

# W205 Final Project

Predicting Economic Indicators Using Open Data Sources

*Evyyatar, Daniel, Konniam*



# Problem

## **Economic indicators**

Crucial for business decision making  
(unemployment, consumer confidence)

## **Issues**

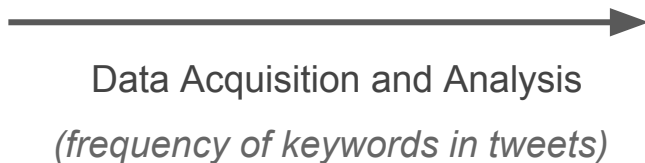
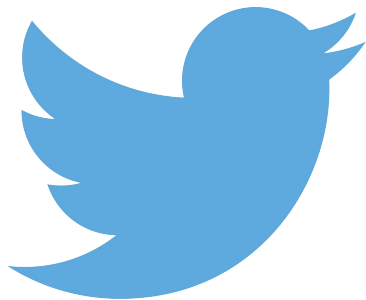
Infrequent updates  
Rigid schema  
Backward facing

**Opportunity to create real-time and predictive indicators that provide value to stakeholders in the economy**

# Proposal

## Data Sources

Twitter  
Stock Prices  
Initial Claims



## Outputs

Unemployment metrics  
Dashboard of metrics



**For this project we are implementing a subset of the overall proposal due to complexity and time constraints of the Twitter APIs**

# Data and Processing Dimensions

## DATA

**Size**  
6GB per day

**Structure**  
JSON (semi-structured)

**Velocity**  
1000 tweet/min

**Source Latency**  
hourly updates

## PROCESSING

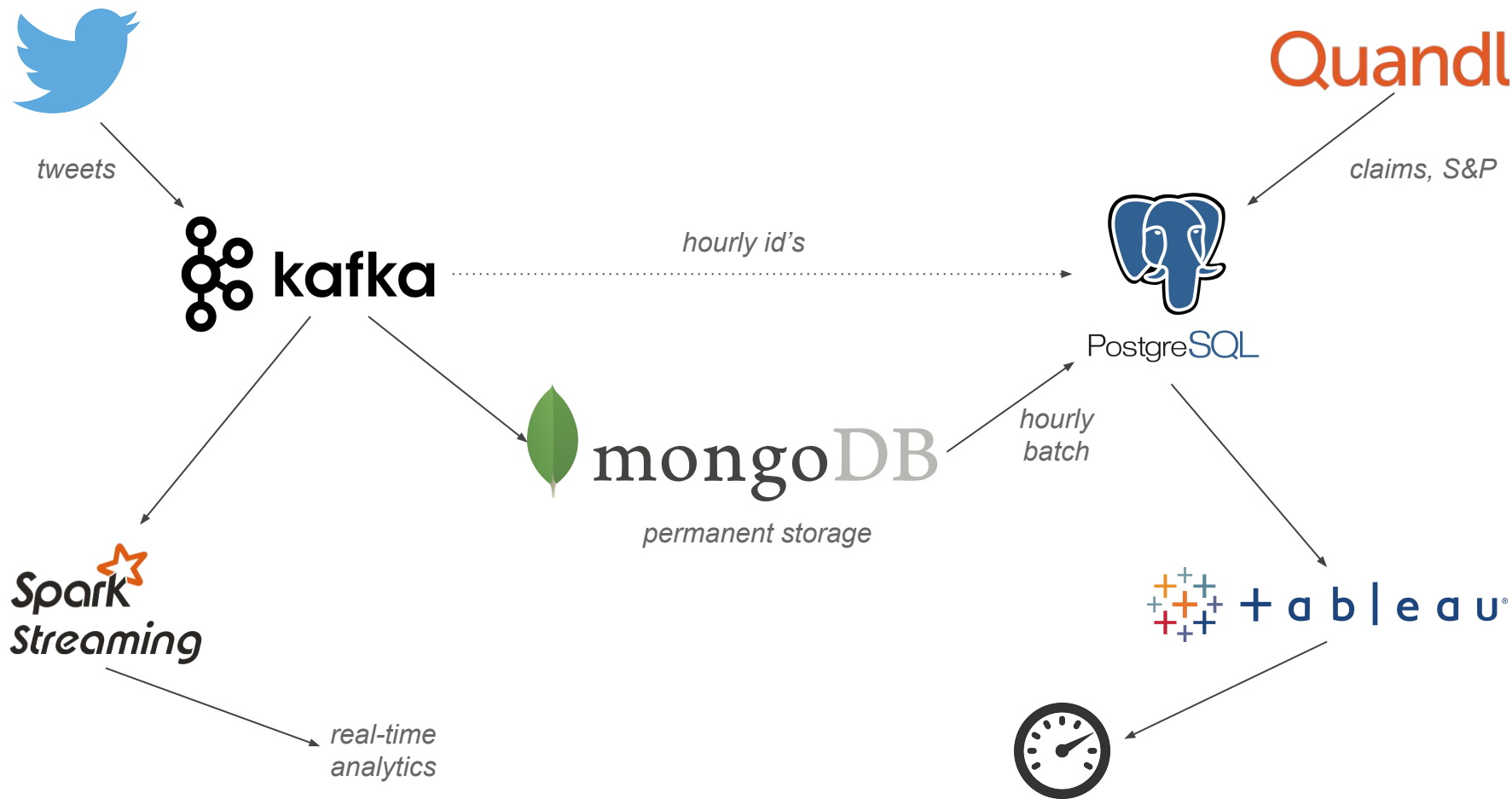
**Selectivity**  
process every tweet

**Processing Time**  
semi-batch

**Aggregation**  
keyword counting

**Precision**  
approximate is fine

# Architecture



# Choice of Tools

## **Kafka**

Backbone of ingestion  
Decoupled reading/writing  
Distributed

## **MongoDB**

Documents format  
Horizontal sharding available

## **PostgreSQL**

Convenient as serving layer  
Great for stocks/economic data

## **Spark Streaming**

Small batch approach

# Setup

## **Amazon Web Services**

m3.xlarge instance  
15GB ram  
400GB EBS (1.2K/3K IOPS)  
Amazon Linux  
Data Collection: 12/3-12/10

## **Software**

Kafka  
MongoDB  
PostgreSQL  
Spark

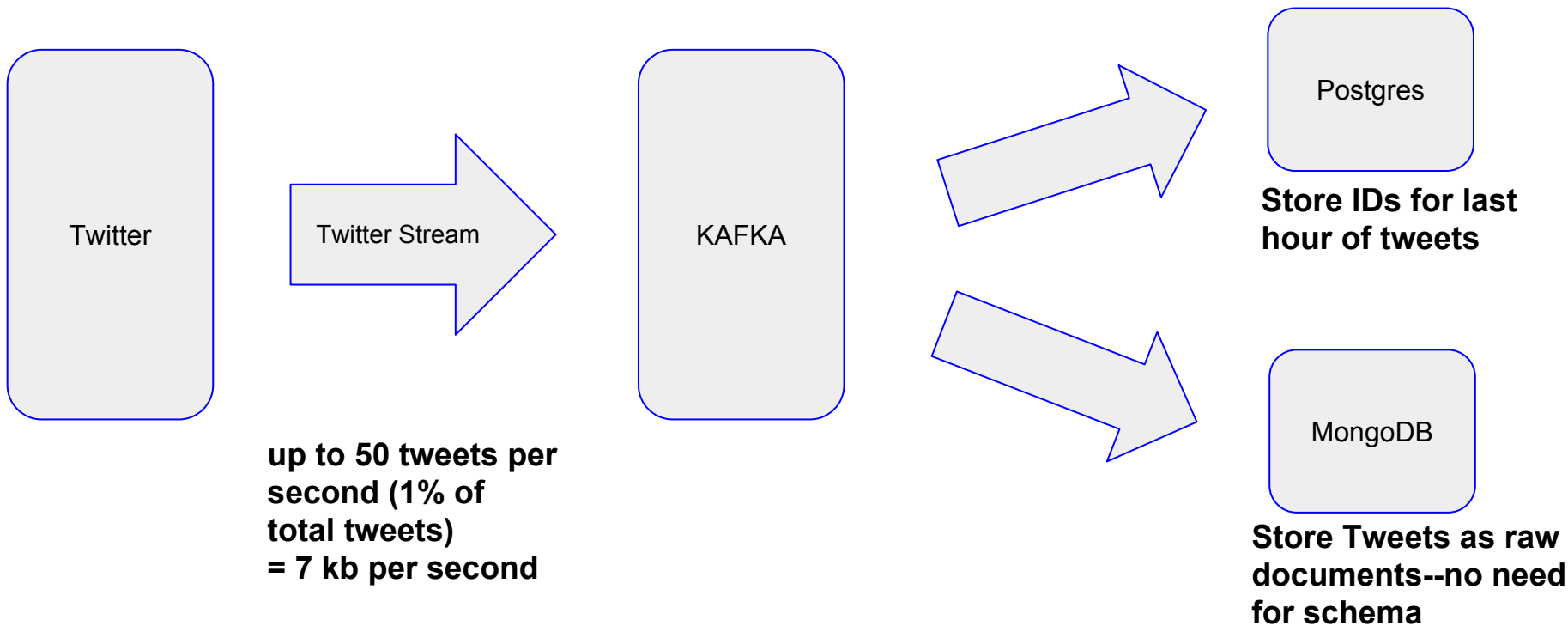
## **Python Packages**

tweepy  
kafka-python  
pymongo  
psycopg2  
Quandl

## **Data Collected**

10 Million Tweets  
42GB

# Data Ingestion - Twitter

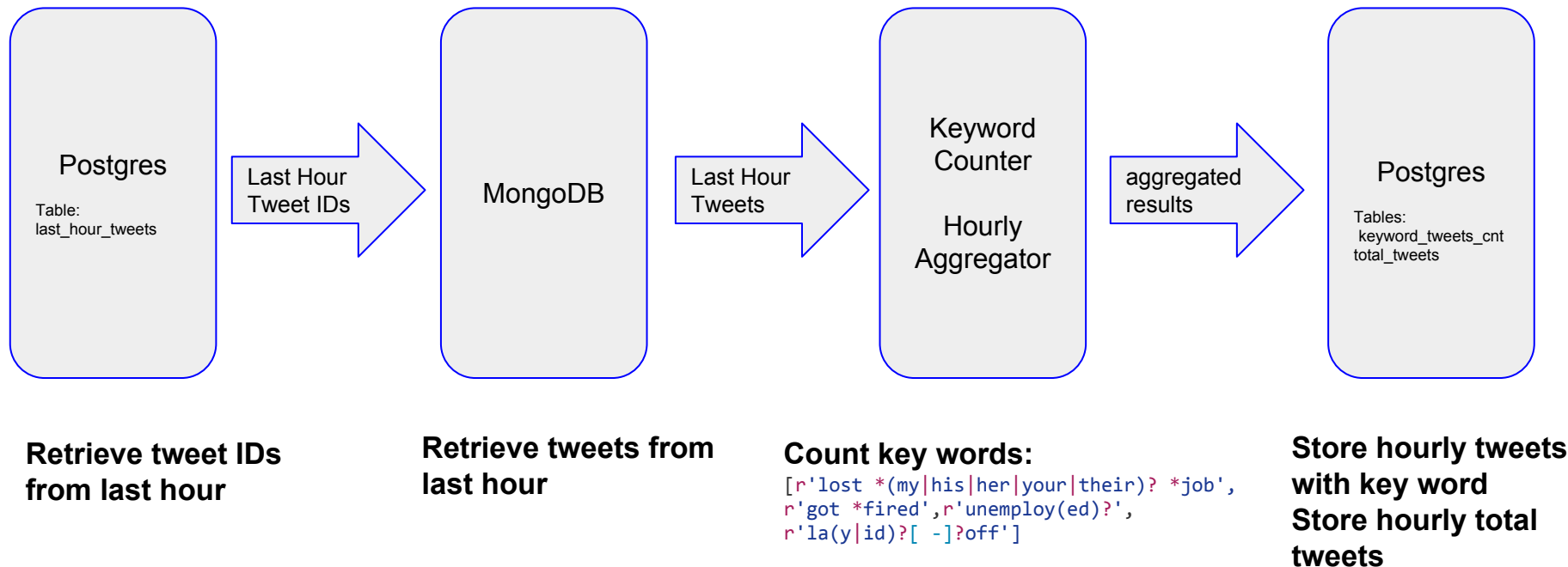




[illegible]

