**Big Data Management and Analysis**
**Spring 2017**

# PULSE OF NEW YORK TIMES

Alexey Kalinin, Enrique Gonzalez, Priyanshi Singh,
Shalmali Kulkarni, Sunny Kulkarni

# AGENDA

- ❖ Introduction

- ❖ What's Popular

- ❖ Project Objectives

- ❖ Data Set

- ❖ Methodology

- ❖ Sentiment Analysis

- ❖ Big Data Challenges

- ❖ Descriptive Statistics

- ❖ Visualization
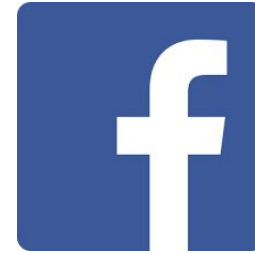
- ❖ Next Steps

- ❖ Appendix (scripts and job log)

# INTRODUCTION



- **Sentiment Analysis** identifying positive and negative opinions, emotions and evaluations in text.

# WHAT'S POPULAR...

Popular Google Search

Popular on Facebook

Popular on Twitter

Popular Emailed

# PROJECT OBJECTIVES

- ❖ Understand NYTimes trends on Twitter across US

- ❖ Reveal Spatial clusters for the tweeted subjects

- ❖ Sentiment analysis of most popular NY Times articles tweets

- ❖ Visualizations

# DATASETS

- **Twitter**

Geocoded tweets

- **NYTimes Articles**

Popular Articles

- **NLTK**

Sentiment Word List

# METHODOLOGY

**Machine Learning Pipeline**

**Sentiment Analysis using NLTK**

# SENTIMENT ANALYSIS

# BIG DATA CHALLENGES

❖ **Scalability**

❖ **Data Volume & Quality**
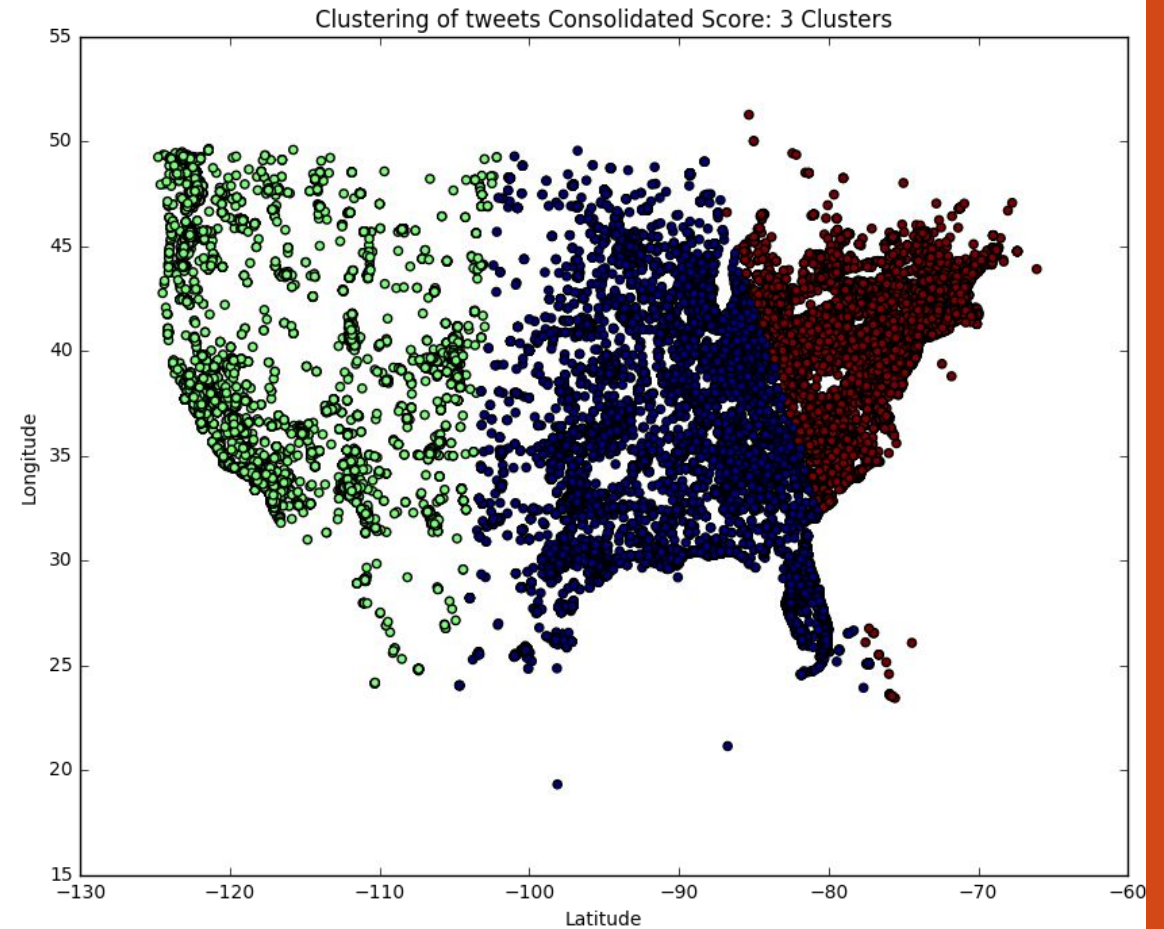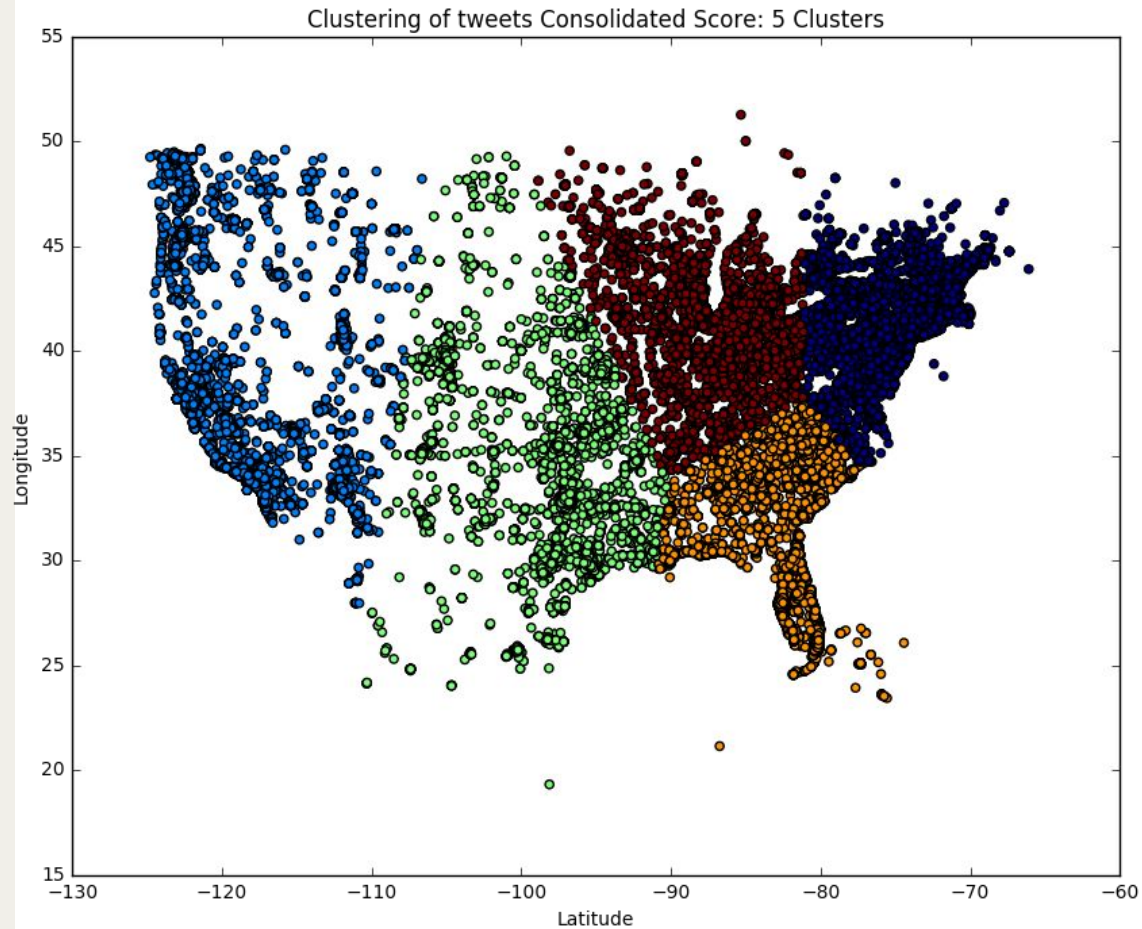
❖ **Data Structure**

