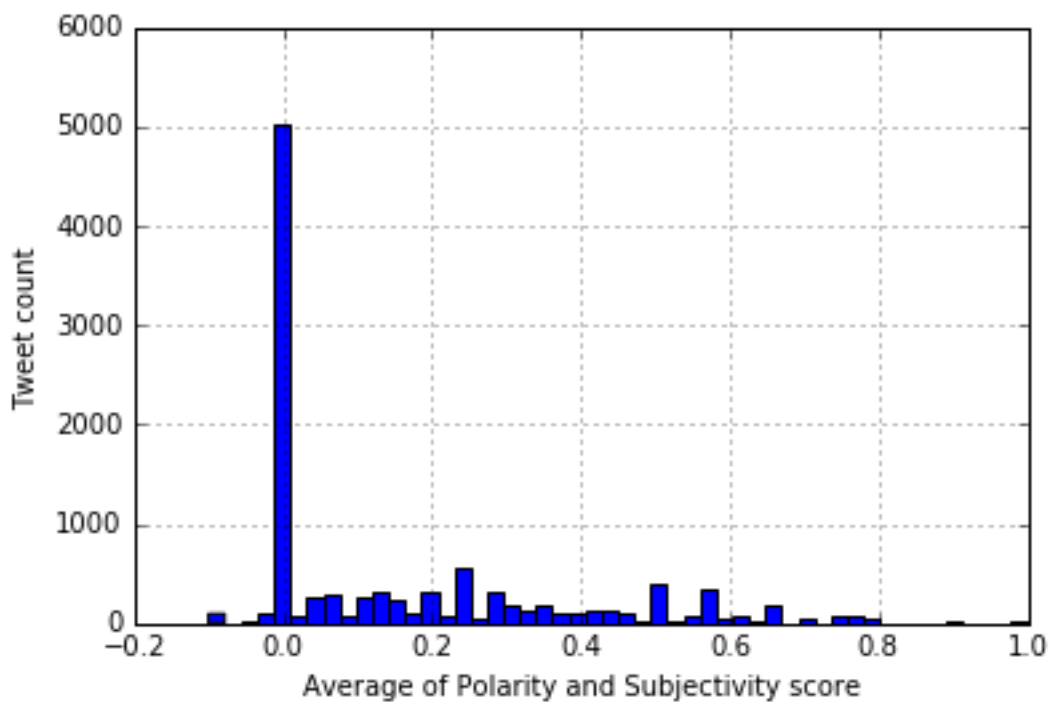
*a.) subjectivity score:*



The subjectivity score distribution is showing that– though good portion of our corpus is not really subjective towards any topic/opinion, most of it have medium to very strong opinions pushing the Avg.weight of subjectivity higher, thus revealing the fact that majority tweets have specific opinion than a general/common opinion.

*b.) Polarity Score:*



        The polarity distribution of the corpus is showing that most of the tweets take a neutral stand, though being subjective.

*c.) Average of Subjectivity & Polarity:*

By plotting the average of subjectivity & polarity we can visualize that most of the users are taking a neutral stand w.r.t their emotion while not making a specific opinion. On the other hand, those who have positive opinions are also specific in expressing their opinions.

**Word Cloud for 10k Tweets:**



From the word cloud, it looks like people are more speaking on the subject: Trump's **claim**, **Obama** had **wiretapped** him during the **campaign.**
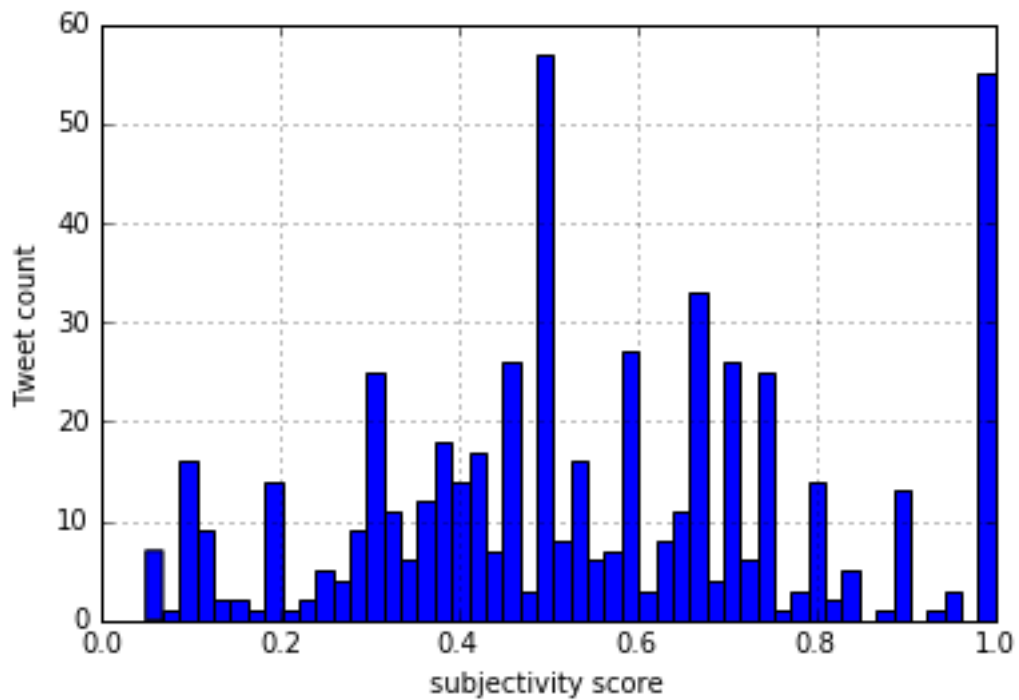
**Topic Modeling:**

a.) LDA:

```
Topic #1 (0, u'0.018*"one" + 0.012*"russia" + 0.009*"another" + 0.008*"racist" + 0.008*"conspiracy" +
0.008*"fmr" + 0.008*"shit" + 0.007*"republicans" + 0.007*"put" + 0.007*"house"')
Topic #2 (1, u'0.033*"obama" + 0.023*"wiretap" + 0.020*"clapper" + 0.018*"tapped" + 0.016*"amp" +
0.016*"claim" + 0.015*"james" + 0.013*"help" + 0.013*"evidence" + 0.013*"much"')
Topic #3 (2, u'0.018*"people" + 0.016*"calls" + 0.015*"wiretap" + 0.014*"things" + 0.013*"see" +
0.012*"news" + 0.012*"obama" + 0.011*"us" + 0.011*"saying" + 0.011*"wont"')
Topic #4 (3, u'0.044*"obama" + 0.028*"de" + 0.025*"right" + 0.021*"barack" + 0.016*"surveillance" +
0.016*"la" + 0.014*"imagine" + 0.013*"former" + 0.013*"potus" + 0.012*"lynch"')
Topic #5 (4, u'0.027*"obama" + 0.024*"russia" + 0.023*"evidence" + 0.019*"team" + 0.016*"amp" +
0.015*"spied" + 0.013*"done" + 0.012*"connections" + 0.011*"ive" + 0.011*"american"')
Topic #6 (5, u'0.031*"media" + 0.013*"obama" + 0.012*"obamas" + 0.012*"steel" + 0.012*"last" +
0.011*"pipeline" + 0.011*"keystone" + 0.010*"case" + 0.010*"political" + 0.010*"russian"')
Topic #7 (6, u'0.024*"staff" + 0.019*"remember" + 0.016*"doesnt" + 0.015*"well" + 0.014*"prison" +
0.014*"house" + 0.013*"went" + 0.013*"white" + 0.012*"dear" + 0.010*"word"')
Topic #8 (7, u'0.050*"obama" + 0.030*"wiretapped" + 0.029*"president" + 0.021*"tower" + 0.020*"order" +
0.020*"wiretapping" + 0.020*"proof" + 0.017*"amp" + 0.016*"levin" + 0.016*"news"')
Topic #9 (8, u'0.036*"says" + 0.028*"gutfeld" + 0.020*"yet" + 0.019*"week" + 0.019*"show" +
0.019*"presidential" + 0.016*"clapper" + 0.016*"fisa" + 0.015*"would" + 0.014*"greg"')
Topic #10 (9, u'0.024*"march" + 0.020*"great" + 0.017*"fisa" + 0.017*"supporters" + 0.016*"mar" +
0.016*"lago" + 0.014*"rally" + 0.013*"request" + 0.013*"reminder" + 0.013*"budget"')
```

