

# **Big-Data Systems and Intelligent Analytics**

## **Twitter based Movie Ratings**

**Team 4:**  
**Samir Sharan**  
**Harshit Shah**  
**Jeevan Reddy**

# Goal

The goal of this project is to provide a user-sourced and always up-to-date movie ratings. The ratings will be updated as much as possible by incorporating data from the newest tweets available.

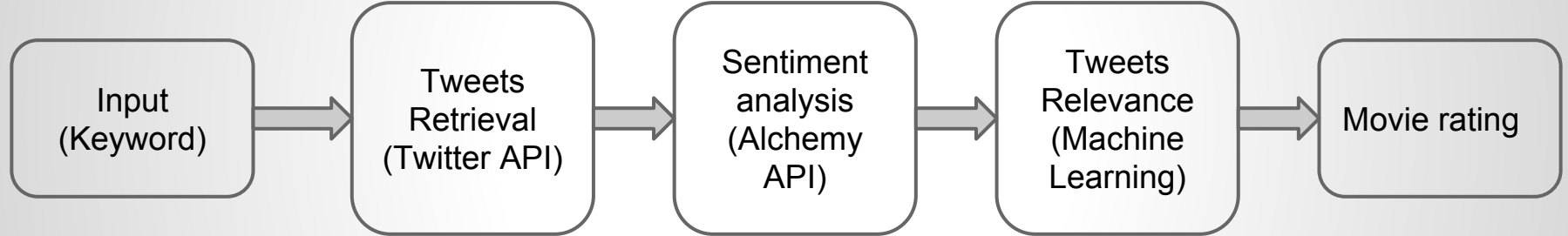
# Why Twitter?

- Twitter has 302 million active users
- 500 million tweets per day
- Tweets reflect user's real thoughts and opinions
- Real time information
- Can capture general audiences opinion

# Methods:

- Alchemy API
- AWS Sentiment Analysis

# Flow



# Tweet Sentiment



```
{
  "verified": false,
  "screen_name": "hiyaimdemi",
  "RT_count": 0,
  "text": "southpaw is such a good film\u263a\u2013\u263a",
  "coordinates": null,
  "sentiment": "positive",
  "language": "en",
  "score": 0.629288,
  "location": null,
  "time": "Wed Aug 19 21:30:24 +0000 2015",
  "movie": "Southpaw",
  "user_followers_cnt": 1078,
  "id": 634115221662134272,
  "fav_count": 1
}
```

- **Score:** Degree of sentiment

# Tweet Relevance



→

```
{"verified": false,  
"screen_name": "hiyaimdemi",  
"RT_count": 0,  
"text": "southpaw is such a good film\u2639\ufe0f",  
"coordinates": null,  
"sentiment": "positive",  
"language": "en",  
"score": 0.629288,  
"location": null,  
"time": "Wed Aug 19 21:30:24 +0000 2015",  
"movie": "Southpaw",  
"user_followers_cnt": 1078,  
"id": 634115221662134272,  
"fav_count": 1}
```

**Relevance Score:**  $W1 * RT\_count + W2 * fav\_count + W3 * user\_followers\_count + W4 * verified$

# Machine Learning

- Train on a set of 21,000 tweets
  - features - RT\_count, fav\_count, user\_follower\_count, verified
  - label - Relevance score
- Test on 14,000 tweets
  - Predict Relevance Score
  - Calculate Mean Squared Error

