# Tweets Sentiment Analysis of 2016 U.S. Presidential Election
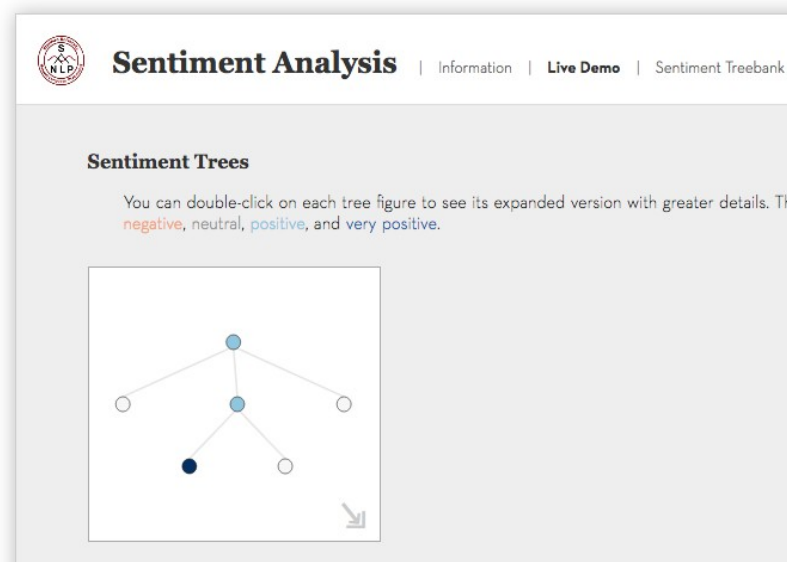
**Bin Li**
**Computer Science Department**
**New York University**

# What is sentimental analysis?

- It is a special case of text mining generally focused on opinion polarity

- The goal is to determine if the text opinion is positive, negative or neutral

# Why sentimental analysis?

- Microblogging has been a central place for people to express their thoughts and opinions on political parties and candidates

- Provide investment guidance

- A company wants to know the reviews of their products. (restaurant reviews)

http://sentiment.vivekn.com/

http://nlp.stanford.edu:8080/sentiment/rntnDemo.html

# Opinion Polarity

**Conor McGregor**
@TheNotoriousMMA    4h

I'm so fresh and so driven

↩    ♺ 3K    ♥ 13K    •••

**Positive**

**Neutral**

**World Surf League**
@wsl    Oct 10

Opening weekend on the North Shore... #Hawaii
Video by Eric Sterman

↩    ♺ 464    ♥ 775    •••

**Negative**

**The Guardian** ✓ @guardian · 6h

Barack Obama: I made bad decisions in my youth – video | US news | The Guardian

# Problem Statement

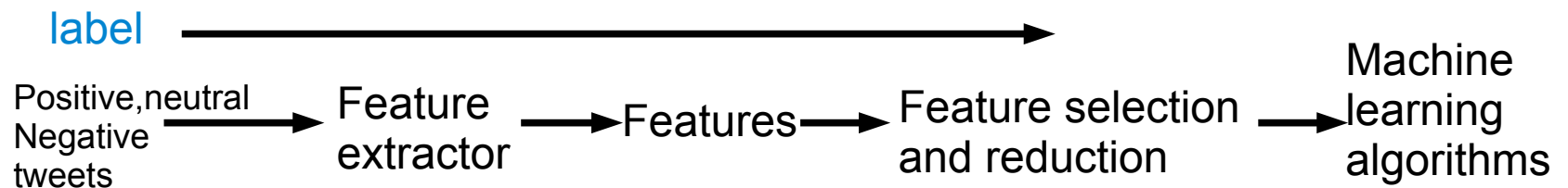- Given a tweet, determine if this tweet is positive, negative or neutral sentiment to a candidate.

# Design Diagram

Tweets collection → Aggregate tweets to candidates → Data cleaning

- Tweets Collection: tweets related to election through Twitter APIs

- Data Pre-processing: tokenization, feature reduction (remove non-english words, URLs, stop words and quotation, stemming, repeated characters etc.)