*b.) NMF:*

```
Topic 0: roger stone wikileaks admits deletes
Topic 1: house white hes staff know
Topic 2: mar lago cost weekends annual
Topic 3: clapper wiretap james activity mounted
Topic 4: gutfeld presidential week greg says
Topic 5: lynch imagine barack potus loretta
Topic 6: obama wiretapping president evidence claim
Topic 7: tower intel wiretapped dueling chiefs
Topic 8: russia team connections web conspiracy
Topic 9: amp obama tapping friends hacking
```
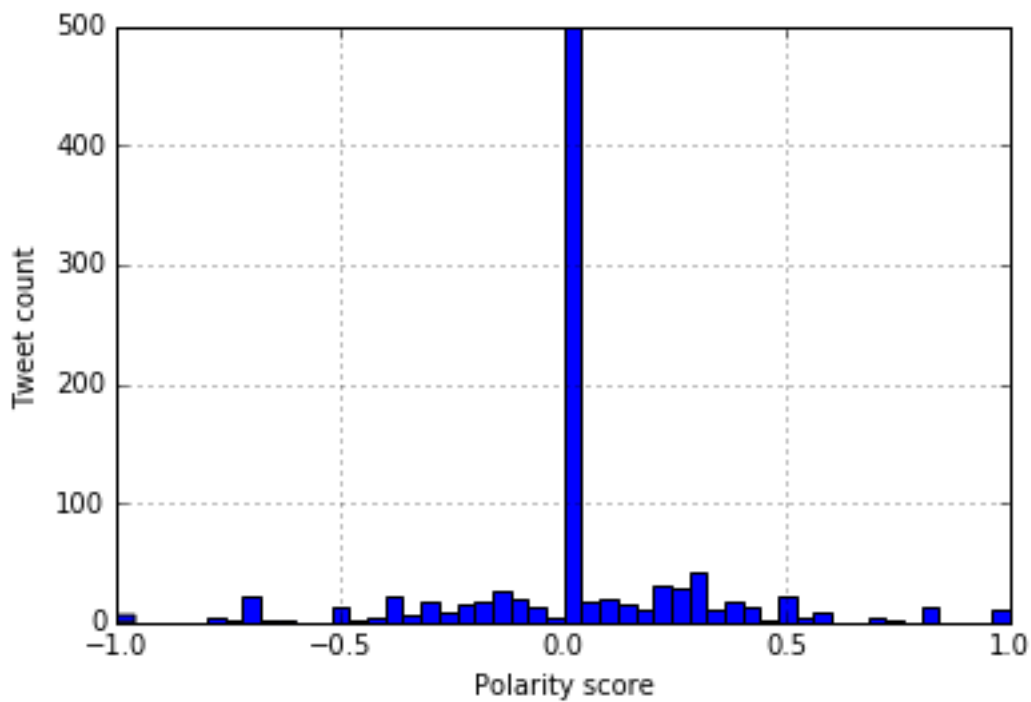
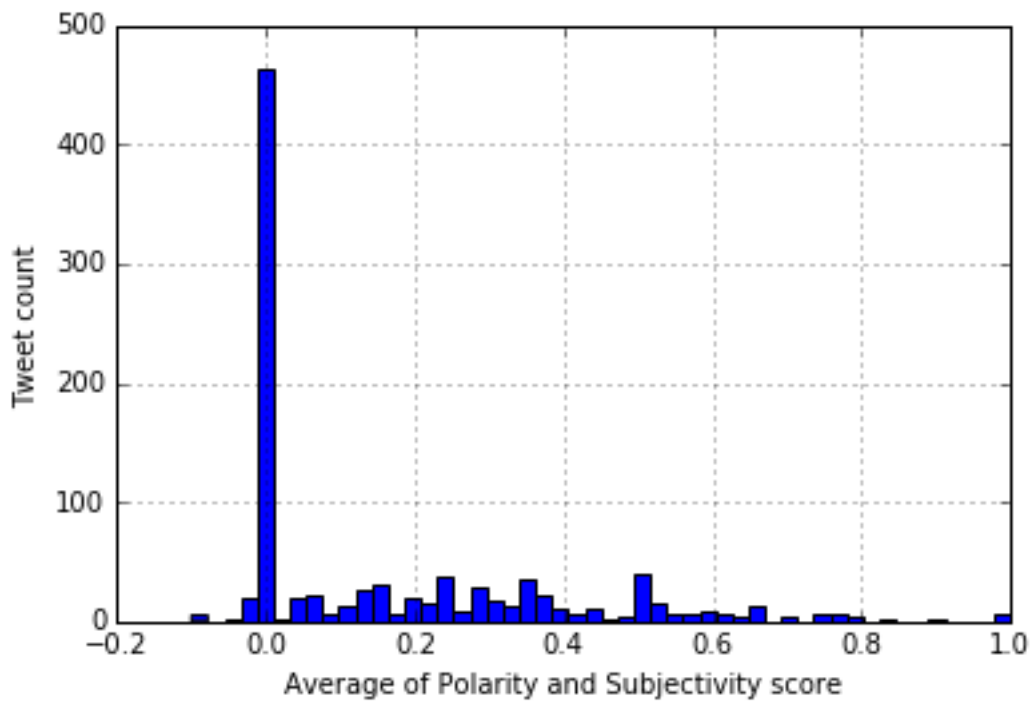# Texas State

*a.) subjectivity score:*



The above plot is showing that we have a fair distribution of subjectivity scores, which states that users from Texas states are making opinions but not very strong opinions.
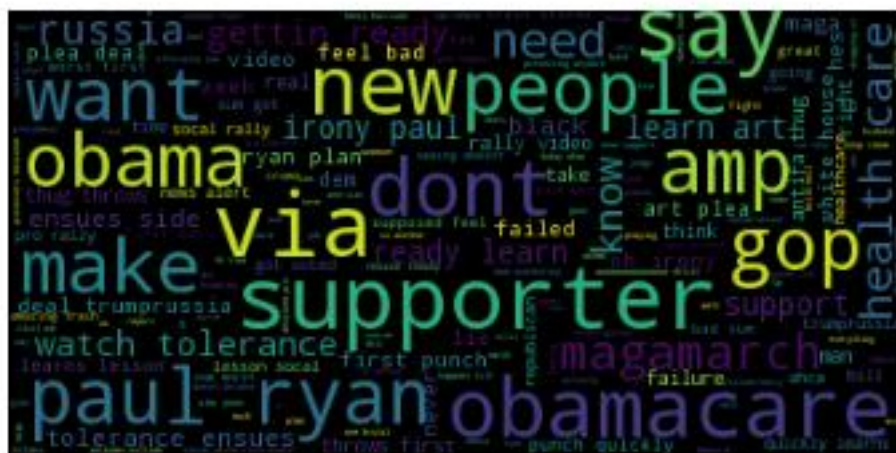
*b.) Polarity Score:*

The polarity distribution of the corpus is showing that most of the tweets take a neutral stand.