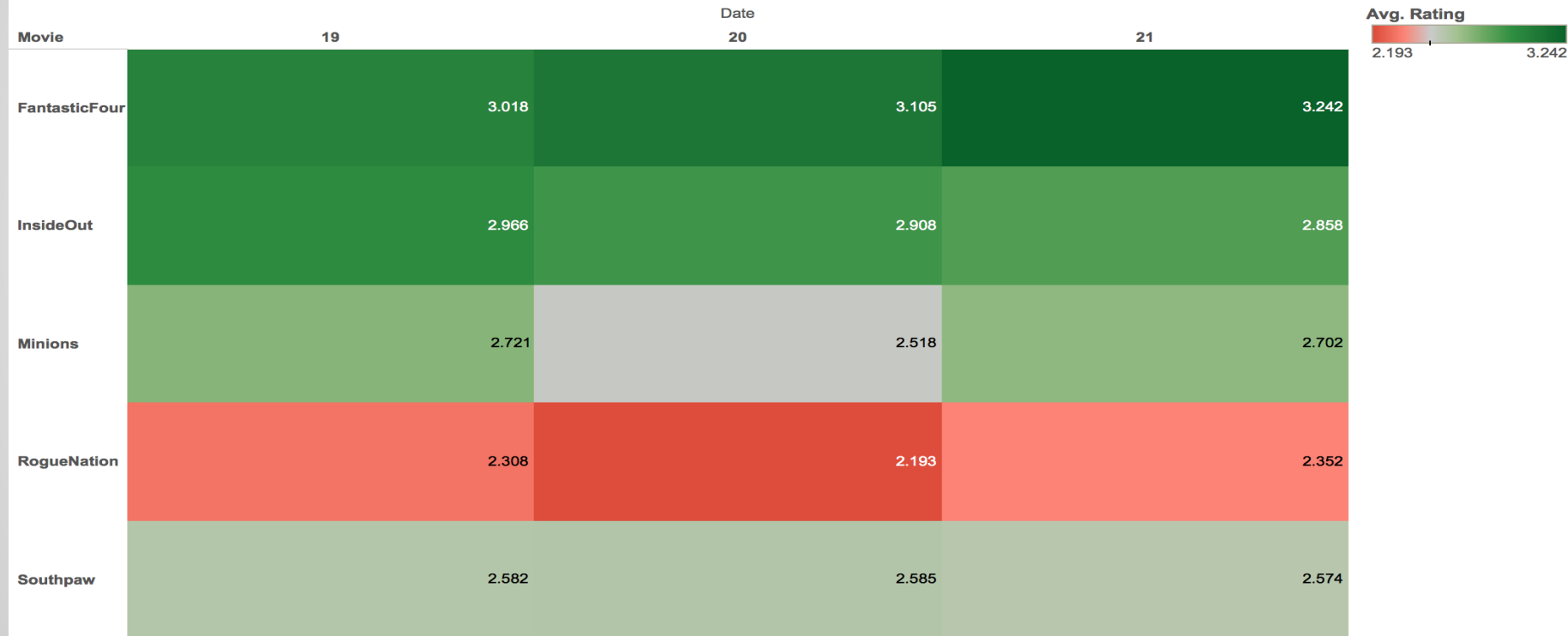


# Rating

- Depends on Sentiment Score (SS) and Relevance Score (RS)
- $R = (SS * RS) \rightarrow \text{scale on } 5 \rightarrow \text{add to } 2.5$
- $R = [(0.1632) * 5] + 2.5$
- $R = 3.316$

# Rating

## Day Wise Rating



Average of Rating broken down by Date Day vs. Movie. Color shows average of Rating. The marks are labeled by average of Rating. The view is filtered on Movie, which keeps FantasticFour, InsideOut, Minions, RogueNation and Southpaw.

