
Project 2B - Spark Machine Learning Project

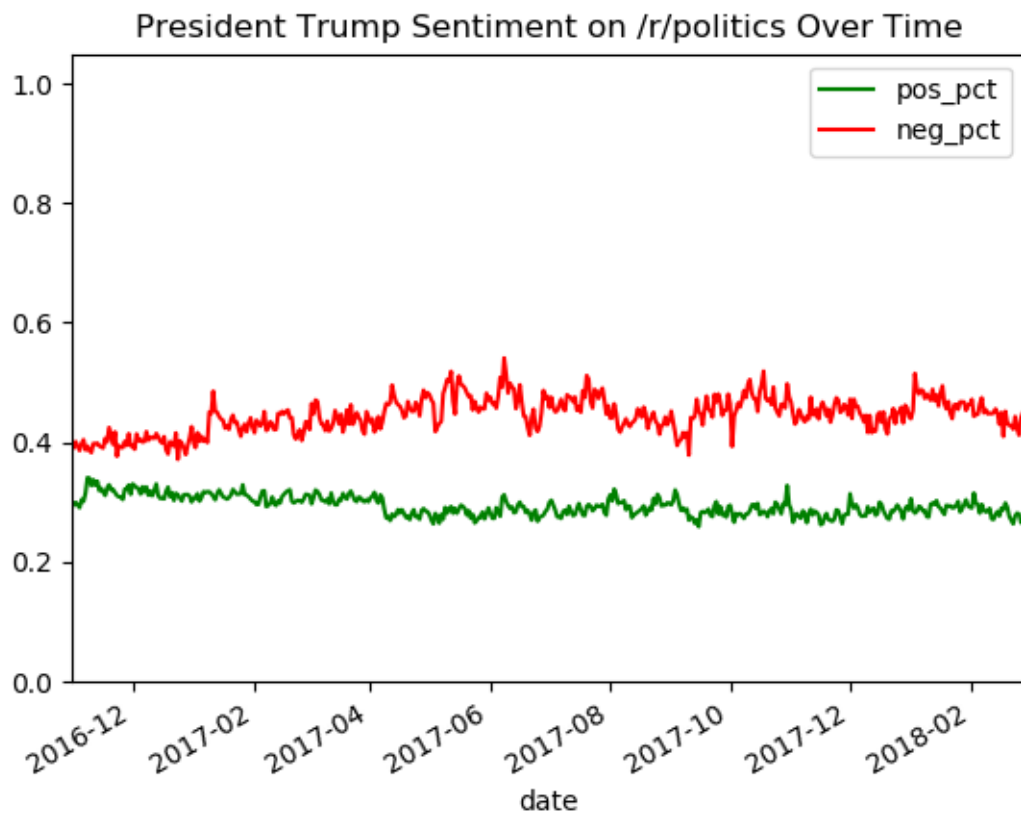
June 3, 2019

Wilson Jusuf
Student ID: 404997407
University of California, Los Angeles

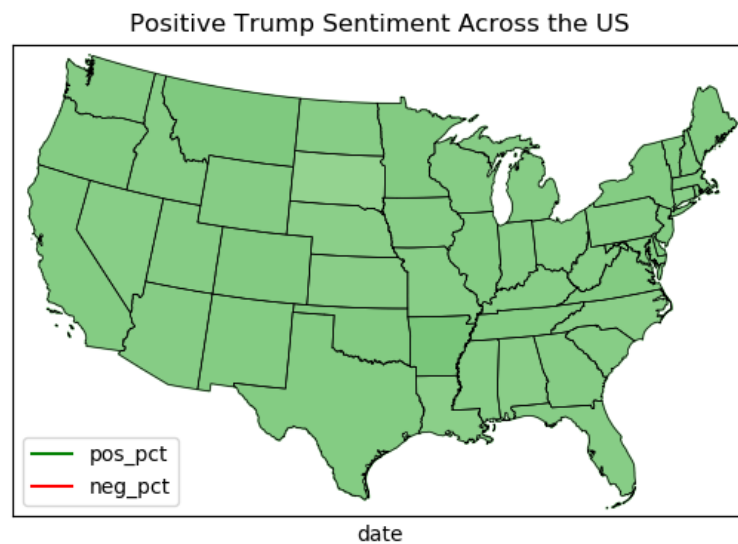
Hanif Leoputera Lim
Student ID: 504971751
University of California, Los Angeles