*c.) Average of Subjectivity & Polarity:*
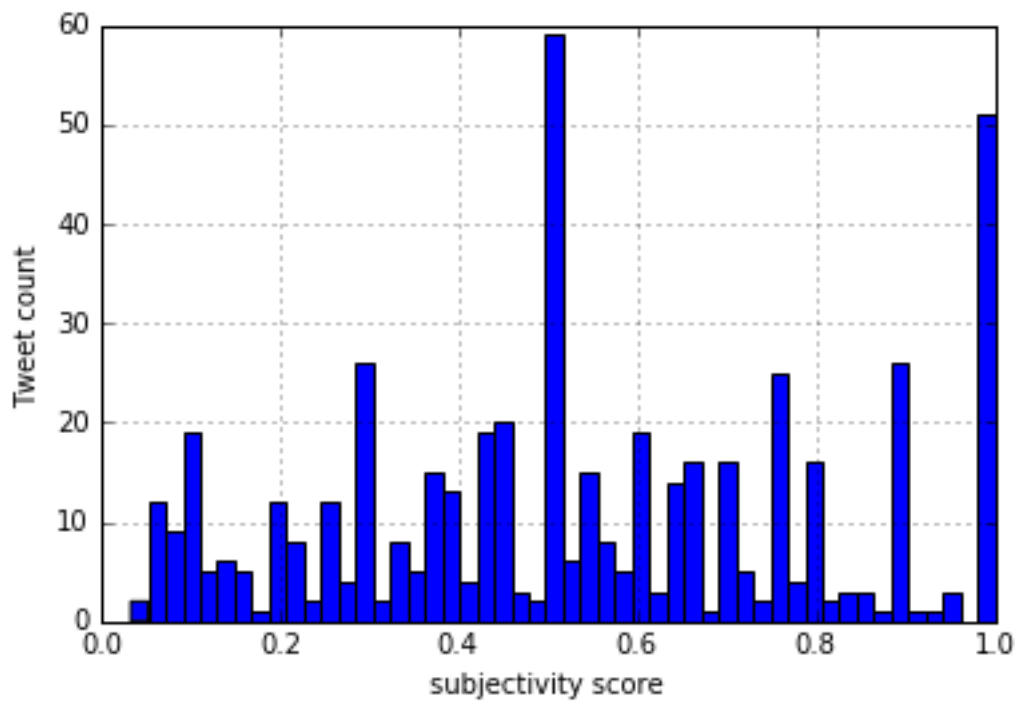


**Word Cloud:**

**Topic Modeling:**

a.) LDA:

```
Topic #1 (0, u'0.022*"like" + 0.020*"ever" + 0.018*"news" + 0.018*"seen" + 0.016*"one" + 0.016*"california" +
0.016*"alert" + 0.016*"ive" + 0.016*"brutal" + 0.016*"quotes"')

Topic #2 (1, u'0.024*"rally" + 0.022*"video" + 0.018*"first" + 0.017*"antifa" + 0.016*"punch" + 0.016*"learns" +
0.016*"quickly" + 0.016*"socal" + 0.016*"lesson" + 0.016*"thug"')

Topic #3 (2, u'0.025*"white" + 0.025*"house" + 0.019*"today" + 0.019*"support" + 0.019*"show" + 0.017*"amazing" +
0.017*"train" + 0.017*"storms" + 0.010*"someone" + 0.009*"people"')

Topic #4 (3, u'0.013*"russiagate" + 0.013*"hes" + 0.011*"supporters" + 0.009*"obama" + 0.009*"report" +
0.009*"order" + 0.008*"antifa" + 0.007*"things" + 0.007*"wide" + 0.007*"open"')

Topic #5 (4, u'0.015*"obamacare" + 0.013*"shiiit" + 0.011*"healthcare" + 0.009*"week" + 0.006*"words" +
0.006*"yes" + 0.006*"news" + 0.006*"hes" + 0.006*"cnn" + 0.006*"plan"')

Topic #6 (5, u'0.016*"magamarch" + 0.014*"supporters" + 0.014*"anti" + 0.012*"getting" + 0.012*"club" +
0.012*"said" + 0.012*"billionaire" + 0.009*"way" + 0.009*"hell" + 0.009*"beat"')

Topic #7 (6, u'0.033*"ryan" + 0.027*"paul" + 0.023*"oh" + 0.019*"plan" + 0.019*"irony" + 0.016*"husband" +
0.014*"president" + 0.013*"golfing" + 0.010*"times" + 0.010*"weeks"')

Topic #8 (7, u'0.022*"another" + 0.021*"one" + 0.019*"side" + 0.019*"watch" + 0.018*"tolerance" + 0.018*"ensues"
+ 0.015*"war" + 0.011*"pro" + 0.011*"becomes" + 0.010*"explode"')

Topic #9 (8, u'0.024*"bill" + 0.016*"democrats" + 0.014*"obama" + 0.013*"happens" + 0.013*"kill" + 0.013*"jones"
+ 0.013*"alex" + 0.012*"failure" + 0.012*"health" + 0.011*"via"')

Topic #10 (9, u'0.018*"magamarch" + 0.018*"first" + 0.017*"president" + 0.016*"days" + 0.016*"worst" + 0.016*"im"
+ 0.016*"like" + 0.015*"doesnt" + 0.015*"harrison" + 0.015*"william"')
```
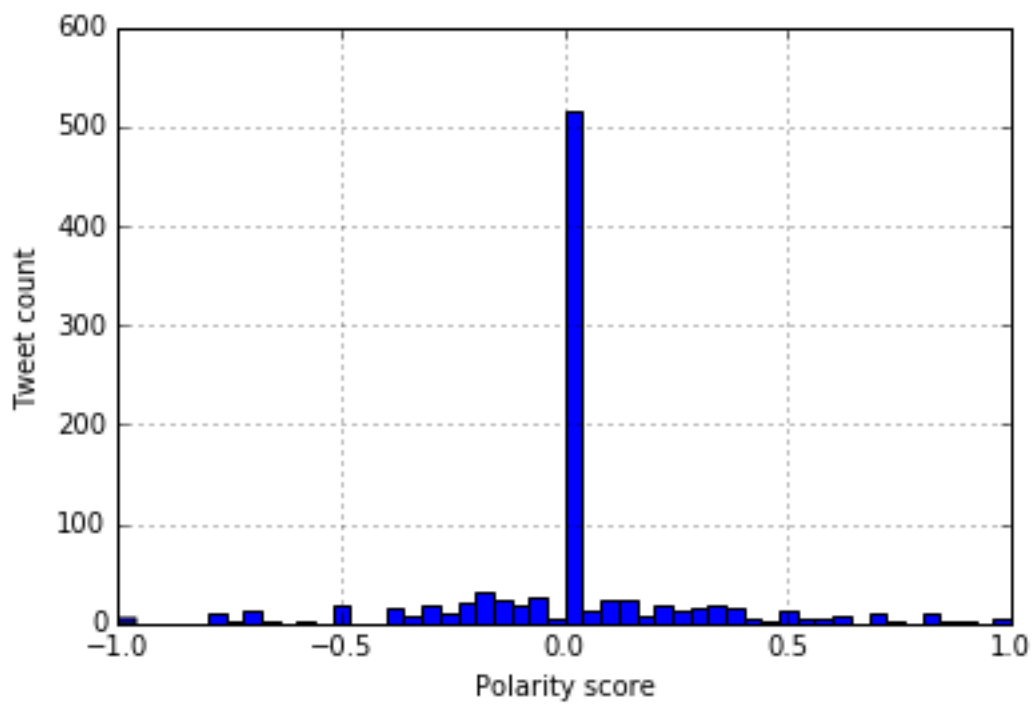
b.) NMF:

```
Topic 0: tolerance ensues watch hollywood rally
Topic 1: rally antifa video lesson thug
Topic 2: deal art gettin plea learn
Topic 3: support today white house train
Topic 4: ryan oh paul irony plan
Topic 5: like doesnt im protesting seeing
Topic 6: president days worst william harrison
Topic 7: magamarch supporters turnout strong maga
Topic 8: happens jones alex kill things
Topic 9: feel voted bad got sum
```
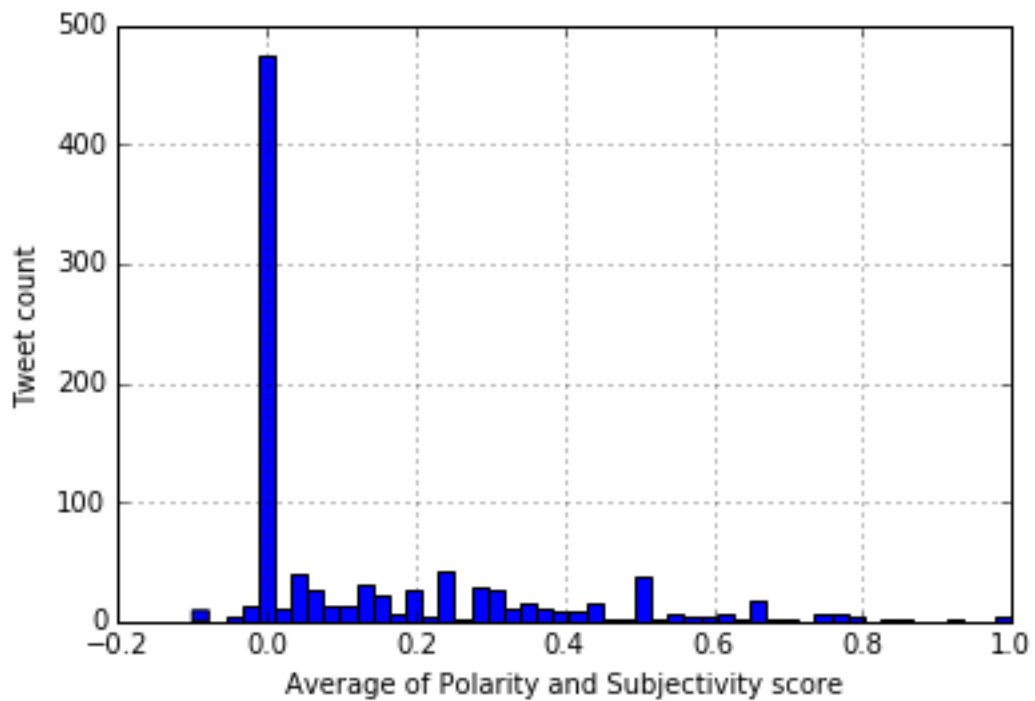
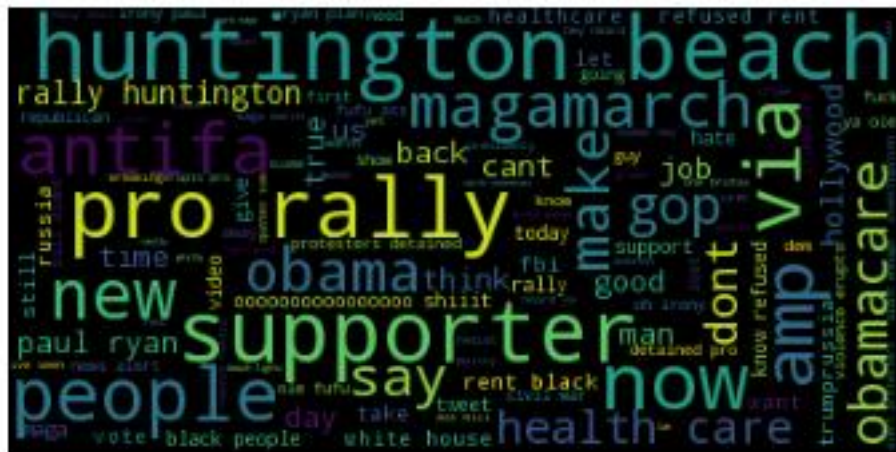# California State

