# Correlating Poll Results with Twitter Sentiment on Obamacare

W205
Tingwen Bao
Jay Cordes
Alex Jamar
Michelle Liu
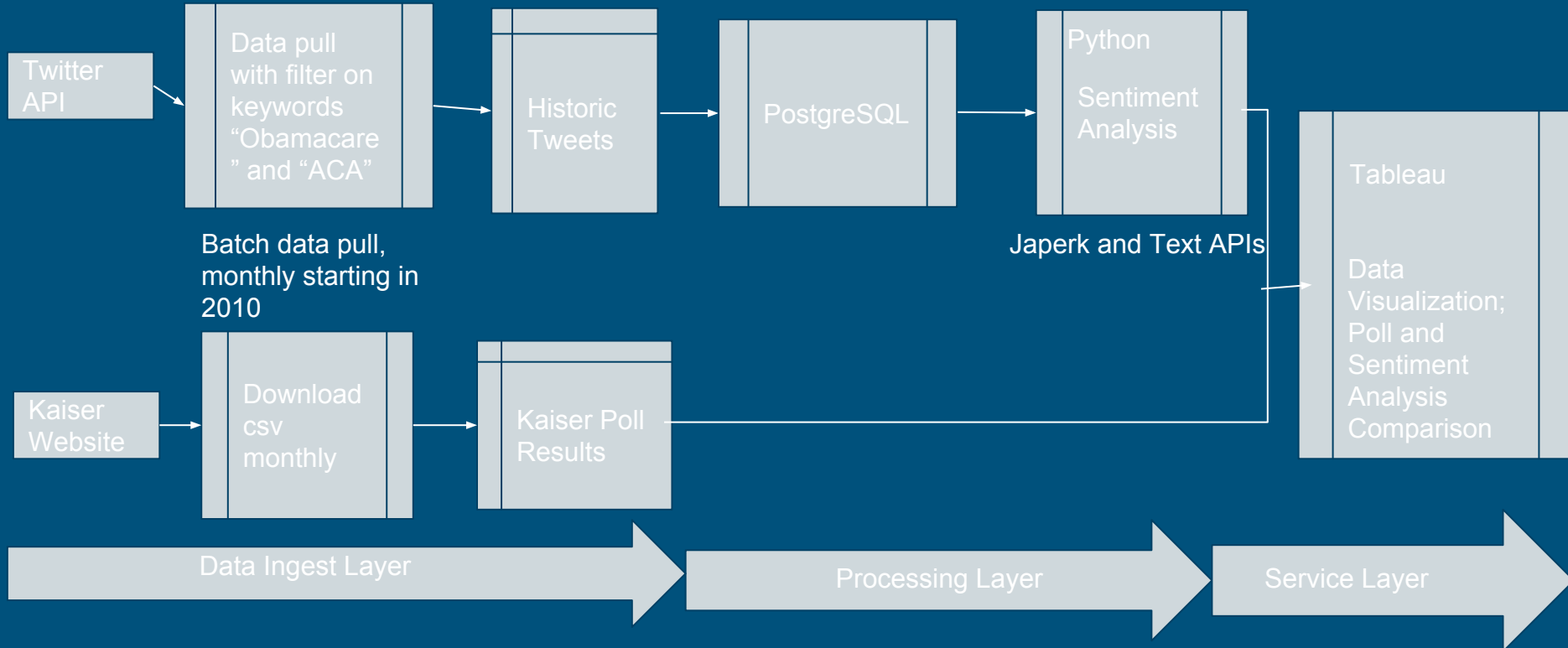
# Overview of the problem

- Develop a data infrastructure to hold both ACA poll results from Kaiser and Twitter posts

- Analyze the correlation between the result of public polls and sentiment on Twitter

- Twitter claims that the real-time, public orientation of its social network makes it a reliable barometer of the public's constantly changing moods and interests

# Overall Architecture

# Acquisition of Twitter Data

- Adopted Python library, GetOldTweets

    - Bypass the time constraint limitations of Twitter API

    - Search for tweets in English about Obamacare or ACA since its introduction in 2010

    - Collects 1000 tweets per day every 7 days.