- Aggregate tweets to candidate: map tweets to the corresponding candidate

**Training**

label →

Positive, neutral Negative tweets → Feature extractor → Features → Feature selection and reduction → Machine learning algorithms

**Prediction**

tweets → Feature extractor → Features → Prediction model → label

# Data understanding

*Tweets to be labeled*

How? Keyword search using Twitter Search API and Stream API. On the second and last presidential debate date

Volume: 9.6 million tweets collected and hosted on HPC cluster

Type: highly unstructured text data

# Data understanding

## Training dataset

10,000 labeled tweets, Combination of the following two dataset:

- Stanford sentiment140

    Label classified based on emoji, no neutral labeled data

```
"4","2192559582","Tue Jun 16 07:12:59 PD 2009","NO_QUERY","yogilo","@EliciaKoay
you went overboard for the girl's birthday again "
```

- Sanders-twitter

    Manually labelled by human whose english is very well.
    Based on specific topics, IT company like apple, google

```
"apple", "neutral","125828984293425152","Mon Oct 17 07:02:22 +0000 2011","ok it
is back @iphone @apple"
```

My training dataset includes positive and negative tweets from sentiment140 and neutral tweets from Sanders.

# Data Preparation

- Remove irrelevant info such as URLs, quotes, citations and numbers

- Tokenization using tweet-preprocessor

- Remove punctuation

- Remove stop words

- Apply porter stemmer

- Finally, remove any word whose length is smaller than 3

For tweet,

rt @weneedfeminlsm: the way donald trump treats women is disgusting and repulsive. nothing about this is ok. https://t.co/lmbjgopgnt

After preprocessing,

way donald trump treat women disgust repuls noth ok

Do all above steps for both training tweets and tweets to be labeled.

# Data Preparation

Implement Naive Bayes algorithm

Text model: unigram or bigram

n*m

| | f1 | f2 | f3 | . | . | . | . | . | . | . | . | . | . | fm-1 | fm |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| t1 | 2 | 4 | | | | | | 7 | | | | | | | |
| t2 | | | 6 | | | | | | | 6 | | | | | |
| t3 | | 4 | | | 6 | | | 3 | 9 | | e | | | | |
| t4 | | | | | | | | | 4 | | | | | | |
| t5 | 2 | | | 7 | | 2 | 13 | | | 6 | | | | | 65 |
| . | 3 | | | 25 | | | | | 45 | | | | | | |
| . | | | | | | | 12 | | 8 | | | 5 | 16 | | |
| . | | 65 | | | 2 | 9 | | | | 4 | | | | | |
| . | | | 1 | | | 32 | 13 | | | | 78 | | | | |
| . | 87 | | | 45 | | | | 4 | | | | | | | |
| tn-1 | | | | | | 7 | | 23 | | | | | | | 65 |
| tn | | | | 6 | | | | | | | | 8 | | | |

# Model training

$$P(yi|x) = \frac{P(x|yi) * P(yi)}{P(x)}, x = a1, a2, ..., an$$

$$\log(P(yi|x)) = \log(P(x|yi)) + \log(P(yi)) - \log(P(x))$$

The goal is to get three prior probabilities that will be used by classification.

- class_prior (P(yi))
   ------------------
   the probability of each label

- feature_prior (P(x))
   -------------------
   the probability of each feature in all feature space

- label_feature_matrix (P(x | yi))
   ----------------------------
   the probability of each feature under certain label