*a.) subjectivity score:*



*b.) Polarity Score:*

*c.) Average of Subjectivity & Polarity:*



**Word Cloud:**



From the word cloud, we can summarize that people tweeted more about the **Trump'sProRally** in **Huntingtonbeach.**

**Topic Modeling:**

a.) LDA:

```
Topic #1 (0, u'0.013*"russiagate" + 0.012*"watch" + 0.009*"amp" + 0.009*"never" + 0.009*"side" + 0.009*"thread" +
0.009*"tolerance" + 0.009*"ensues" + 0.009*"true" + 0.007*"like"')

Topic #2 (1, u'0.013*"rented" + 0.011*"gets" + 0.009*"hes" + 0.009*"thought" + 0.007*"caught" + 0.007*"potus" +
0.007*"husband" + 0.007*"war" + 0.007*"party" + 0.007*"run"')

Topic #3 (2, u'0.024*"rally" + 0.021*"supporters" + 0.019*"pro" + 0.018*"california" + 0.016*"news" +
0.015*"ryan" + 0.015*"oh" + 0.015*"protesters" + 0.014*"plan" + 0.013*"antifa"')

Topic #4 (3, u'0.027*"magamarch" + 0.016*"hollywood" + 0.016*"us" + 0.013*"theres" + 0.011*"like" + 0.009*"says"
+ 0.009*"party" + 0.009*"come" + 0.009*"join" + 0.007*"russia"')

Topic #5 (4, u'0.019*"war" + 0.013*"deal" + 0.013*"gettin" + 0.013*"art" + 0.011*"russia" + 0.011*"congress" +
0.011*"learn" + 0.011*"syria" + 0.011*"ready" + 0.011*"plea"')

Topic #6 (5, u'0.018*"anti" + 0.016*"political" + 0.014*"real" + 0.014*"bill" + 0.013*"russia" + 0.013*"anything"
+ 0.013*"journalists" + 0.013*"nyt" + 0.013*"hacks" + 0.013*"likes"')

Topic #7 (6, u'0.016*"one" + 0.009*"ryan" + 0.009*"another" + 0.007*"art" + 0.007*"yet" + 0.007*"people" +
0.007*"policy" + 0.007*"usa" + 0.007*"paul" + 0.007*"chants"')