- Ran on an Amazon Web Services Elastic Cloud Compute virtual machine
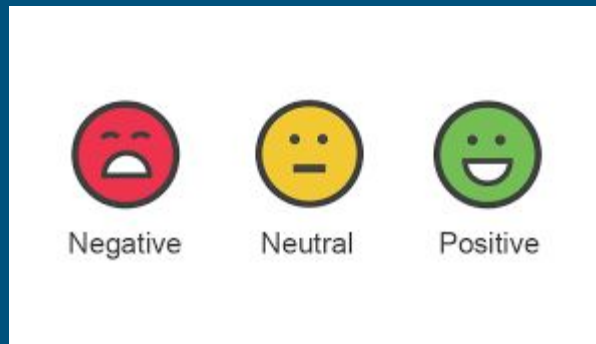
- Stored about 400,000 Tweets

# The Database

- Created a python script to import tweets tab-delimited file into PostgreSQL

- Used command-line argument with filename for flexibility

- Added surrogate primary key (bigserial data type is auto-incrementing) and indexes

- Exceptions handled gracefully and script keeps importing

- Skipped transactions and kept everything simple and maintainable

- PostgreSQL is very easy to use and facilitated later sentiment scoring / analysis steps

# Sentiment Analysis

- Randomly selected 10% of the tweets for analysis

- Used 2 sentiment analysis APIs

    - Sent POST requests containing the Tweet text to the APIs

    - Received JSON objects in return with either "positive," "negative," or "neutral" labels

    - Mapped positive → 1, negative → -1, neutral → 0

- Updated Tweets database with the results


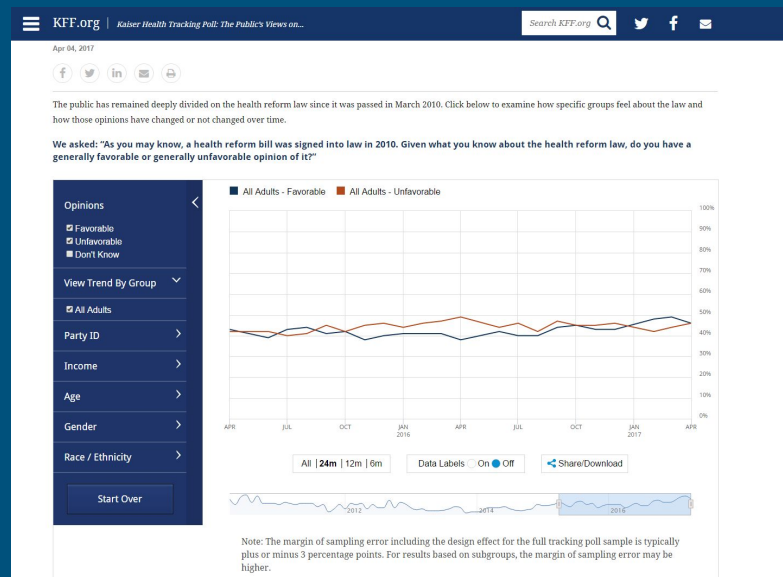
Negative    Neutral    Positive

# Evaluation & Visualization

Public Poll Data: Kaiser Health Tracking Poll

Why Tableau:

1. Can easily connect both databases and flat files
2. Relative strong visualization with interactive features
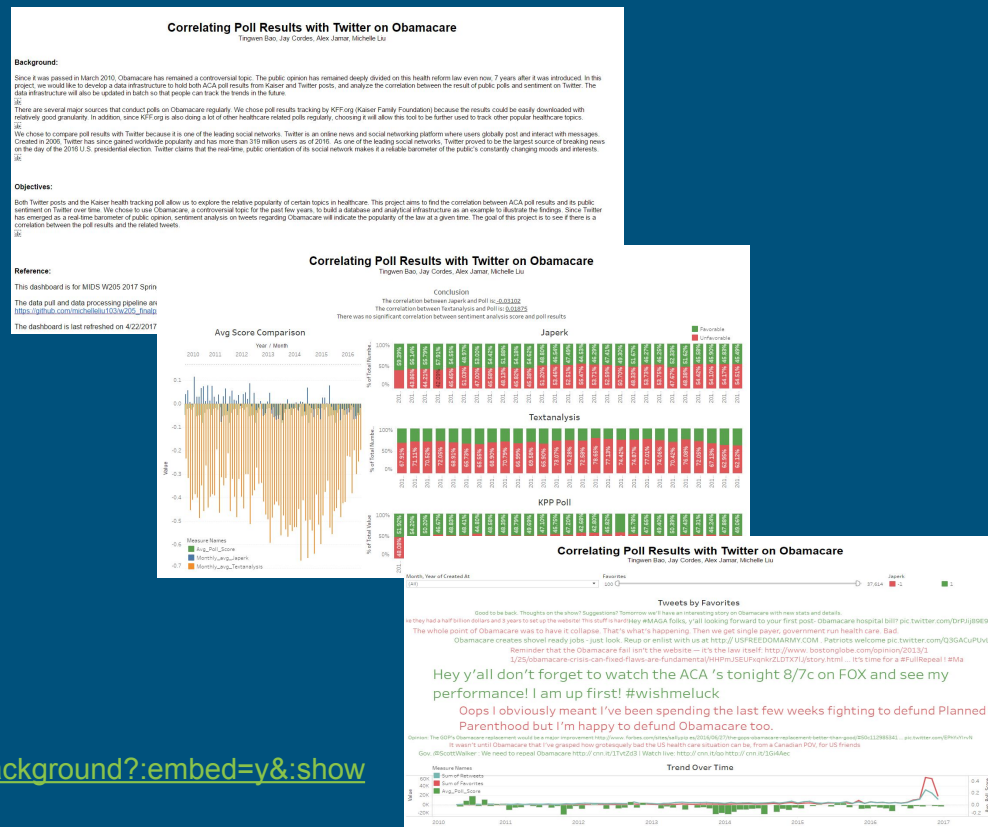3. Various ways to share/access

Drawbacks:

1. Not strong at stats
2. Could become slow after scale up

# Tableau Dashboard

To evaluate and visualize the result, we create a dashboard with 3 tabs:

- Background: a general introduction of the background and objective of this project

- Comparison: visuals used to compare the 2 different sentiment analysis method and understand its correlation with public poll data

- Explore Tweets: word cloud of top tweets by favorites and its trend by time with interactive filters

https://us-east-1.online.tableau.com/t/tingwen/views/Obamacare/Background?:embed=y&:show ShareOptions=true&:display_count=no&:showVizHome=no

# Results

- Using different sentiment analysis method yields different conclusion of Obamacare related tweets
    - Japerk's score indicates that the positive and negative tweets of Obamacare are close to a tie
    - Textanalysis's score indicates that there are about twice as many negative tweets as positive.

- KPP poll shows public has a nearly even spread on their opinion on Obamacare.
    - More positive at the start and then more negative, and most recently, almost even.

- No significant correlation between sentiment analysis score and poll results

- Tweets with most favorites and retweets happened in 2016 rather than when it was first introduced.
    - Right before election when people mentioned Obamacare to express their political stands.

# Future Work

Roadmap for improving the solution with increased usage and increasing data size: improving the ways in acquiring the data, storing and processing it more efficiently, and overall, scaling up our solution.

- Instead of pulling the same number of tweets every week, grab all related tweets or a proportional number of them
- Default the sentiment scores to null instead of 0 and/or score all of the tweets in our database
- Automate and import the polling data into a PostgreSQL table as well.
- Create aggregation table to reduce the processing time of Tableau.
- Strip out all of the tab characters in the tweets before writing it to CSV.

# Conclusion

- Text sentiment analysis is nowhere near perfect
- We did not find correlation between poll results and two types of Twitter sentiment scores
- The 3Vs were important challenges for our project: long term needs of the solution will have to deal with the sheer *volume* of the data coming in from Twitter, and to increase the *velocity* of doing sentiment analysis and processing of the data.
- Our program could provide a useful test for any new sentiment analysis approaches that are developed by validating them against polling data.
- If our app finds a strong correlation between Twitter sentiment scores and polling data, we would have provided evidence for the ability to do real-time polling on any subject.

# References

http://kff.org/interactive/kaiser-health-tracking-poll-the-publics-views-on-the-aca/#?aRange=twoYear

https://github.com/Jefferson-Henrique/GetOldTweets-python

https://market.mashape.com/textanalysis/sentiment-analysis

https://market.mashape.com/japerk/text-processing