# Data Intensive Computing – Lab 2

Santosh Phani Vishnu Aita, UBID: saita

Harsha Kosta, UBID: harshako

## Data Extraction:

Data extraction is the most important step in any computing or machine learning problem. Even though it is distributed in decent amounts by various organization for the developers to use. It is generally messy and unorganized. Every data has some kind of gibberish and noisy information. Noise is a hindrance to using data directly after extracting from the API or any kind of source.

For lab 2, Raw data is needed to feed to the Map Reduce application which removes all unrequired information like stopwords, punctuations, symbols etc.

Types of data:

For this lab, data from Twitter using RtweetAPI and NYTimes articles using nytimesarticleapi was extracted. This step is quite hideous. Tweets data comprised of location, dates, usernames and much more information. Only Tweet  data had to be extracted. So tweet data was fetched and stored in separate files as it contains important information.

In case of NYTimes articles. The nytimes article API only provided JSON response which contained information like urls, dates, pages, location etc. Using the urls from the JSON response the article data was fetched.

Each url was passed to the BeautifulSoup (bs4), a package used in python for web scrapping and all the article data available in the paragraph tags were extracted and appended to form a larger data set.

We collected Tweets and NYTimes articles for Short-Span(one day) and Long-Span(one week). MR produced different results when visualized for weekly and daily data even though the key words were same.

The data from both Twitter and NYTimes was collected for three different keywords.

## Hadoop Map Reduce:

As we know, Hadoop is a collection of open-source software utilities that facilitate using a network of many computers to solve problems involving massive amounts of data and computation. It provides a software framework for distributed storage and processing of big data using the MapReduce programming model.

We implemented Map Reduce functionality using Python3 as programming language.

Data collected from Twitter and NYTimes is separately fed to the Mapper and each word in the data set is selected if it is not in the stopwords, all the punctuations and unwanted symbols are removed and all the words are stemmed to eliminated similar meaning words using nltk lemmatize.

The list of stopwords was extracted from the nltk package available in python.

The output from the mapper is emitted as <Word, Count> for each word and this is fed as input to reducer after shuffle and sort.
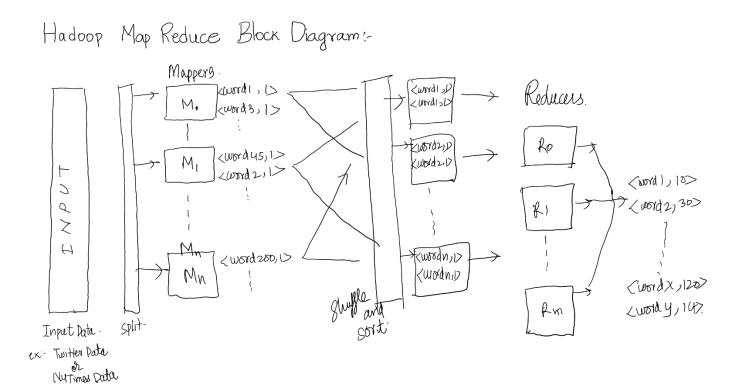
In shuffle and sort phase all the similar key words are paired together and sent to reduce as input.

The reducer aggregates or sums all the word count and emits the total number of occurrences of a word in the data set.
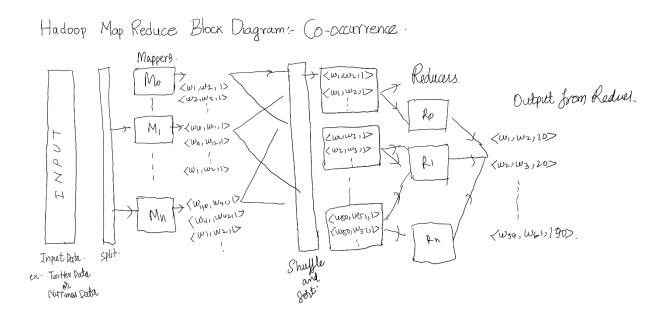
This is then fed to the d3.js to display as a word cloud on a webpage to visualize the data.