Topic #8 (7, u'0.016*"like" + 0.016*"im" + 0.014*"doesnt" + 0.014*"would" + 0.012*"look" + 0.012*"anymore" +
0.012*"seeing" + 0.012*"protesting" + 0.011*"believed" + 0.007*"nunes"')

Topic #9 (8, u'0.020*"oooooooooooooooo" + 0.020*"shiiit" + 0.011*"got" + 0.011*"administration" + 0.011*"bad" +
0.009*"best" + 0.009*"voted" + 0.009*"feel" + 0.009*"wiretapped" + 0.009*"sum"')

Topic #10 (9, u'0.028*"president" + 0.017*"dont" + 0.011*"change" + 0.011*"climate" + 0.010*"white" +
0.008*"back" + 0.008*"america" + 0.008*"stepstoreverseclimatechange" + 0.007*"house" + 0.007*"impeach"')
```
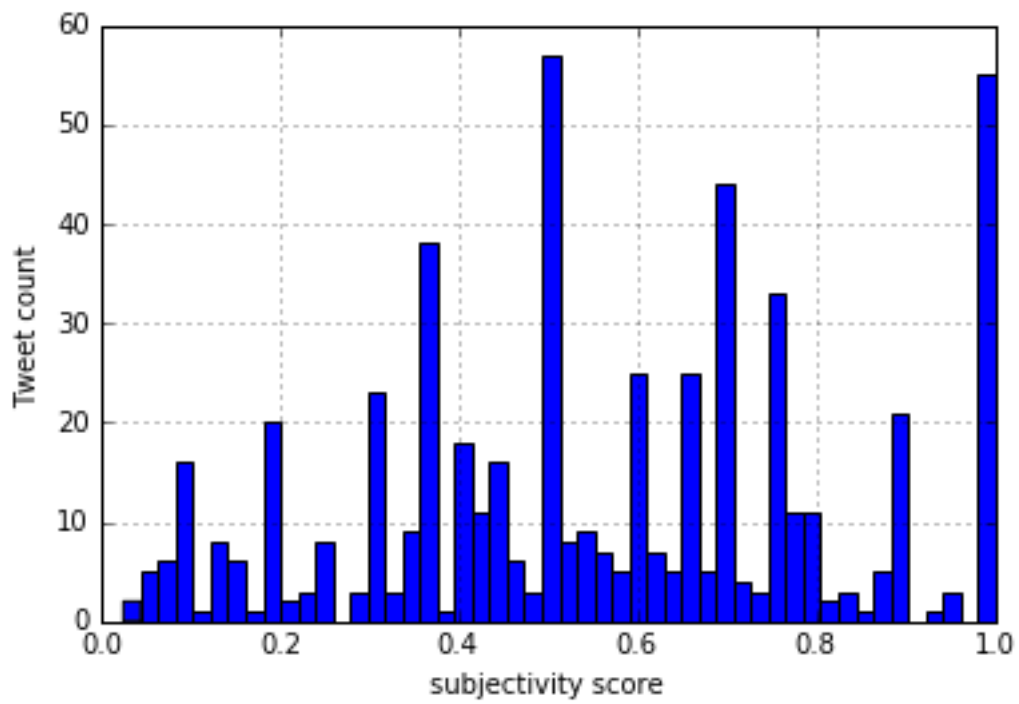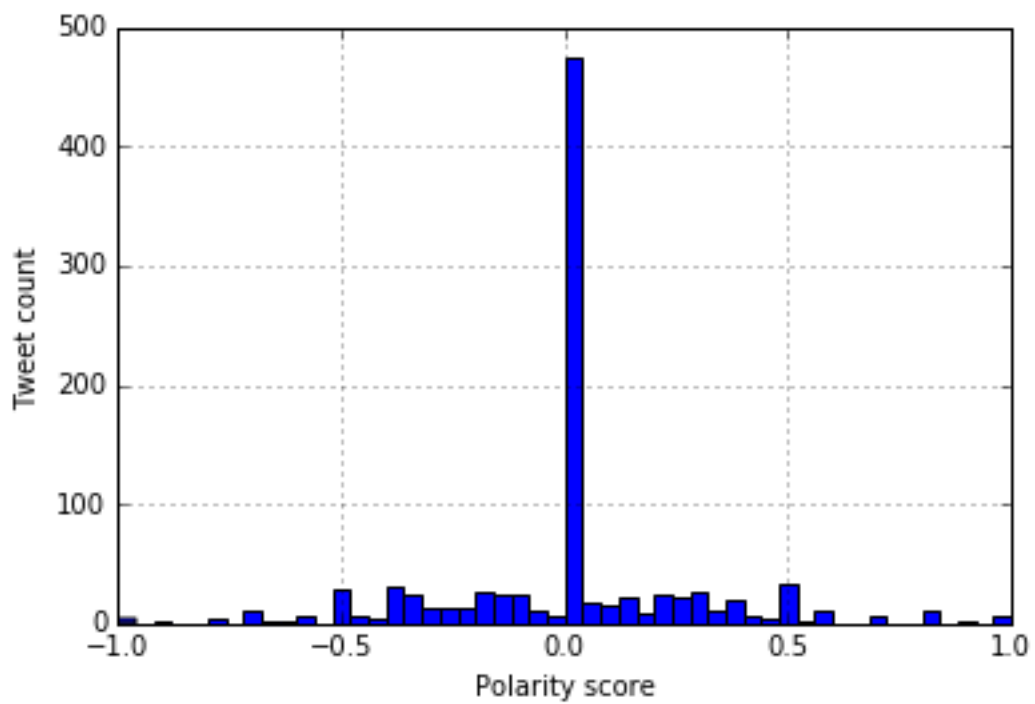
*b.) NMF:*

```
Topic 0: huntington beach pro rally violence
Topic 1: oh ryan paul irony plan
Topic 2: holy shit thread fuck got
Topic 3: brutal quotes ive seen russia
Topic 4: news detained protesters alert california
Topic 5: magamarch supporters anti way assembly
Topic 6: know black people rent refused
Topic 7: president health care doesnt democrats
Topic 8: obama didnt heard ass ya
Topic 9: deal art ready plea gettin
```
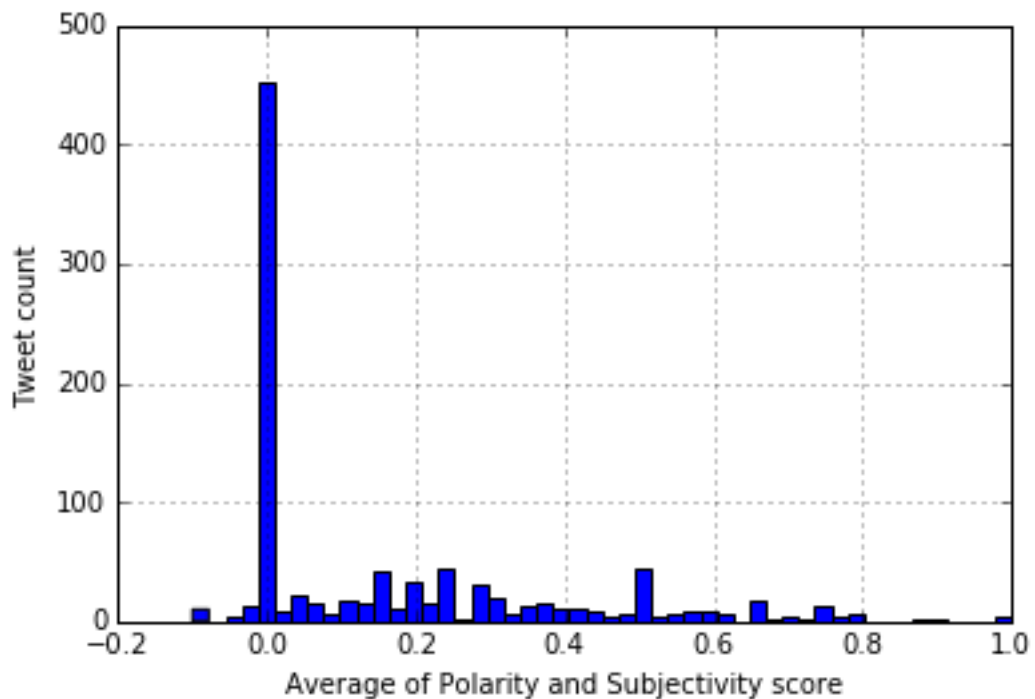
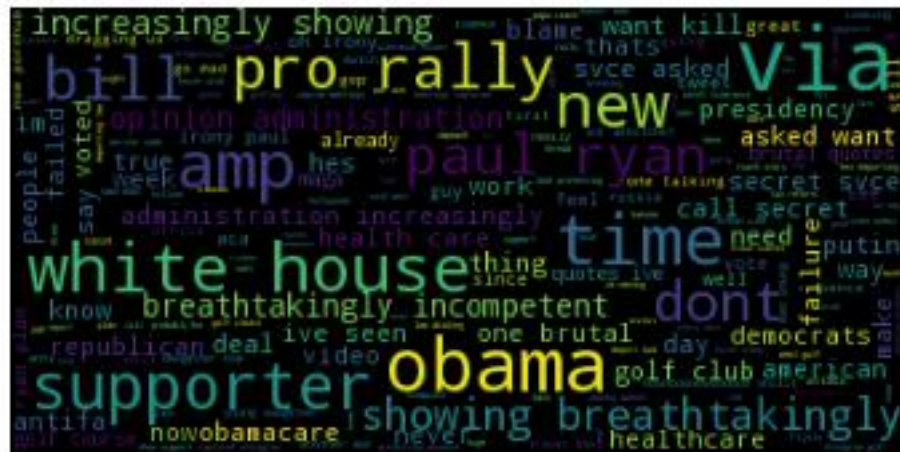# New York State

*a.) subjectivity score:*



*b.) Polarity Score:*

*c.) Average of Subjectivity & Polarity*



**Word Cloud:**



From the New york's word cloud: While people are talking about the pro rally, they look like tweeting more about the Obama care healthcare reform bill and it's failure.

**Topic Modeling:**

a.) LDA:

```
Topic #1 (0, u'0.021*"enough" + 0.018*"one" + 0.018*"seen" + 0.017*"real" + 0.017*"ive" + 0.017*"ever" +
0.015*"quotes" + 0.015*"brutal" + 0.014*"america" + 0.014*"ryan"')

Topic #2 (1, u'0.016*"news" + 0.013*"watch" + 0.013*"side" + 0.013*"ensues" + 0.013*"tolerance" + 0.011*"one" +
0.009*"yet" + 0.009*"art" + 0.009*"california" + 0.009*"violation"')

Topic #3 (2, u'0.035*"pro" + 0.035*"rally" + 0.013*"violence" + 0.009*"magamarch" + 0.009*"philadelphia" +
0.008*"supporters" + 0.007*"black" + 0.007*"maga" + 0.007*"huntington" + 0.007*"beach"')

Topic #4 (3, u'0.028*"administration" + 0.028*"opinion" + 0.028*"showing" + 0.026*"increasingly" +
0.026*"breathtakingly" + 0.026*"incompetent" + 0.017*"president" + 0.010*"says" + 0.010*"treason" +
0.010*"russiagate"')

Topic #5 (4, u'0.023*"nsa" + 0.016*"congress" + 0.014*"project" + 0.014*"please" + 0.012*"says" +
0.012*"whistleblower" + 0.012*"spied" + 0.012*"court" + 0.012*"supreme" + 0.012*"true"')

Topic #6 (5, u'0.016*"case" + 0.015*"worse" + 0.012*"bannon" + 0.012*"believed" + 0.012*"watergate" +
0.011*"wish" + 0.011*"health" + 0.010*"lol" + 0.010*"president" + 0.010*"care"')

Topic #7 (6, u'0.049*"president" + 0.032*"want" + 0.027*"kill" + 0.026*"secret" + 0.025*"asked" + 0.023*"svce" +
0.023*"call" + 0.010*"like" + 0.009*"huge" + 0.009*"sincerely"')

Topic #8 (7, u'0.011*"like" + 0.011*"made" + 0.011*"weekend" + 0.011*"deal" + 0.011*"consecutive" +
0.011*"property" + 0.011*"president" + 0.010*"th" + 0.008*"remember" + 0.008*"days"')
```
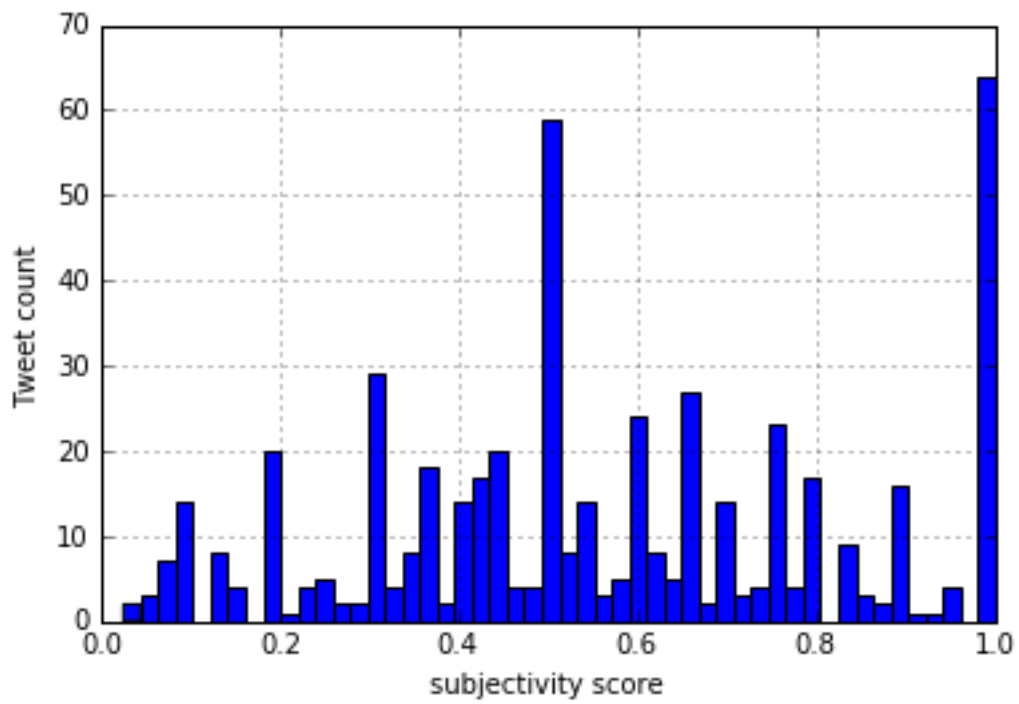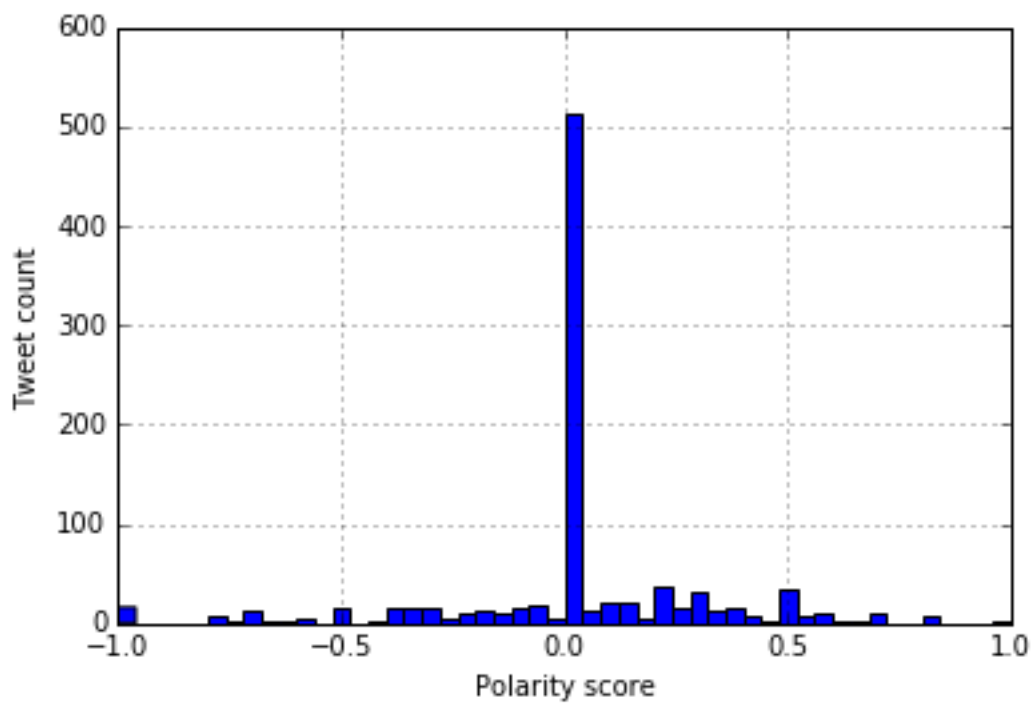
b.) NMF:

```
Topic 0: breathtakingly incompetent increasingly opinion showing
Topic 1: svce secret kill asked want
Topic 2: golf meetings house white course
Topic 3: ryan oh paul irony plan
Topic 4: ive brutal quotes seen word
Topic 5: like think rally pro im
Topic 6: war talking dragging time crime
Topic 7: ooooooooooooooo shiiit week earlier millions
Topic 8: president pornhub mad young daughter
Topic 9: watch ensues tolerance hes bad
```
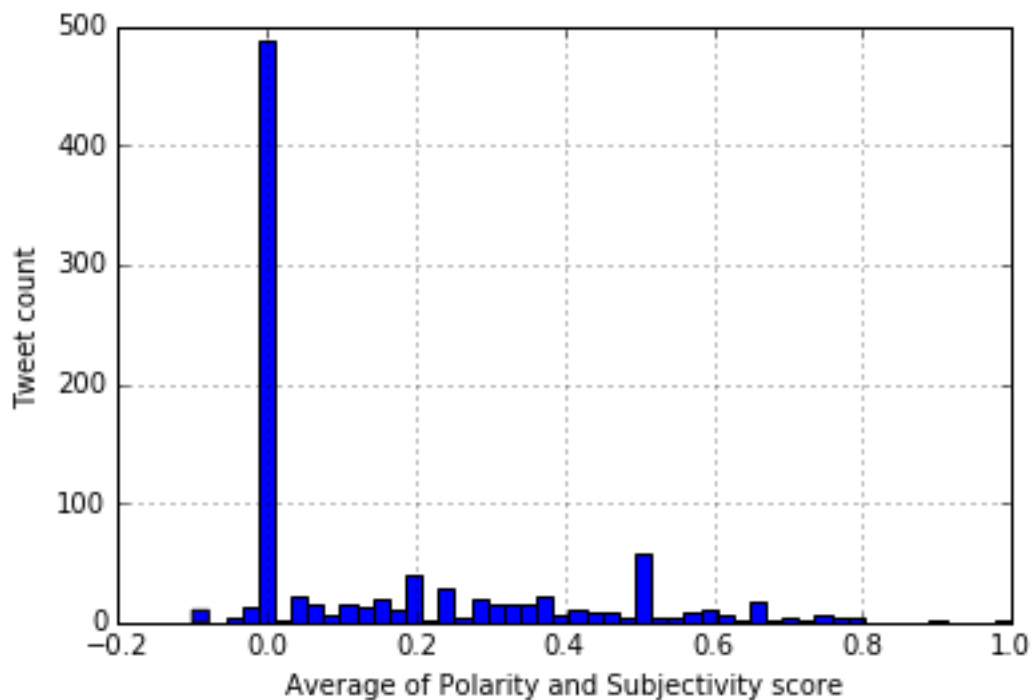
# Florida State

*a.) subjectivity score:*



*b.) Polarity Score:*

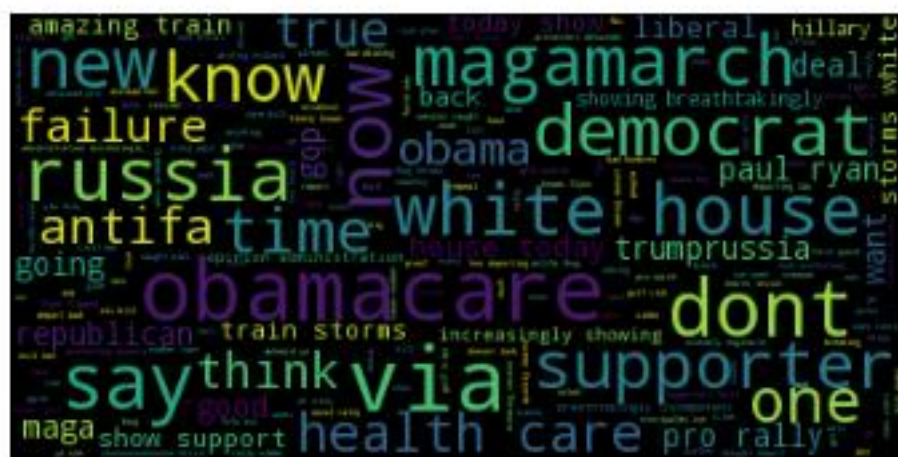The polarity distribution of tweets from Florida shows some strong negative emotions tweeted by the people.

*c.) Average of Subjectivity & Polarity:*



**Word Cloud:**



Tweets by the users from Florida have more references to the word "Russia", "antifa", "magamarch".

**Topic Modeling:**

a.) LDA:

```
Topic #1 (0, u'0.023*"opinion" + 0.021*"administration" + 0.020*"increasingly" + 0.020*"incompetent" +
0.020*"showing" + 0.020*"breathtakingly" + 0.015*"ever" + 0.013*"ive" + 0.013*"seen" + 0.013*"one"')

Topic #2 (1, u'0.017*"supporters" + 0.017*"president" + 0.014*"magamarch" + 0.014*"anti" + 0.014*"way" +
0.014*"force" + 0.014*"assembly" + 0.012*"remove" + 0.012*"democrats" + 0.012*"wall"')

Topic #3 (2, u'0.026*"president" + 0.025*"ryan" + 0.024*"doesnt" + 0.024*"im" + 0.020*"paul" + 0.020*"oh" +
0.019*"like" + 0.017*"protesting" + 0.017*"plan" + 0.017*"seeing"')

Topic #4 (3, u'0.019*"care" + 0.015*"man" + 0.015*"via" + 0.013*"today" + 0.013*"evil" + 0.013*"dem" +
0.013*"trumpcare" + 0.011*"number" + 0.011*"halls" + 0.011*"speeches"')

Topic #5 (4, u'0.021*"flynn" + 0.018*"breaking" + 0.018*"flipped" + 0.015*"senator" + 0.015*"caught" +
0.015*"plot" + 0.013*"team" + 0.013*"gets" + 0.013*"president" + 0.013*"likely"')

Topic #6 (5, u'0.025*"support" + 0.024*"white" + 0.023*"house" + 0.022*"today" + 0.021*"show" + 0.021*"amazing" +
0.021*"train" + 0.021*"storms" + 0.011*"hes" + 0.011*"russiagate"')

Topic #7 (6, u'0.017*"march" + 0.016*"pro" + 0.014*"law" + 0.014*"hes" + 0.013*"obamacare" + 0.013*"would" +
0.013*"bad" + 0.013*"thought" + 0.013*"husband" + 0.013*"deporting"')

Topic #8 (7, u'0.032*"rally" + 0.022*"antifa" + 0.019*"golf" + 0.016*"first" + 0.016*"thread" + 0.015*"video" +
0.013*"true" + 0.011*"lesson" + 0.011*"thug" + 0.011*"throws"')
```
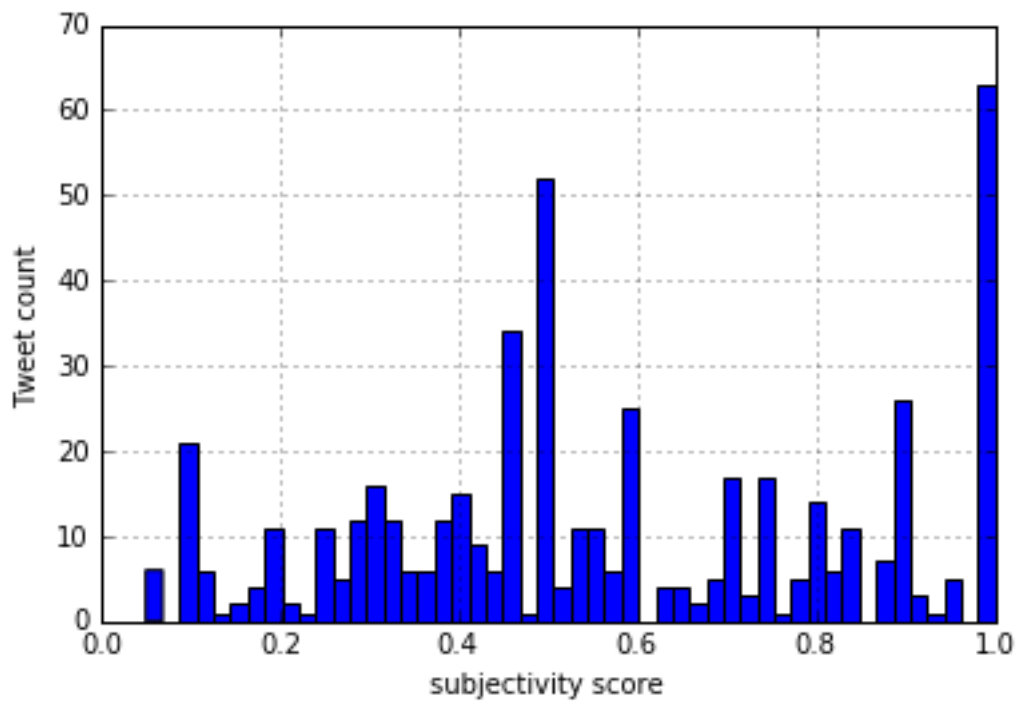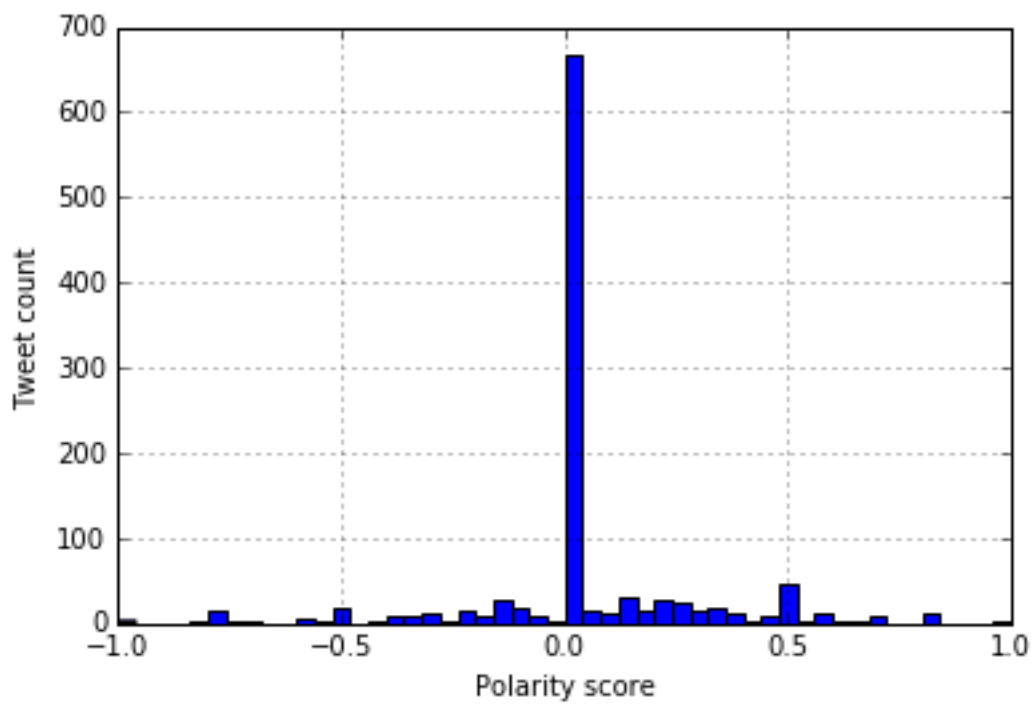
b.) NMF:

```
Topic 0: house white today support storms
Topic 1: president golf number trips taken
Topic 2: oh ryan paul irony plan
Topic 3: like im doesnt protesting anymore
Topic 4: senator breaking plot caught releases
Topic 5: team flynn flipped known likely
Topic 6: breathtakingly showing increasingly incompetent opinion
Topic 7: hes bad deporting husband thought
Topic 8: pro rally news alert california
Topic 9: magamarch supporters assembly way anti
```
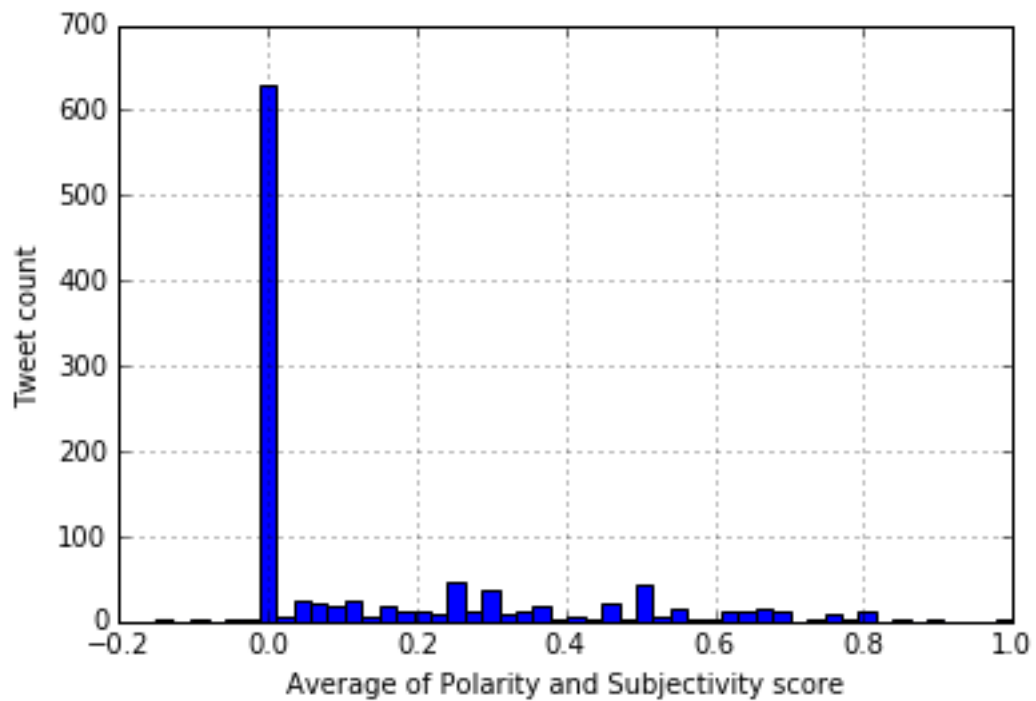
# Arizona State