# Model training

sample

| f1 | f2 | f3 | . | . | . | . | fn |
|----|----|----|---|---|---|---|----|

$$P(yi|x) = \frac{P(x|yi) * P(yi)}{P(x)}$$

$\log(P(x|yi))$      X                              X

|     | l1 | l2 | . | . | . | lm |
|-----|----|----|---|---|---|----|
| f1  | 2  | 4  |   |   |   |    |
| f2  |    |    | 6 |   |   |    |
| .   |    | 4  |   |   | 6 |    |
| .   |    |    |   |   |   |    |
| fn  |    |    |   |   |   |    |

| f1 |
|----|
| f2 |
| f3 |
$\log(P(x))$ | f4 |
| . |
| . |
| fn |

argmax
↓

| l1 | l2 | l3 | . | . | lm |
|----|----|----|---|---|----|

↑

↓

| l1 | l1 | l3 | . | . | lm |
|----|----|----|---|---|----|

−

| | + |

$\log(P(yi))$

| l1 | l1 | l3 | . | . | lm |
|----|----|----|---|---|----|

# Stats

**K-fold validation, 4216 features based on 500 training tweets using freq=2**

<span style="color:red">************Feature Freq=2********</span>

unigram stats

----------label:neutral------

tp: 86, tn: 301, fp: 4, fn: 4

accuracy:       0.980

precison:       0.956

recall:         0.956

----------------------------------

----------label:positive------

tp: 129, tn: 253, fp: 7, fn: 4

accuracy:       0.972

precison:       0.949

recall:         0.970

----------------------------------

----------label:negative------

tp: 121, tn: 280, fp: 2, fn: 5

accuracy:       0.983

precison:       0.984

recall:         0.960

----------------------------------

bigram stats

----------label:neutral------

tp: 110, tn: 157, fp: 25, fn: 0

accuracy:       0.914

precison:       0.815

recall:         1.000

----------------------------------

----------label:positive------

tp: 63, tn: 273, fp: 4, fn: 9

accuracy:       0.963

precison:       0.940

recall:         0.875

----------------------------------

----------label:negative------

tp: 59, tn: 300, fp: 0, fn: 20

accuracy:       0.947

precison:       1.000

recall:         0.747

----------------------------------

# Stats

**K-fold validation, 2612 features based on 500 training tweets using freq=3**

**\*\*\*\*\*\*\*\*\*\*\*\*Feature Freq=3\*\*\*\*\*\*\*\***

```
unigram stats
------------label:neutral-------
tp: 91, tn: 290, fp: 12, fn: 3
accuracy:    0.962
precison:    0.883
recall:      0.968
--------------------------------

------------label:positive------
tp: 128, tn: 266, fp: 2, fn: 11
accuracy:    0.968
precison:    0.985
recall:      0.921
--------------------------------

------------label:negative------
tp: 124, tn: 285, fp: 3, fn: 3
accuracy:    0.986
precison:    0.976
recall:      0.976
--------------------------------
```

```
bigram stats
----------label:neutral-------
tp: 116, tn: 106, fp: 22, fn: 3
accuracy:    0.899
precison:    0.841
recall:      0.975
------------------------------

----------label:positive-------
tp: 46, tn: 287, fp: 5, fn: 9
accuracy:    0.960
precison:    0.902
recall:      0.836
------------------------------

-----------label:negative------
tp: 43, tn: 310, fp: 1, fn: 16
accuracy:    0.954
precison:    0.977
recall:      0.729
------------------------------
```

**Performance: freq(2) > freq(3), unigram > bigram**

# Stats

Feature selection

Based on 10,000 training tweets, compare feature frequency

### neutral



### positive



### negative



Performance:
Freq=1 and freq=2 very close
so we choose freq=2

# PCA

| | f1 | f2 | f3 | . | . | . | . | . | fm |
|---|---|---|---|---|---|---|---|---|---|
| t1 | 2 | 4 | | | | | | | 7 |
| t2 | | | 6 | | | | | | |
| t3 | | 4 | | | 6 | | | 3 | 9 |
| t4 | | | | | | | | | |
| . | 2 | | | 7 | | 2 | 13 | | |
| . | 3 | | | 25 | | | | | |
| tn | | | | | | | | 12 | |

| | f1 | f2 | . | fk |
|---|---|---|---|---|
| t1 | | | | |
| t2 | | | | |
| t3 | | | | |
| t4 | | | | |
| . | | | | |
| . | | | | |
| tn | | | | |

$\longrightarrow$

n*m (m = 5000)

n*k (k = 20, 30, 40)