# DESCRIPTIVE STATISTICS

Time frame:  April 28th to May 3rd 2017

No. of Tweets - approx. 600,000 records (3.2 GB)

No.of Geocoded Tweets - 100,000 records

NYTimes - 10 popular Articles per day

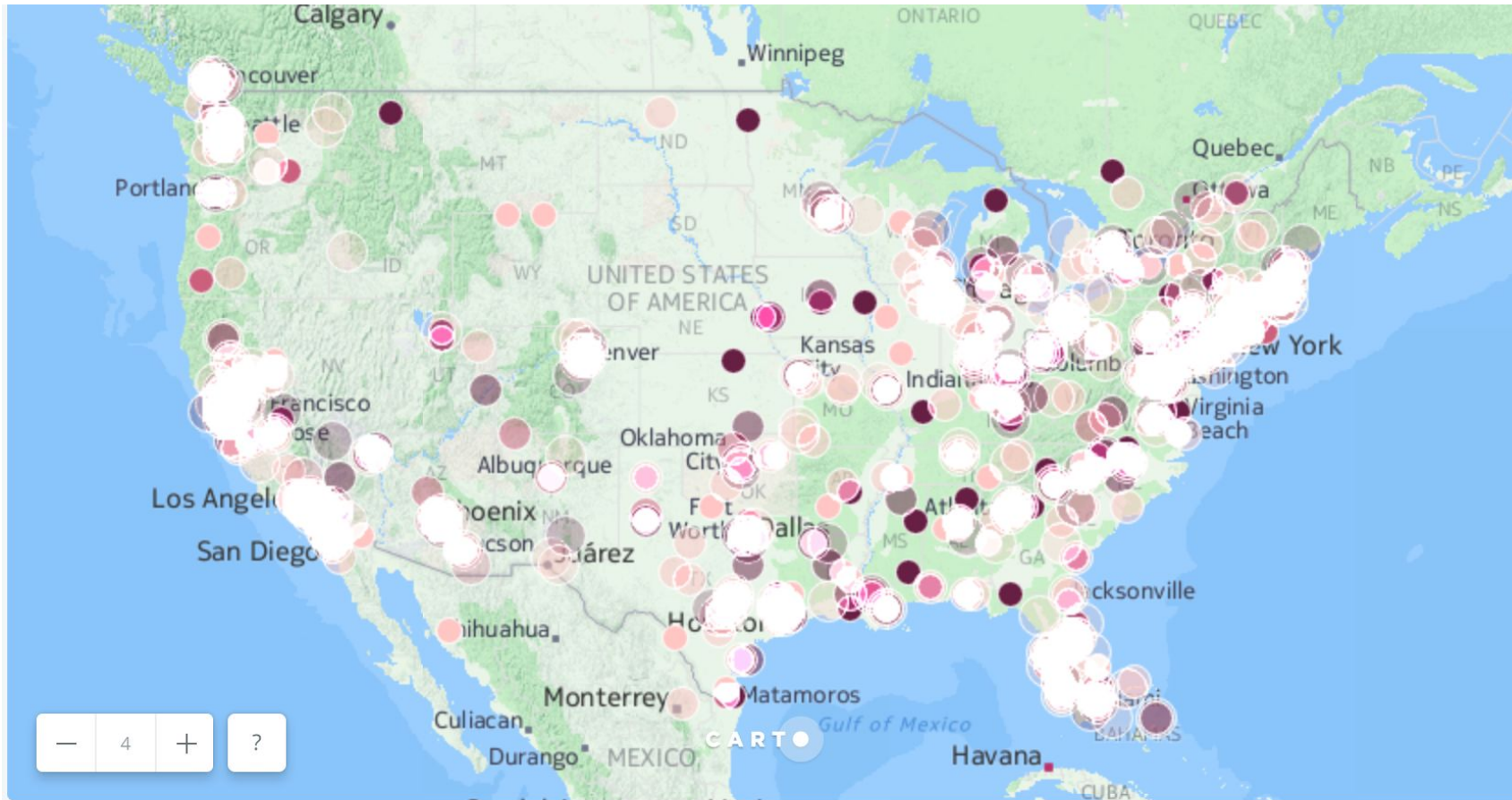# VISUALIZATION - SPATIAL CLUSTERING



**Clusters found after spatial clustering**
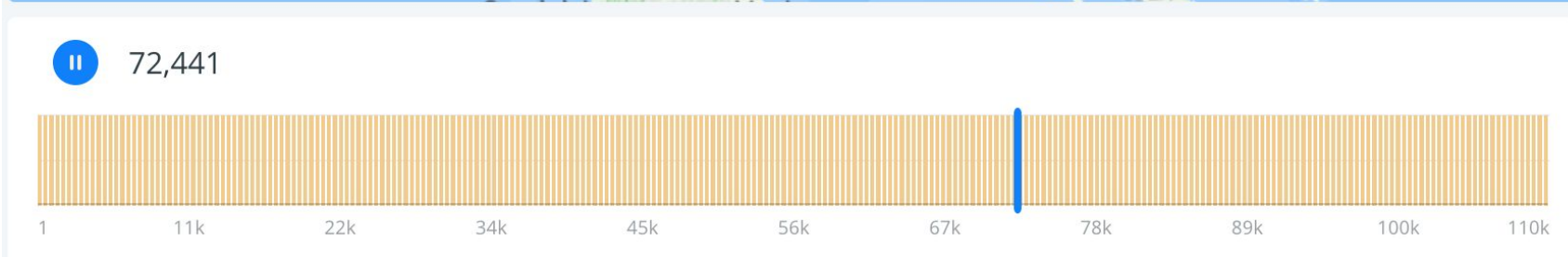
# VIS. - SENTIMENTS ACROSS USA



**Sentiments attached to the tweets across US with Time**

Real Time

72,441

← **Timeline**

# TRENDING TOPICS



28th April, 2017

1st May, 2017

# PROCESSING TIME



Data Size: 3.23 GB

Processing Time : 2m 30s

# NEXT STEPS…

Process Pipeline On Live Data Using Spark Streaming.



Spark.MLlib to Train Models and Perform Prediction

# APPENDIX – Spark Job Log

**Spark** 1.6.0 | Jobs | Stages | Storage | Environment | Executors | twitter_US_v2.py applic

## Executors

### Summary

| | RDD Blocks | Storage Memory | Disk Used | Cores | Active Tasks | Failed Tasks | Complete Tasks | Total Tasks | Task Time (GC Time) | Input | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Active(3) | 0 | 0.0 B / 1590.8 MB | 0.0 B | 2 | 0 | 0 | 2 | 2 | 7.1 s (132 ms) | 0.0 B | 0.0 B | 0.0 B |