## Block Diagrams
Word Count MR Problem:

Co-occurrences MR Problem:

Hadoop Map Reduce Block Diagram:- Co-occurrence.

Mappers.

INPUT

M₀ → $\langle w_1, w_2, 1\rangle$
$\langle w_2, w_3, 1\rangle$

M₁ → $\langle w_{10}, w_{11}, 1\rangle$
$\langle w_{11}, w_{12}, 1\rangle$
$\langle w_1, w_2, 1\rangle$

Mₙ → $\langle w_{40}, w_{41}, 1\rangle$
$\langle w_{41}, w_{42}, 1\rangle$
$\langle w_1, w_2, 1\rangle$

Shuffle and Sort

$\langle w_1, w_2, 1\rangle$
$\langle w_1, w_2, 1\rangle$

$\langle w_2, w_3, 1\rangle$
$\langle w_2, w_3, 1\rangle$

$\langle w_{50}, w_{51}, 1\rangle$
$\langle w_{50}, w_{51}, 1\rangle$

Reducers.

R₀
R₁
Rₙ

Output from Reducer.

$\langle w_1, w_2, 20\rangle$
$\langle w_2, w_3, 20\rangle$
$\langle w_{59}, w_{61}, 190\rangle$.

Input Data.
ex:- Twitter Data
or
NYTimes Data

Split-

## Visualization using d3.js:

The output from Map Reduce for various key words and types of data like Twitter and NYTimes are plotted as word clouds using d3.js.

D3.js intakes <word_text, word_count> tuples stored in either a text file or a csv file as input. From this the word_text is the word which is supposed to be displayed in the word cloud and the word_count is used as the size of the word to be displayed on the web page.

The words having larger size mean that the particular word have comparatively more occurrence than other words. Lesser prevalent words are displayed in smaller sizes.

The two types of data NYTimes and Twitter resulted in different rampant words which is expected even though the key words were same because the are two different forms of media, one is social and other is professional new media.

There were many types of visualization mechanisms available in d3.js and are very widely used by various organization to visualize data.

We chose word cloud representation and it is really easy to understand the data by quantizing it effectively.

We chose the Avengers, Haunted and Psychopath as key words and extracted data from both Twitter and NYTimes articles and word counts, co-occurrences of words as pairs was evaluated using Hadoop map reduce environment. The output were fetched and stored as csv files and fed to d3.js draw method to produce word cloud.

The following is an example of comparison between Twitter and NYTimes article data for a keyword : Avengers and Span of data as week.



## Exploration:

There are unlimited possibilities as the data available is enormous, the companies like facebook can make smart emojis/ trends using the top words or showcase the popular words used by various audiences.

When a movie or a book is released the reviews data or posts on social media can be fetched and the word counts can be analyzed to check whether the movie has positive or negative ratings.

There are various possibilities in medical fields to determine commonly occurring molecules or gene data.

## Usage in Organization:

The Map Reduce is used by various organizations. A video game statistics company called Open Dota displays word clouds for each user using the data from the conversations he had while interacting with his team mates during all the video games.

The following is the word cloud generated by them for me:



## Directory Structure:

The following is the HDFS Structure for lab 2.

The following is the directory structure of the submission for lab 2.

- ⌄ 📁 lab2
  - 〉 📁 Saita_Lab2_Part1
  - ⌄ 📁 Saita_Lab2_Part2
    - 📁 D3_Files
    - 📁 DATA_EXTACTION_CODE
    - ⌄ 📁 MR_CODE
      - 📁 COOCCURENCE
      - 📁 WORDCOUNT
    - ⌄ 📁 MR_OUTPUT
      - 📁 COOCCURENCE
      - 📁 WORDCOUNT
    - ⌄ 📁 RAW_DATA
      - 📁 NYTIMES
      - 📁 TWITTER