First k PCs

normalization $\longrightarrow$ Covariance matrix $\longrightarrow$ Eigen vectors

# PCA

# PCA

# PCA

Train(freq=2, model=unigram, PCs=20), Based on 10,000 training data

accuracy: 55.21%

|  | true 1 | true 0 | true −1 | class precision |
|---|---|---|---|---|
| pred. 1 | 539 | 148 | 190 | 61.46% |
| pred. 0 | 76 | 723 | 60 | 84.17% |
| pred. −1 | 2383 | 1371 | 3949 | 51.27% |
| class recall | 17.98% | 32.25% | 94.05% | |

Train(freq=2, model=unigram, PCs=30), based on 10,000 training data

accuracy: 55.44%

|  | true −1 | true 1 | true 0 | class precision |
|---|---|---|---|---|
| pred. −1 | 3914 | 2339 | 1355 | 51.45% |
| pred. 1 | 201 | 584 | 152 | 62.33% |
| pred. 0 | 84 | 75 | 735 | 82.21% |
| class recall | 93.21% | 19.48% | 32.78% | |

# PCA

Train(freq=2, model=bigram, PCs=20), Based on 10,000 training data

**accuracy: 48.83%**

|  | true 0 | true −1 | true 1 | class precision |
|---|---|---|---|---|
| pred. 0 | 389 | 14 | 17 | 92.62% |
| pred. −1 | 1792 | 4099 | 2860 | 46.84% |
| pred. 1 | 61 | 86 | 121 | 45.15% |
| class recall | 17.35% | 97.62% | 4.04% |  |

Train(freq=2, model=bigram, PCs=30), Based on 10,000 training data

**accuracy: 49.25%**

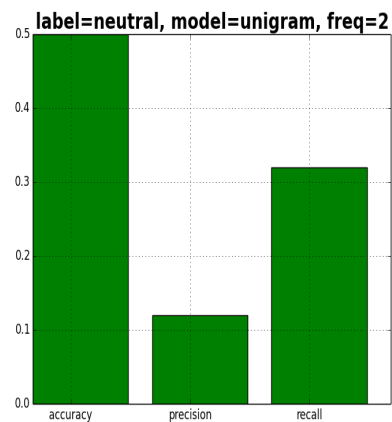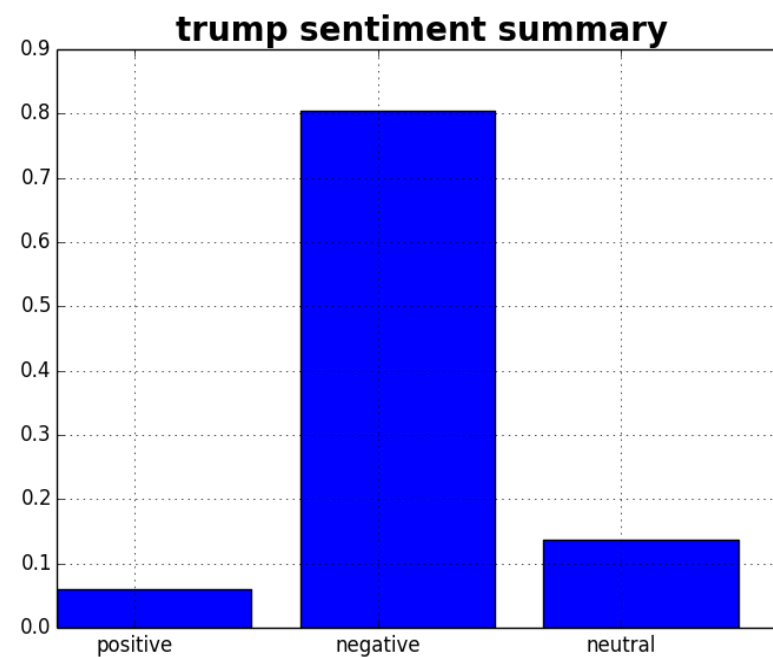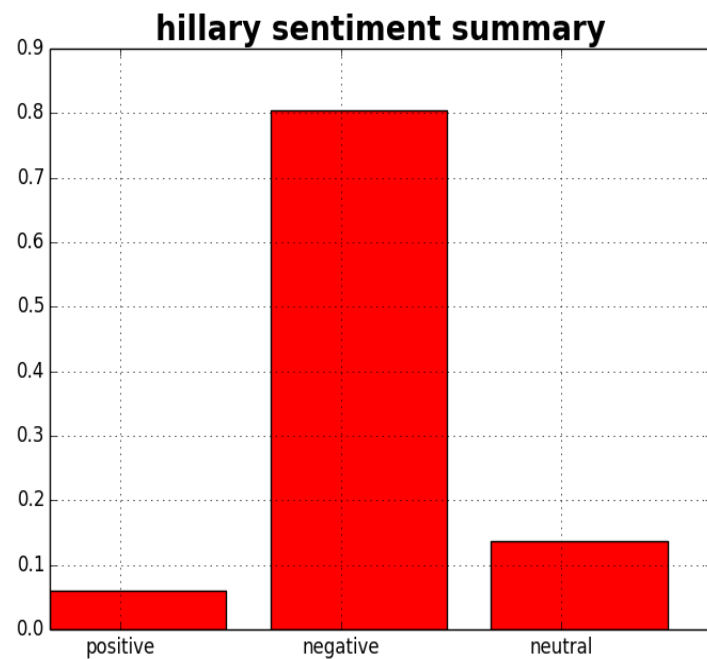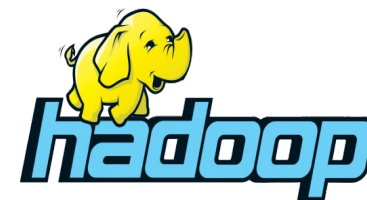|  | true 1 | true −1 | true 0 | class precision |
|---|---|---|---|---|
| pred. 1 | 111 | 20 | 3 | 82.84% |
| pred. −1 | 2869 | 4162 | 1863 | 46.80% |
| pred. 0 | 18 | 17 | 376 | 91.48% |
| class recall | 3.70% | 99.12% | 16.77% |  |

# Evaluation

Model=unigram, feature frequency=2

Trump, compare with baseline, based on 10972 tweets



Hillary, compare with baseline, based on 8941 tweets

Sentiment summary computed by Hadoop Map/Reduce, based on 9.6 million tweets

# Deployment

- Write tweets and their predicted sentiment to HBase
- Front end UI to use JavaScript to make REST call to HBase and virtualize tweet sentiment

hillari:2016-11-25 16:36:50_1657    column=data:sentiment, timestamp=1480136553022, value=positive

hillari:2016-11-25 16:36:50_1657    column=data:tweet, timestamp=1480136553022, value=rt @foxnews: .@judgenap:
new fbi docs show 'bribe offer to agents in #hillaryclinton email probe https://t.co/49cj9rgdun https://t.co/svfd\x0A

# Current standing

*Data preparation*

```
┌─────────────────────┐      ┌─────────────────────┐      ┌─────────────────────┐
│                     │      │                     │      │      Aggregate      │
│  Tweets collection  │ ───> │        Data         │ ───> │     tweets to       │
│                     │      │   pre-processing    │      │     candidates      │
└─────────────────────┘      └─────────────────────┘      └─────────────────────┘
```

- Collected 4,000,000 tweets using keywords search. Starting from the second debase. Data is hosted on NYU HPC clusters.

- Downloaded labelled tweets (1.6 million) as training data from Stanford Sentiment140

- Performed tokenization, removed stop words, URLs, repeated characters and quotation, stemming

# Acknowledgement

- Stanford sentiment140

- Sanders-twitter

- NYU HPC

helpful links:

http://www.laurentluce.com/posts/twitter-sentiment-analysis-using-python-and-nltk/

http://adilmoujahid.com/posts/2014/07/twitter-analytics/