| Dead(24) | 0 | 0.0 B / 12.4 GB | 0.0 B | 24 | 0 | 0 | 148 | 148 | 8.4 m (16.5 s) | 147.0 MB | 69.3 KB | 72.3 KB |
| Total(27) | 0 | 0.0 B / 14.0 GB | 0.0 B | 26 | 0 | 0 | 150 | 150 | 8.5 m (16.6 s) | 147.0 MB | 69.3 KB | 72.3 KB |

### Executors

| Executor ID | Address | Status | RDD Blocks | Storage Memory | Disk Used | Cores | Active Tasks | Failed Tasks | Complete Tasks | Total Tasks | Task Time (GC Time) | Input | Shuffle Read | Shuffle Write | Logs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | compute-2-1.local:50791 | Dead | 0 | 0.0 B / 530.3 MB | 0.0 B | 1 | 0 | 0 | 5 | 5 | 16.9 s (556 ms) | 256.0 KB | 0.0 B | 0.0 B | stdout stderr |
| 10 | compute-1-5.local:32974 | Dead | 0 | 0.0 B / 530.3 MB | 0.0 B | 1 | 0 | 0 | 10 | 10 | 28.8 s (1.2 s) | 576.0 KB | 0.0 B | 3.0 KB | stdout stderr |
| 11 | compute-2-12.local:33388 | Dead | 0 | 0.0 B / 530.3 MB | 0.0 B | 1 | 0 | 0 | 4 | 4 | 16.0 s (479 ms) | 192.0 KB | 0.0 B | 3.0 KB | stdout stderr |
| 12 | compute-3-2.local:49479 | Dead | 0 | 0.0 B / 530.3 MB | 0.0 B | 1 | 0 | 0 | 6 | 6 | 21.2 s (678 ms) | 320.0 KB | 0.0 B | 3.0 KB | stdout stderr |

# QUESTIONS?

---

# THANK YOU!