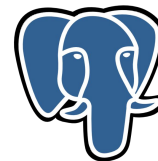
# Data Processing- Twitter

## Hourly Batch Job



# Claims and S&P500 Data

Quandl



PostgreSQL

```
w205project=> SELECT * FROM claims ORDER BY date DESC;
```

date	initial_claims
2015-12-05	384481
2015-11-28	262628
2015-11-21	305424
2015-11-14	264816
2015-11-07	291097
2015-10-31	258436
2015-10-24	245360
2015-10-17	232860
2015-10-10	256522
2015-10-03	227176
2015-09-26	215116
2015-09-19	219342
2015-09-12	198903
2015-09-05	232507
2015-08-29	230079
2015-08-22	226649
2015-08-15	229251
2015-08-08	239326
2015-08-01	224104



# Serving

## **PostgreSQL**

SQL users can query and get keyword trends  
Join Twitter, S&P500, and claims tables

## **Tableau**

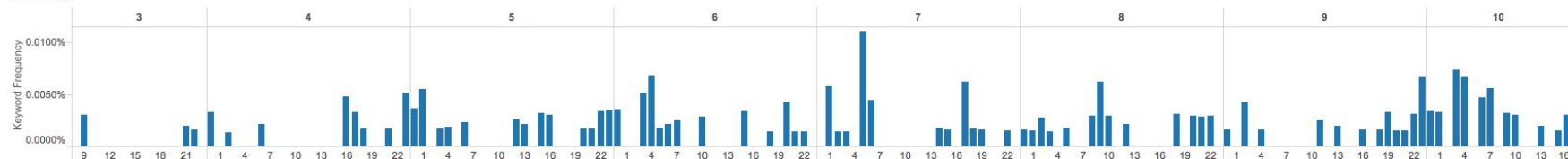
Dashboard for quick overview

## **MongoDB**

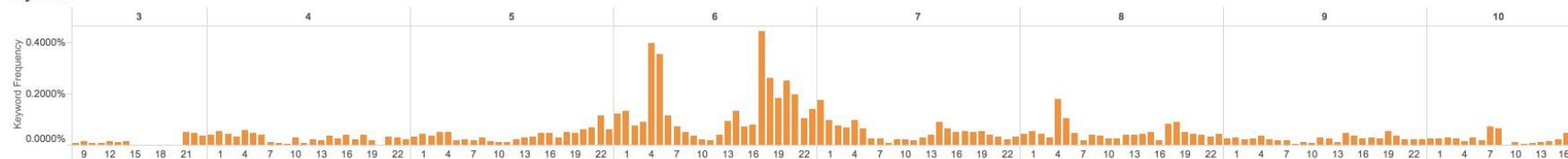
Use MongoDB for other specific queries (sort by # of hashtags)