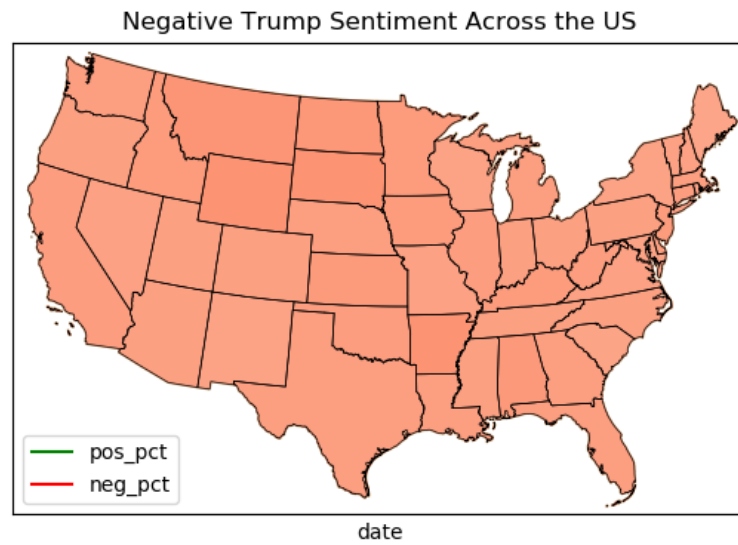
1 Findings

1.1 Sentiment Time series plot by day

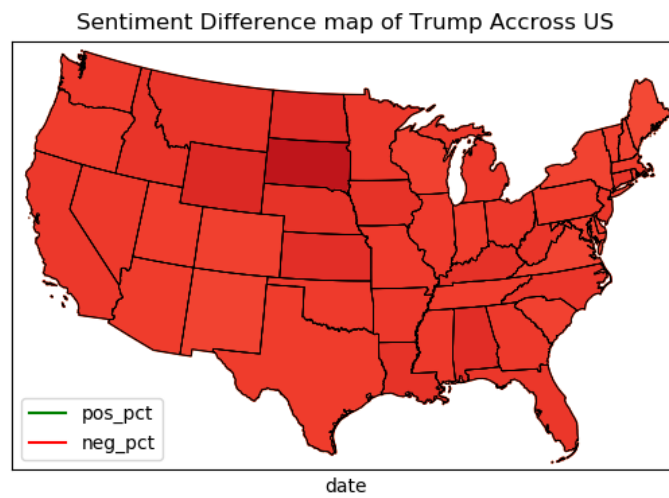


1.2 Sentiment by State





1.3 Sentiment Difference by State



1.4 Top 10 Positive Stories

```

title,pct_pos
Sean Spicer finally gets to meet Pope Francis,1.0
"Breitbart under Bannon: Breitbarts comment section reflects alt-right, anti-Semitic
↪ language",1.0
"After Trump tweets, pressure grows for full Russia investigation",1.0
Confirmed: Donald Trump Says He Will Take $1 Salary as President,1.0
Is President Trumps Immigration Order a 'Muslim Ban'?,1.0
Kamala Harris says she'll co-sponsor single-payer health care bill,1.0
The Left Has Lost its Mind: Cory Booker to Testify Against Sessions,1.0
EU says to stick to Iran rapprochement despite Trump's criticism,1.0
"Tom Perez: Win or Lose, I Will Work to Unite the Democratic Party",1.0
Maxine Waters: Trump is the most deplorable person Ive ever met,1.0