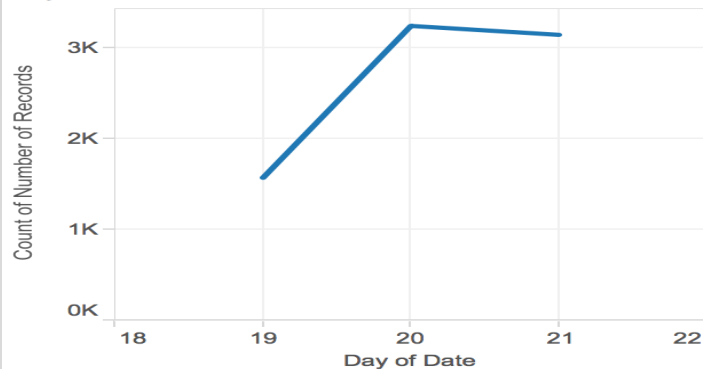
# Rating

## Average Rating Overall

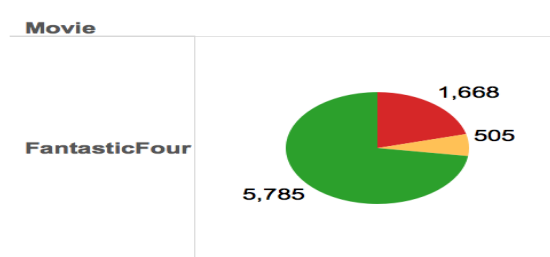


# Dashboard

## Days vs Tweets



## Movie vs Sentiments



### Movie

- ☒ FantasticFour
- ☐ InsideOut
- ☐ Minions
- ☐ RogueNation
- ☐ Southpaw

### Movie

- ☒ FantasticFour

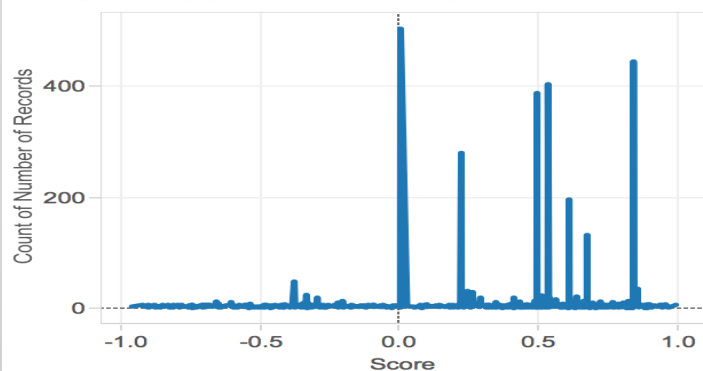
### Sentiment

- ☒ negative
- ☐ neutral
- ☐ positive

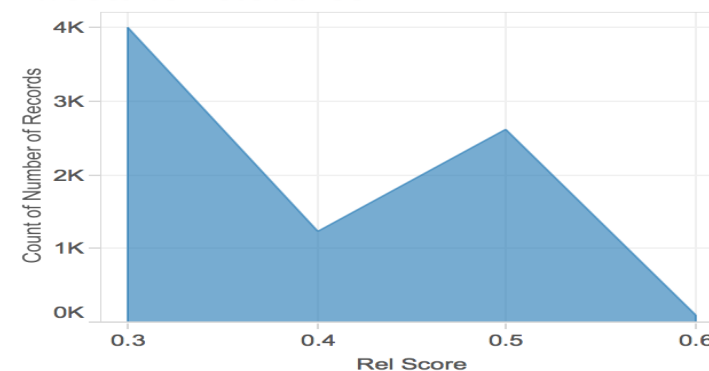
### Count of Number of Rec..

7,958

## Tweets vs Sentiment Score



## Tweets vs Relevance



# AWS Sentiment Analysis

- Create an Amazon S3 Bucket
- Collect and Store the data
- Build the mapper
- Create an Amazon EMR cluster
- Examine the output

# Steps in EMR:

Add Step

×

Step type


Streaming program

Name\*

Streaming program


Mapper\*

sentiment.py

 S3 location of the map function or the name of the Hadoop streaming command to run.


Reducer\*

aggregate

 S3 location of the reduce function or the name of the Hadoop streaming command to run.


Input S3 location\*

s3://aws-logs-202509000957-us-east-1/input/  
s3://<bucket-name>/<folder>/



Output S3 location\*

s3://aws-logs-202509000957-us-east-1/final\_output/  
s3://<bucket-name>/<folder>/




Arguments

```
-files  
s3://awsdocs/gettingstarted/latest/sentiment/classifier.p#classifier.p,s3://aws-logs-202509000957-us-east-1/mapper/sentiment.py
```

Action on failure

Continue

 What to do if the step fails.

Cancel

Add

## S3 Bucket:

[Upload](#) [Create Folder](#) [Actions](#)

[All Buckets](#) / [aws-logs-202509000957-us-east-1](#)

	Name
<input type="checkbox"/>	<a href="#">elasticmapreduce</a>
<input type="checkbox"/>	<a href="#">input</a>
<input type="checkbox"/>	<a href="#">mapper</a>

## Output:

	Name	Storage class	Size	Last modified
<input type="checkbox"/>	<a href="#">_SUCCESS</a>	Standard	0 bytes	Mon Aug 17 19:53:03 GMT-400 2015
<input type="checkbox"/>	<a href="#">part-00000</a>	Standard	0 bytes	Mon Aug 17 19:52:53 GMT-400 2015
<input type="checkbox"/>	<a href="#">part-00001</a>	Standard	0 bytes	Mon Aug 17 19:52:54 GMT-400 2015
<input type="checkbox"/>	<a href="#">part-00002</a>	Standard	23 bytes	Mon Aug 17 19:52:55 GMT-400 2015
<input type="checkbox"/>	<a href="#">part-00003</a>	Standard	0 bytes	Mon Aug 17 19:52:56 GMT-400 2015
<input type="checkbox"/>	<a href="#">part-00004</a>	Standard	15 bytes	Mon Aug 17 19:52:58 GMT-400 2015
<input type="checkbox"/>	<a href="#">part-00005</a>	Standard	0 bytes	Mon Aug 17 19:53:01 GMT-400 2015
<input type="checkbox"/>	<a href="#">part-00006</a>	Standard	0 bytes	Mon Aug 17 19:53:02 GMT-400 2015

No match: 12

minions: positive 364

# References

- Alchemy API: <http://www.alchemyapi.com/products/alchemylanguage/sentiment-analysis>
- AWS Sentiment Analysis: <http://docs.aws.amazon.com/gettingstarted/latest/emr/getting-started-emr-sentiment-tutorial.html>



**Thank You!**