# Twitter Keyword Frequencies (Per Hour, 12/3-12/10)

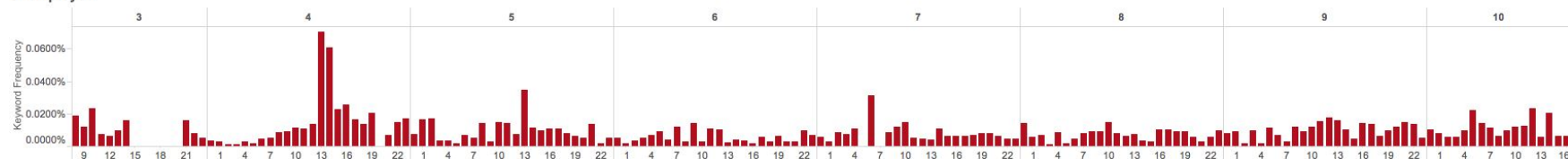
## Got Fired



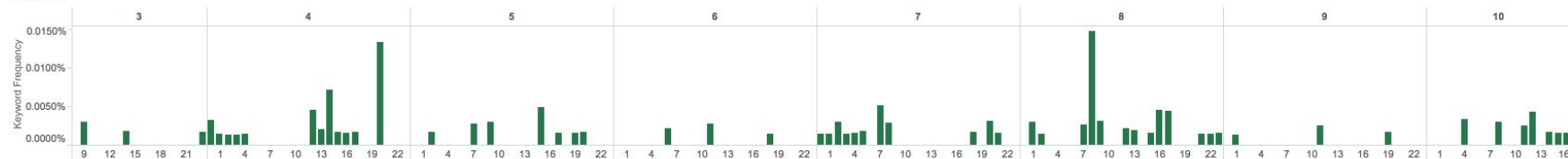
## Layed Off



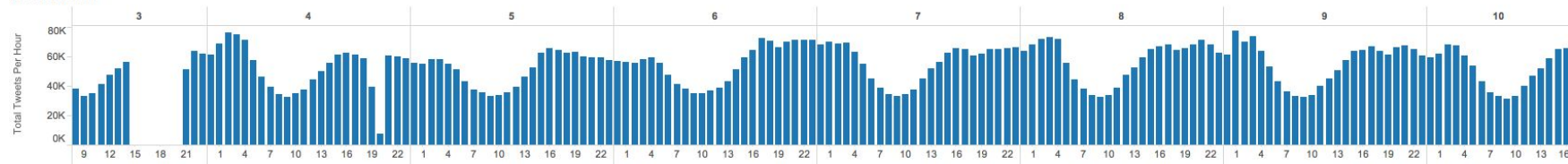
## Unemployed



## Lost Job

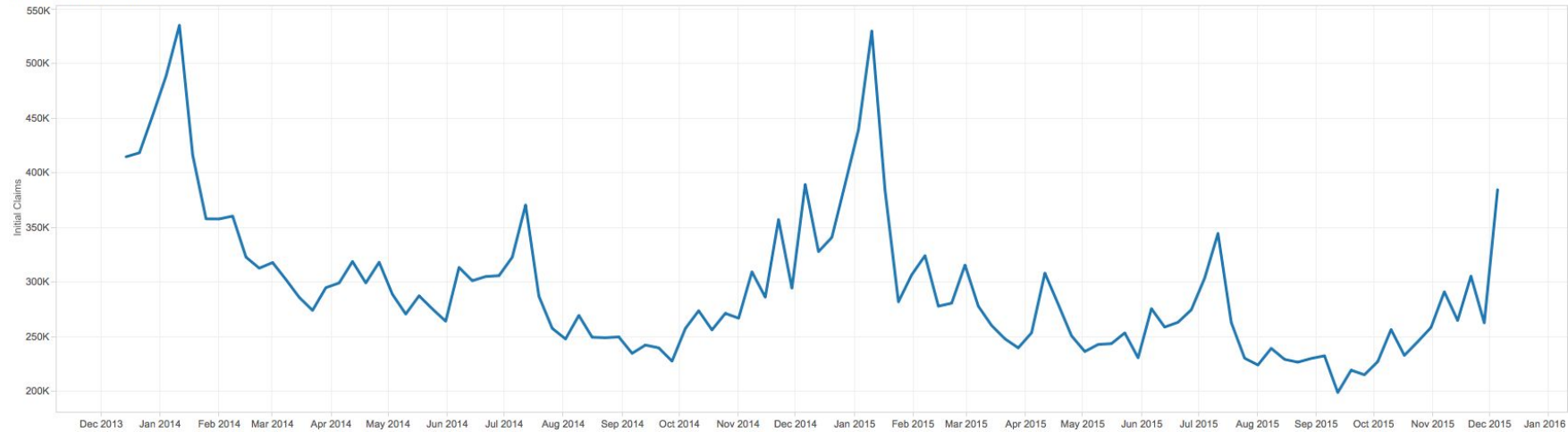


## Total Tweets

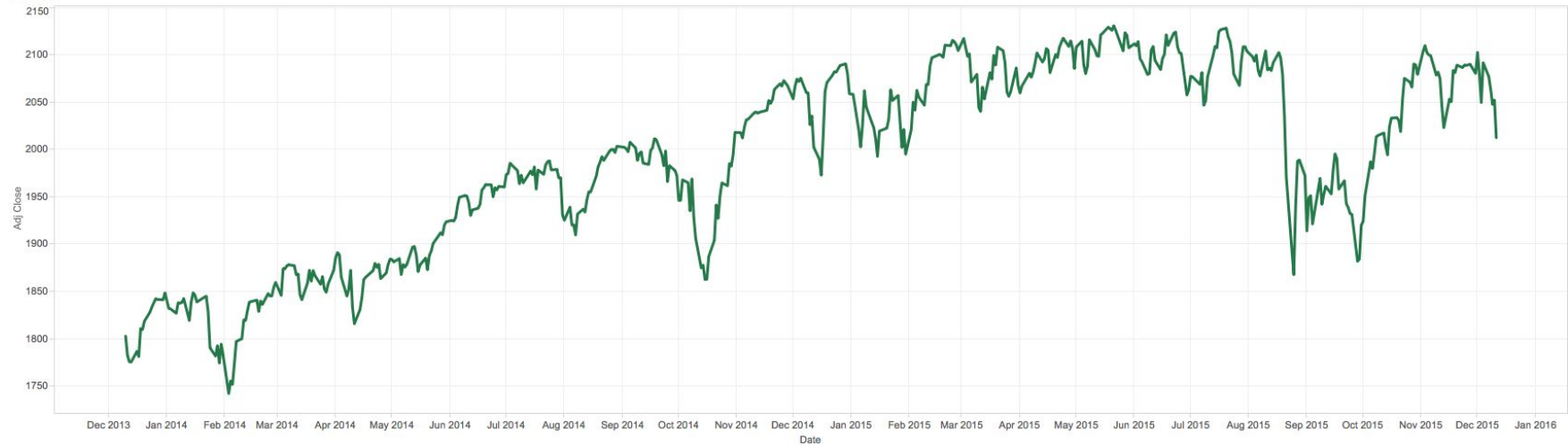


## Economic Indicators

### Initial Jobless Claims



### S&P500



# Spark Streaming



kafka

*tweet*

Spark  
Streaming

*json parsing*  
*'text' extraction*  
*filter for keywords*

*"unemployed", 3*  
*"fired", 2*  
*...*  
*total, 100*

*search terms = [*  
*'iphone',*  
*'ipad',*  
*'samsung',*  
*'android']*

```
Time: 2015-12-15 00:31:30
heavy duty hybrid rugged hard phone case cover for iphone 5c c+stylus+film us - bid now! o- https://t.co/ubhhcoukjk https://t.co/polulyavb7
rt https://t.co/nvyriVbjq bell/virgin iphone 5s 32g mint https://t.co/frgx6dogoj
luxury original leather flip cover card wallet case for apple iphone 6 6s plus - bid now! - https://t.co/ubheu1oyib https://t.co/u0mnhwchf
https://t.co/tnevp83kw #7443 apple iphone 5c 16gb "factory unlocked" 4g lte smartphone https://t.co/cyjbbbr15e
```

```
Time: 2015-12-15 00:31:30
```

```
4
```

```
Time: 2015-12-15 00:31:30
```

```
Time: 2015-12-15 00:31:30
```

```
0
```

```
Time: 2015-12-15 00:31:30
```

```
samsung appealing $548m patent infringement bill in apple case https://t.co/uucxnyw27
samsung galaxy s7 to pack in pressure sensitive display, improved low-light ... - firstpost https://t.co/ic6altlliu #tech
```

```
Time: 2015-12-15 00:31:30
```

```
2
```

```
Time: 2015-12-15 00:31:30
```

```
i've collected 17,200 gold coins! https://t.co/nmbdq0chv #android, #androidgames, #gameinsight
```

```
Time: 2015-12-15 00:31:30
```

```
1
```

```
Time: 2015-12-15 00:31:30
```

```
497
```

# Next Steps

## **Machine Learning**

Continue collecting tweets as training data

Predict jobless claims

## **Incorporate More Data**

NYTimes articles, weather data