```

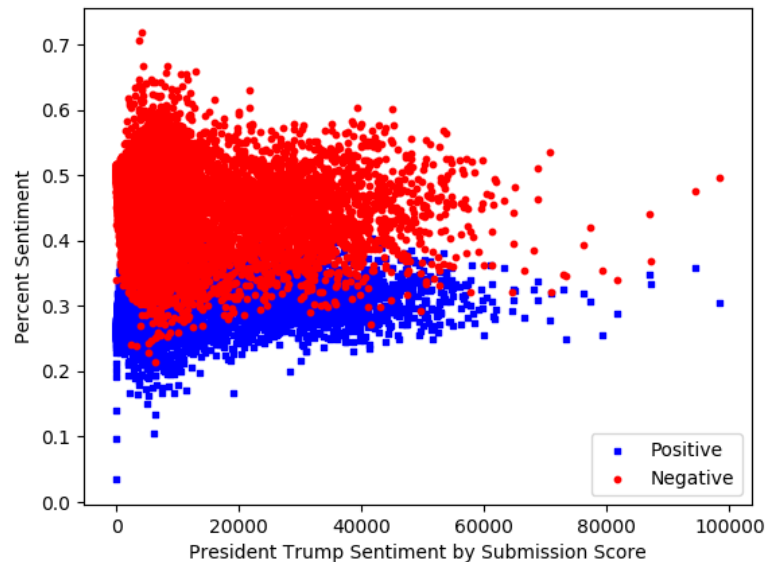
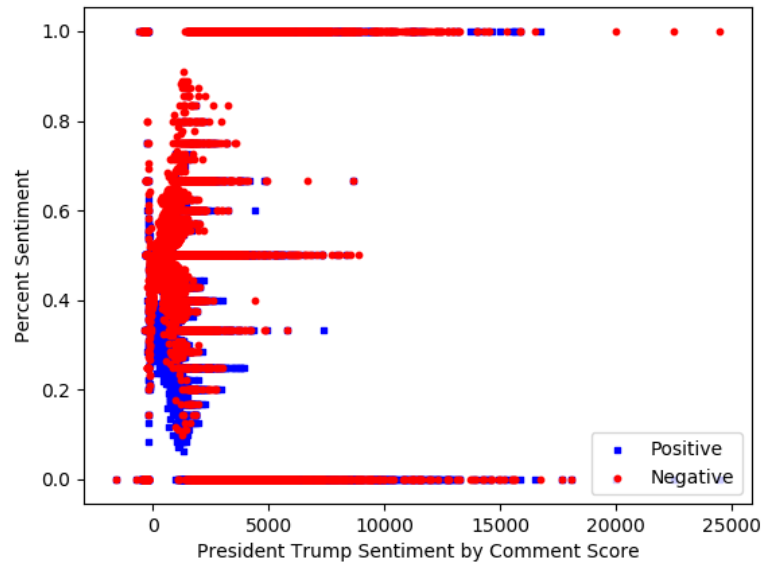
1.5 Top 10 Negative Stories

```

title,pct_neg
"Trump's Government wants to punish people who speak out against him, especially if
↪ they're black women",1.0
"Mnuchin Asked for Tax Dollars to Fund His Honeymoon, Despite $300 Million Net
↪ Worth",1.0
Trump and North Korea war of words escalates,1.0
Trump's 'Loser' Speech in Bethlehem,1.0
Trump revives criticism of 'both sides' in Charlottesville,1.0
How Deutsche Bank Enabled A Dirty Offshore Bank To Move Dark Money,1.0
Trump again seizes on terror incident to call for travel ban,1.0
Trump says hes fairly close to deal on young immigrants,1.0
Judge questions tossing out Arpaio conviction after pardon,1.0
Trump repeats 'both sides' Charlottesville rhetoric,1.0

```

1.6 Scatterplot of Comment Score and Submission Score vs Sentiment



1.7 Summary of findings

Overall, the sentiment of trump was overall more negative in all states. The positive percentage was considerably lower than the negative percentage. However, the intensity of the negative sentiment varies by state.

Both negative and positive stories involve Trump in the title. It relates to either Trump or his administration.

Looking at the scatterplot, we see that the sentiment percentages converge to neutral (0.5) as the comment score increases. However, there are also extreme stories (sentiment of 1.0 (100%) or 0.0 (0%)) getting high comment scores.

For submission score, the sentiment weakly converges to around 0.4 as submission score increases. The submissions that are more 'neutral' therefore tend to have higher scores.

2 Questions

1. Take a look at `labeled_data.csv`. Write the functional dependencies implied by the data. Set $I = \text{Input_id}$, $D = \text{labeldem}$, $G = \text{labelgop}$, $T = \text{labeldjt}$.

The dependencies are:

- $I \rightarrow D$
- $I \rightarrow G$
- $I \rightarrow T$

2. Take a look at the schema for comments. Forget BCNF and 3NF. Does the data frame look normalized? In other words, is the data frame free of redundancies that might affect insert/update integrity? If not, how would we decompose it? Why do you believe the collector of the data stored it in this way?

The data frame isn't normalized. There are redundancies like `(subreddit_id, subreddit)`, so updating the subreddit would require us to update `subreddit_id` too. This could lead to insert/update integrity if the DBA isn't careful. To decompose it, we could break relation R into $R_2(\text{subreddit_id}, \text{subreddit})$, while R_1 will have everything in R except for `subreddit`. However, one reason why the previous schema was used is because joining two separate tables are expensive. Querying which subreddit a comment come from, is probably so common that these two tables are frequently joined, hence computationally wasteful.

