# SENTIMENT ANALYSIS OF PRESIDENT TRUMP'S TWEETS

Akshay Ketkar
ahketka2@illinois.edu

# Table of Contents

# 1. INTRODUCTION

In today's world of social media and networking, there is a tremendous creation and exchange of information. It is much easier today to express one's opinions on the public platform than what it was in the era gone by. One such social media platform, **twitter.com** has become an important media through which people can express or shape opinions on several matters of social, political, national or international standing.

The project analyzes a dataset released on kaggle.com. The dataset consists of over 7000 of the current US President Mr. Donald Trump's tweets from 2015 to 2016. The project endeavors to run a sentiment analysis on the dataset to determine the proportion of positive and negative words apart from some other very frequent references using text analysis techniques.

# 2.LITERATURE REVIEW

## 1. Finding Similar Items in Structured Datasets

*Reference: IE531-Algorithms for Data Analysis- Lesson 3: Finding Similar Items in Structured Data sets by Prof. R.S. Sreenivas*

This methodology is used for finding similar items among collections. While most cases dealt with are in terms of measuring the similarity between texts, the procedures here can be extended to others as well. Multimedia such as images and audio would require a feature-extraction step, which then proceeds to analyze features on similar grounds

### Similarity measures:

The Jaccard similarity of sets S and T is SIM(S, T) = card(S∪T) / card(S∩T) , is the ratio of the number of common elements between two sets to the total number of elements in both sets. For this, the notion of shingling is introduced. That is, we slide a window of size k across the sentence and record what we see. The k-shingles are computed using a window of size k. The common practice, when it comes to text files is to use a value of k = 9. To compute the similarity between two text files, we compute Jaccard Similarity of their 9-shingle set. We use this technique to improve our computation of the proportion of positive and negative terms in President Trump's tweets.

## 2. Map-Reduce

*Reference: IE531-Algorithms for Data Analysis- Lesson 5: Basics of MapReduce and some related topics by Prof. R.S. Sreenivas*

Although an attempt has been made to get in the code in the Mapper and Reducer form, the execution was not successful on the Cloudera Emulator's Hadoop file system. Instead the codes were successfully run on Cloudera local machine using pipe, to stream twitter words from Mapper code to the Reducer code, whereby the Reducer code does all the analysis.

### MapReduce operations:

All **MapReduce** operations have two parts to it a mapper, which produces intermediate values in the form of (key; value) the output of the mapper is sorted based on the key values and presented to the reducer, which does the needful to complete the task. The mapper code and reducer code run on different nodes on different blocks of data.

# 3. PROJECT FLOW

**Step1**:

**Data Cleaning:** The twitter data set is cleaned for the removal of quotes and other special characters such as commas, hashes, question marks, exclamations and other characters such as smileys. All the tweets are converted to lower case and outputted to a file.

**Step2**:

**Mapper Code:** The Mapper code parses each line of the twitter text column from the file mentioned above. Then each of the tweets is split to generate words, which are streamed for the Reducer code.

**Step3:**

**Reducer Code:** The Reducer code takes the words streamed by the Mapper and does the following:

a) Creates a dictionary of the distinct words used and creating a list of streamed words

b) Using the stream list and dataset containing Positive and Negative words from

https://github.com/jeffreybreen/twitter-sentiment-analysis-tutorial-201107/tree/master/data/opinion-lexicon-English

, it creates a list of positive and negative words, creates 9-character shingles of positive, negative and twitter words.
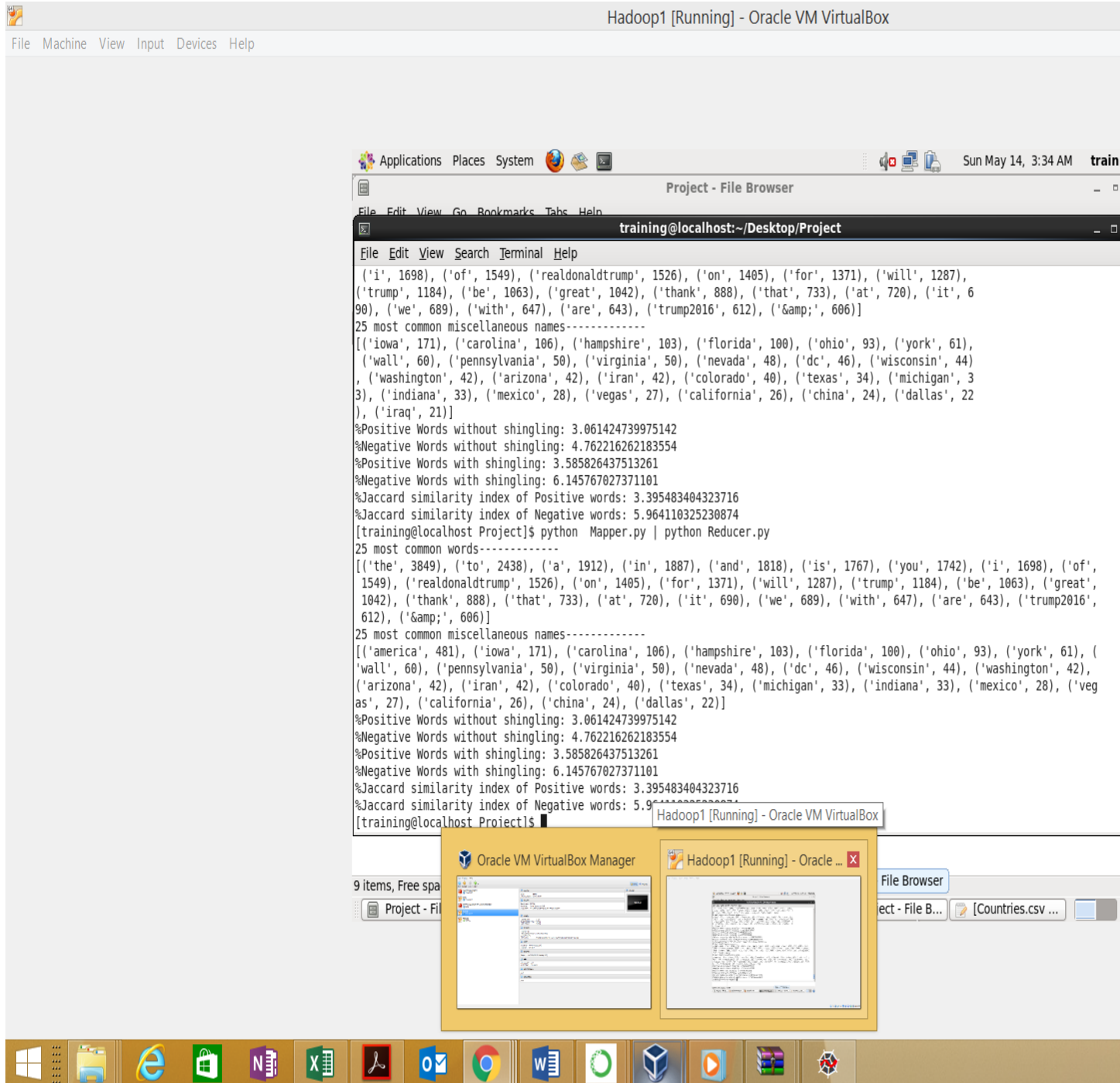
c) A dataset consisting mostly of country and US State names, but also a few terms such as 'wall' has been used to determine how frequently President Trump uses them.

d) The Jaccard similarity and proportion of Positive and Negative words is computed. Also, the dictionary of wordcounts and miscellaneous terms is sorted to reveal the top 25 most frequent entries.

e) The respective details are printed.

# 4. RESULTS

### 1. Screenshot:

**2. Interesting Results:**

**a) '25 most common words' and '25 most common miscellaneous terms:**

Among the 25 most common words, the words that stand out are **'great'(Count:1042) and 'thank'(Count:888) and 'trump'(Count:1184).** Some of the 'great' may refer to **'Let's make America great again!'** as the count for the term **'a(A)merica' is uttered 481** times as well. The utterance of the word **'thank'** may refer to President Trump's thank you for several compliments and good wishes.

Among American states, the state **Iowa** features **171** times. Among countries, the country **Iran** features **42 times.** Among other words, the word **'wall' features 60 times**. Some of the sentences consisting of 'wall' could point towards building a wall on the border with Mexico.

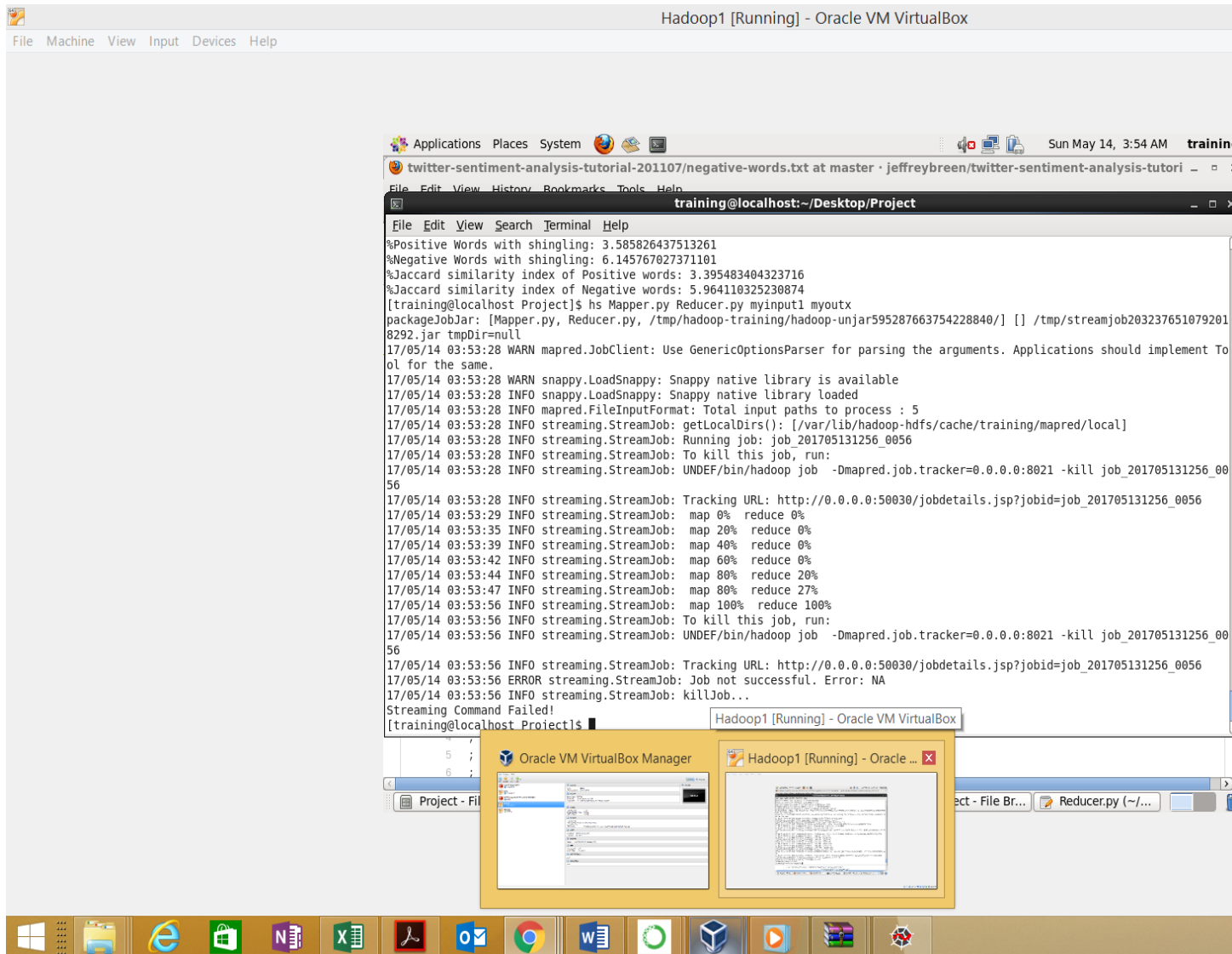**b) Results of analysis using Similarity Measures: Shingling and Jaccard Similarity index:**

**-** Taking **direct ratio** of the number of **Positive** and **Negative** words with the **total number of distinct twitter words** yields **3.06%** and **4.76%** words as Positive and Negative respectively.

- However, **using 9-character shingles** (as per convention), we get a slight improvement in the results, with the percentage of **Positive** and **Negative** word estimate increasing to **3.58%** and **6.15%.** The Jaccard similarity index is **3.39%** and **5.96%** for **Positive** and **Negative** words respectively.

# 5. CONCLUSION, HADOOP AND SCOPE FOR IMPROVEMENT

It can be concluded from the Sentiment Analysis that there is a consistently higher percentage of Negative words (With or without shingling), indicating a marginally higher negative sentiment than a positive sentiment.

**Hadoop Screenshot**



**Scope for Improvement:**

It was not possible to run the MapReduce on Cloudera emulator. If there is a way to directly connect to a Hadoop cluster and determine the reasons for the failure, there could be a successful run possible.

# 6. REFERENCES

1. https://www.kaggle.com/kingburrito666/better-donald-trump-tweets

2. https://github.com/jeffreybreen/twitter-sentiment-analysis-tutorial-201107/tree/master/data/opinion-lexicon-English

3. *IE531-Algorithms for Data Analysis- Lesson 3: Finding Similar Items in Structured Data sets by Prof. R.S. Sreenivas*

4. *IE531-Algorithms for Data Analysis- Lesson 5: Basics of MapReduce and some related topics by Prof. R.S. Sreenivas*