*a.) subjectivity score:*



*b.) Polarity Score:*

*c.) Average of Subjectivity & Polarity:*



**Word Cloud:**

**Topic Modeling:**

a.) LDA:

```
Topic #1 (0, u'0.056*"president" + 0.051*"america" + 0.046*"people" + 0.041*"pay" +
0.038*"defrauded" + 0.038*"agreed" + 0.038*"united" + 0.038*"states" + 0.013*"judge" +
0.013*"settlement"')
Topic #2 (1, u'0.015*"awareness" + 0.014*"assault" + 0.014*"sexual" + 0.012*"month" + 0.011*"would"
+ 0.011*"april" + 0.011*"prevention" + 0.011*"national" + 0.010*"ceremony" + 0.010*"im"')
Topic #3 (2, u'0.055*"pussy" + 0.051*"grab" + 0.045*"may" + 0.039*"em" + 0.038*"month" +
0.017*"high" + 0.017*"official" + 0.017*"unmasked" + 0.017*"intelligence" + 0.013*"says"')
Topic #4 (3, u'0.020*"russia" + 0.011*"plan" + 0.011*"like" + 0.011*"scandal" + 0.009*"via" +
0.008*"please" + 0.008*"end" + 0.008*"immunity" + 0.008*"court" + 0.007*"nunes"')
Topic #5 (4, u'0.068*"executive" + 0.038*"order" + 0.034*"orders" + 0.031*"signing" +
0.031*"without" + 0.030*"ceremony" + 0.029*"signed" + 0.029*"walked" + 0.021*"jordan" +
0.021*"dion"')
Topic #6 (5, u'0.056*"google" + 0.056*"university" + 0.056*"fridayfeeling" + 0.054*"price" +
0.054*"tom" + 0.053*"rikersisland" + 0.053*"tdov" + 0.053*"closerikers" + 0.053*"maps" +
0.053*"ripselena"')
Topic #7 (6, u'0.037*"house" + 0.028*"white" + 0.017*"surveillance" + 0.017*"report" +
0.015*"nominee" + 0.015*"unmasking" + 0.015*"gop" + 0.013*"characters" + 0.013*"cards" +
0.013*"veep"')
Topic #8 (7, u'0.027*"get" + 0.020*"deserve" + 0.019*"via" + 0.018*"destroyed" + 0.018*"nominate" +
0.015*"president" + 0.013*"sean" + 0.012*"spicer" + 0.012*"job" + 0.011*"oh"')
Topic #9 (8, u'0.017*"liar" + 0.015*"dangerous" + 0.015*"confirmed" + 0.011*"fbi" +
0.010*"president" + 0.010*"amp" + 0.010*"would" + 0.010*"conspiracy" + 0.008*"russian" +
0.008*"well"')
Topic #10 (9, u'0.024*"via" + 0.017*"russiagate" + 0.011*"russia" + 0.010*"resist" + 0.009*"world" +
0.009*"amp" + 0.007*"end" + 0.007*"presidency" + 0.007*"pence" + 0.007*"hell"')
```

b.) NMF:

```
Topic 0: rikersisland tdov closerikers maps tom
Topic 1: executive order ceremony signed walked
Topic 2: defrauded states united agreed pay
Topic 3: dion jordan island flashbackfriday cesarchavezday
Topic 4: grab pussy em month awareness
Topic 5: new think going like cover
Topic 6: huckabee et christie cruz al
Topic 7: destroyed deserve nominate myth great
Topic 8: russia hold president scandal plan
Topic 9: irony april heart pair dull
```

## Insights drawn from analysis:

1.) Subjectivity score is highest in the state of Texas, by which we can say that the users in the Texas are being more particular in expressing their opinions.

2.) As polarity scores are more positive in Florida and Arizona we can cay that users from these states are favorable to trump. Which can also be backed by the fact, he has won in both these states.

3.) People from Arizona are paying more tribute to the musical icon "selena" than in any other state, even more than her born state Texas.

4.) "General opinion in the New York about Trump's administration is increasingly incompetent" – this can be concluded from the LDA and NMF topic models.

5.) New York & Texas has more tweets in the range of -0.5 to +0.5 that are not falling on 0.0. This could possibly mean that users from these states are more comfortable in showing their emotions (both positive & negative) on twitter.

6.) California tweeted more about the **Trump's Pro Rally** that happened in **Huntington beach,** CA. But the other states were surprisingly silent about this Rally.

7.) Arizona is talking more about **rikers island**, after the announcement by the city mayor that it is going to get be closed. And the main top keywords are '**rickersisland', 'rickersisland google', 'closerickers', 'google maps'**. From this, we could infer that most of the users from Arizona doesn't know about the location of this jail or the rickers lisland.