3. Pick one of the joins that you executed for this project. Rerun the join with `.explain()` attached to it. Include the output. What do you notice? Explain what Spark SQL is doing during the join. Which join algorithm does Spark seem to be using?

```
== Physical Plan ==
*(4) Project [id#14, link_id#2193 AS submission_id#2399, body#4, author_flair_text#3 AS state#2401, created_utc#10L, t
title#106, score#92L, score#20L AS comm_score#2405L]
+- *(4) BroadcastHashJoin [link_id#2193], [id#69], Inner, BuildRight
   :- *(4) Project [id#14, body#4, author_flair_text#3, pythonUDF0#2415 AS link_id#2193, created_utc#10L, score#20L]
   : +- BatchEvalPython [<lambda>(link_id#16)], [author_flair_text#3, body#4, created_utc#10L, id#14, link_id#16, sco
re#20L, pythonUDF0#2415]
   :   +- *(2) Project [author_flair_text#3, body#4, created_utc#10L, id#14, link_id#16, score#20L]
   :   +- *(2) Project [author#0, author_cakeday#1, author_flair_css_class#2, author_flair_text#3, body#4, can_gi
ld#5, can_mod_post#6, collapsed#7, collapsed_reason#8, controversiality#9L, created_utc#10L, distinguished#11, edited#
12, gilded#13L, id#14, is_submitter#15, link_id#16, parent_id#17, permalink#18, retrieved_on#19L, score#20L, stickied#
21, subreddit#22, subreddit_id#23, subreddit_type#24]
   :   +- *(2) Filter isNotNull(pythonUDF0#2414)
   :   +- BatchEvalPython [<lambda>(link_id#16)], [author#0, author_cakeday#1, author_flair_css_class#2, au
thor_flair_text#3, body#4, can_gild#5, can_mod_post#6, collapsed#7, collapsed_reason#8, controversiality#9L, created u
tc#10L, distinguished#11, edited#12, gilded#13L, id#14, is_submitter#15, link_id#16, parent_id#17, permalink#18, retrie
ved_on#19L, score#20L, stickied#21, subreddit#22, subreddit_id#23, ... 2 more fields]
   :   +- *(1) FileScan parquet [author#0,author_cakeday#1,author_flair_css_class#2,author_flair_text#3,
body#4,can_gild#5,can_mod_post#6,collapsed#7,collapsed_reason#8,controversiality#9L,created_utc#10L,distinguished#11,e
dited#12,gilded#13L,id#14,is_submitter#15,link_id#16,parent_id#17,permalink#18,retrieved_on#19L,score#20L,stickied#21,
subreddit#22,subreddit_id#23,subreddit_type#24] Batched: true, Format: Parquet, Location: InMemoryFileIndex[file:/home
/willyspinner/cs-projects/cs143-projects/2b/comments-minimal.parquet], PartitionFilters: [], PushedFilters: [], ReadSc
hema: struct<author:string,author_cakeday:boolean,author_flair_css_class:string,author_flair_text:string...
+- BroadcastExchange HashedRelationBroadcastMode(List(input[0, string, true]))
+- *(3) Project [id#69, score#92L, title#106]
   +- *(3) Filter isNotNull(id#69)
   +- *(3) FileScan parquet [id#69,score#92L,title#106] Batched: true, Format: Parquet, Location: InMemoryFile
```

Spark first sequentially scanned the parquet, select all the attributes in the relation and then select all the attributes that we wanted. Once projected, the smaller and filtered table will be broadcast hash joined on the attribute we specified. Once joined, the attributes that we want are selected again. In Spark, table will be transmitted through the network with broadcast hash joins. Transmitting a huge table is prohibitively slow. So, we feel this is the reason why spark seems to redundantly select only the attributes we need before it gets broadcasted. Spark is using broadcast hash join. Spark is usually ran across many machines i.e. RDD, thus it usually have to join across partitions in different machines. Thus, the reason why they use broadcast hash join.

3 Classifier

We changed the positive and negative classifiers' thresholds to their optimal values based on how close it is to the point (FPR = 0, TPR = 0). Formally, we want the threshold T as the minimal euclidean distance:

$$\hat{T} = \underset{T}{\operatorname{argmin}} \sqrt{(\operatorname{FPR}_T)^2 + (1 - \operatorname{TPR}_T)^2} \quad (1)$$

Here are our classifier auROC graphs:

