# LAB-2

# DATA AGGREGATION, BIG DATA ANALYSIS AND VISUALIZATION

**Divya Srivastava**
Person number – **50290383**
UBID – **divyasri**
And
**Meghana Vasudeva**
Person number – **50290586**
UBID **- mvasudev**

## Table of Contents

**Abstract:**

In this lab, we expanded our skills in data exploration developed in Lab1 and enhanced them by adding big data analytics and visualization skills. This document describes Lab2: Data Aggregation, Big Data Analysis and Visualization, involves (i) data aggregation from more than one source using the APIs (Application programming interface) exposed by data sources, (ii) Applying classical big data analytic method of MapReduce to the unstructured data collected, (iii) store the data collected on WORM infrastructure Hadoop and (iii) building a visualization data product.
We have leveraged the data collection and exploratory data analysis skills developed in Lab1 to accomplish the goals of Lab2.
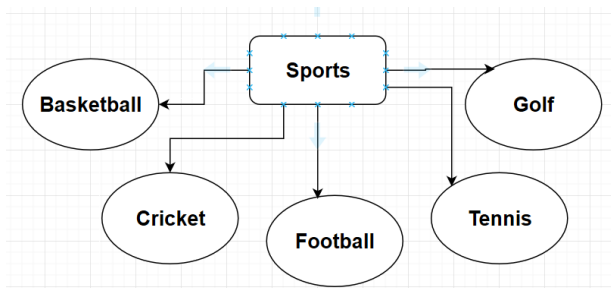
**Part – 1:**

**Data Collection**
An important and critical phase of the data-science process is data collection. Several organizations including the federal government (data.gov) have their data available to the public for various purposes. Social network applications such as Twitter and Facebook collect enormous amount of data contributed by their numerous and prolific user. An API or application programming interface is a standard, secure and programmatic access to data by an organization that owns the data.
We collected the data from the following three sources-
1.  Opinion-based social media in twitter
2.  Research data in New York Times
3.  Common crawl.

The above data was collected based on same set of keywords and topics. Then we processed the three data sets collected individually using classical big data methods and compared the outcomes using popular visualization methods.
We selected the following topic and sub topics for our analysis-



1.  <u>Twitter</u>
    Twitter data is unique from data shared by most other social platforms because it reflects information that users choose to share publicly. We registered for twitter API in our lab1 and used it for our lab2 as well. These APIs are required to register an application. By default, applications can only access public information on Twitter.
2.  <u>New York Times</u>
    We used the article search API from nytimesarticle to look up articles by keyword. We refined our search using various filters. We also used BeautifulSoup to extract the content. Beautiful Soup parses anything we give it, and does the tree traversal stuff for us. We can tell it "Find all the links", or "Find all the links of class externalLink", or "Find all the links whose urls match "foo.com", or "Find the table heading that's got bold text, then give me that text."

3. Common Crawl

   For extracting data from common crawl, we read the warc file that we downloaded from the common crawl website for March 2019 data. Crawling through it, we extracted the data from the links where language was English and matching keywords were found. We did this for all our keywords obtaining the data that we used for our lab.

**Cleaning**

The data obtained after data collection for twitter, New York times and common crawl through their respective APIs, was in the raw format. So our next step was cleaning for which we used a python script. We used natural language toolkit known as **NLTK** package to tokenize and tag the text and for identifying the named entities. We also used **regex** package for cleaning the data. This module provides regular expression matching operations. It is a sequence of character(s) mainly used to find and replace patterns in a string or file, so we used NLTK and regex both to filter the unwanted text from the data.

The text that we filtered out mainly contained the following unwanted elements-

1. Punctuations ( '!"#$%&\'()*+,-./:;⇔?@[\\]^_`{|}~')
2. Stop words ( I, me, my, myself, we, are etc)
3. Stemming words ( running becomes run and so on )
4. Tokenize (
5. Digits (0-9)
6. URL
7. Tags
8. Unicode

**Part – 2:**

**Virtual Machine Installation Steps**

We followed the following steps for the installation of VM and the image to run the wordcount and co-occurrence with Hadoop Map Reduce on HDFS-

1. Installed virtualbox from https://www.virtualbox.org/
2. Download the virtual image provided by TA on piazza
3. For loading the image, opened virtualbox and clicked **New.**
4. Used the following configurations: Type **-** Linux, **Version** -64bit ,**Name -** CSE587
5. Selected RAM as 8192 MB
6. When selecting hard disk, used the one provided by TA on piazza
7. Installed ubuntu by running the image.

Opened the terminal and executed the following commands to execute the word count and co-occurrence using mapper and reducer on Hadoop-

> start-dfs.sh                                        //starting namenodes and datanodes. Setting up the single node cluster.
> sudo rm -r /tmp/*                                        //cleaning everything in tmp directory
> hdfs namenode -format
> hdfs dfs -ls /                                        //check whether there is any directory in hdfs
> hdfs dfs -mkdir /test                                        //if there is no directory create one named test

> hdfs dfs -put /home/cse587/examplehadoop/exam.txt  /test   //putting the exam.txt from local directory into the test directory inside hdfs

> cd hadoop-3.1.2/bin

> hadoop jar /home/cse587/hadoop-3.1.2/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1.2.jar wordcount /test/exam.txt /test/output

> hdfs dfs -ls /test/output                                                          //You should see the output file as part-r0000

**Screenshots of Execution on VM**

**Code Snippets for mapper.py & reducer.py (Word Count)**

```python
#!/usr/bin/env python
"""mapper.py"""

import sys

# input comes from STDIN (standard input)
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # split the line into words
    words = line.split()
    # increase counters
    for word in words:
        # write the results to STDOUT (standard output);
        # what we output here will be the input for the
        # Reduce step, i.e. the input for reducer.py
        #
        # tab-delimited; the trivial word count is 1
        print '%s\t%s' % (word, 1)
```

```python
#!/usr/bin/env python
import sys

# Create a dictionary to map words to counts
wordcount = {}

# Get input from stdin
for line in sys.stdin:
    #Remove spaces from beginning and end of the line
    line = line.strip()

    # parse the input from mapper.py
    word, count = line.split('\t', 1)
    # convert count (currently a string) to int
    try:
        count = int(count)
    except ValueError:
        continue

    try:
        wordcount[word] = wordcount[word]+count
    except:
        wordcount[word] = count

# Write the tuples to stdout
# Currently tuples are unsorted
for word in wordcount.keys():
    print '%s\t%s'% ( word, wordcount[word] )
```

**Top-10 Words Word count**

| Twitter | | NYTimes | | Common Crawl | |
|---|---|---|---|---|---|
| cricket | 16145 | woods | 1537 | game | 7721 |
| golf | 12281 | players | 1501 | highlight | 6366 |
| basketball | 4426 | team | 1447 | sport | 5645 |
| tiger | 4051 | first | 1367 | virginia | 5225 |
| ipl | 3654 | game | 1122 | date | 4721 |
| game | 3271 | last | 989 | time | 4617 |
| play | 2974 | williams | 957 | point | 4589 |
| win | 2966 | year | 902 | win | 4308 |
| baseball | 2876 | hockey | 880 | score | 4204 |
| one | 2826 | play | 854 | team | 4176 |

**Code Snippets for mapper & reducer (Word Co-occurrence)**





Converted the following code into jar files for each data source and ran on Virtual Machine. The following screenshots to run the code is as follows.

**Top 10 words Co-Occurrence**

| New York Times | | | Twitter | | | Common Crawl | |
|---|---|---|---|---|---|---|---|
| **Data** | **Count** | | **Data** | **Count** | | **Data** | **Count** |
| team, year | 25 | | basketball ,game | 238 | | year, game | 6 |
| hockey, play | 26 | | game, baseball | 261 | | hockey , play | 8 |
| team, players | 26 | | play ,basketball | 265 | | play , game | 9 |
| team, first | 31 | | game ,cricket | 340 | | team, last | 11 |
| play,game | 34 | | play ,cricket | 409 | | play , first | 14 |
| woods, first | 43 | | play,golf | 532 | | last , game | 39 |
| players, hockey | 48 | | golf ,game | 607 | | first , year | 55 |
| first, game | 55 | | tiger,golf | 875 | | last ,year | 55 |
| team, hockey | 68 | | ipl,cricket | 980 | | first ,team | 58 |
| year, last | 208 | | ipl,ipl | 1092 | | game ,first | 94 |

**Part – 3:**

**Visualization**
**Tableau Installation:**
1.  Downloaded the setup for student version from https://www.tableau.com/academic/students
2.  Got the product key and registered for it.
3.  Ran the server setup
4.  Configured the essential Tableau server settings
5.  Set the authentication type
6.  Set run as service account
7.  Set the port and continue configuration
8.  Created a tableau server administrator user



This is our menu page where we can choose different visualization using the dropdown button.

Twitter Word Count

Count
2,826      16,145

MENU
▼

golf
12,281

ipl
3,654

one
2,826

play
2,974

baseball
2,876

basketball
4,426

game
3,271

tiger
4,051

cricket
16,145

win
2,966

**Analysis-**
This is the bubble visualization for top 10 words for Twitter. The color here shows the count strength shown in the count section on the right side. We have Menu button in each dashboard to navigate to all other dashboards.
Here the word with the highest count is cricket and the least is one.

New York Times Word Count

Count
854     1,537

MENU

▼

| | |
|---|---|
| year 902 | last 989 |
| play 854 | woods 1,537 |
| first 1,367 | game 1,122 |
| players 1,501 | williams 957 |
| hockey 880 | team 1,447 |

**Analysis-**

This is the bubble visualization for top 10 words for New York Times. The color here shows the count strength shown in the count section on the right side. We have Menu button in each dashboard to navigate to all other dashboards.

Here the word with the highest count is woods and the least is play.

## Common Crawl Word Count

Count
4,176    7,721

MENU

▼



**Analysis-**

This is the bubble visualization for top 10 words for Common Crawl. The color here shows the count strength shown in the count section on the right side. We have Menu button in each dashboard to navigate to all other dashboards.

Here the word with the highest count is game and the least is team.

Common Crawl Word Co-Occurrence

Count
6    94

hockey , play team, last

MENU

last ,year first ,team ▼

year, game

game ,first    play , game

last , game first , year

play , first

**Analysis-**

This is the Textual visualization for top 10 co-occurrence words for Common Crawl. The color here shows the count strength shown in the count section on the right side. We have Menu button in each dashboard to navigate to all other dashboards.

Here the co-occurrence words with the highest count are game and first and the co-occurrence words with the least count are year and game.

New York Times Word Co-Occurrence

Count
25 ▭ 208

team, year
play, game
team, players    players, hockey

first, game

MENU
▼

year, last

woods, first          team, hockey
team, first          hockey, play

**Analysis-**
This is the Textual visualization for top 10 co-occurrence words for New York Times. The color here shows the count strength shown in the count section on the right side. We have Menu button in each dashboard to navigate to all other dashboards.

Here the co-occurrence words with the highest count are year and last and the co-occurrence words with the least count are team and year.

Twitter Word Co-Occurrence

Count
238          1,092

game ,cricket

play ,cricket

ipl,ipl

golf ,game

MENU

ipl,cricket

play ,basketball
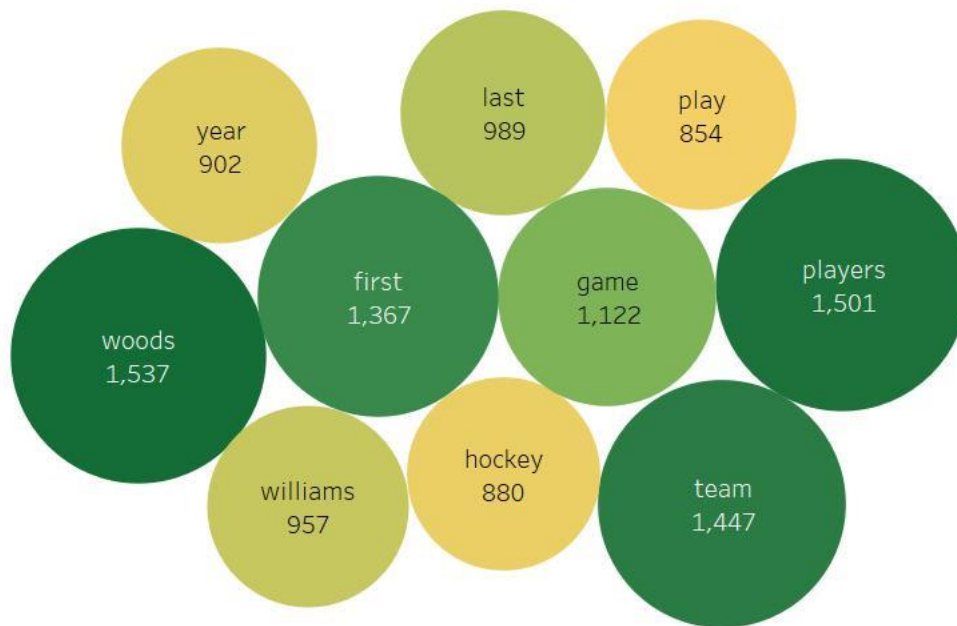
play,golf

tiger,golf

basketball ,game

game, baseball

**Analysis-**

This is the Textual visualization for top 10 co-occurrence words for Twitter. The color here shows the count strength shown in the count section on the right side. We have Menu button in each dashboard to navigate to all other dashboards.

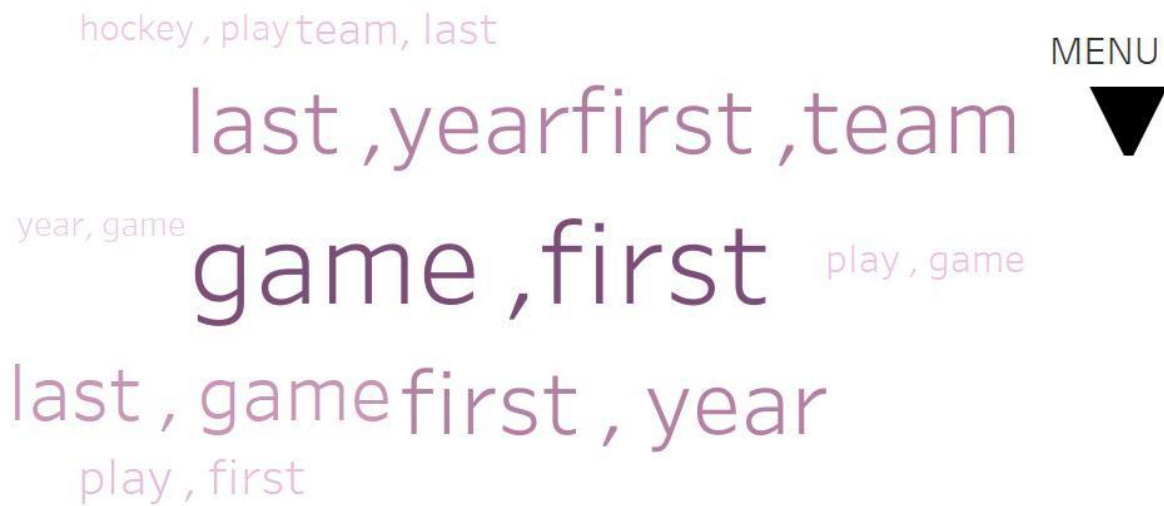Here the co-occurrence words with the highest count are ipl and cricket and the co-occurrence words with the least count are basketball and game.

**Link to published tableau: https://us-east-1.online.tableau.com/#/site/divyasrivastava/workbooks**
**We have published the whole online and the above is the link… Please login to view this..**

**Directory Structure:**

- **Part1**
  - **Code**
    - **Common Crawl**
    - **NYT**
    - **Twitter**
  - **Data**
    - **Common Crawl**
    - **NYT**
    - **Twitter**
- **Part2**
  - **Code**
    - **Co-occurrence**
    - **Word Count**
  - **Data**
    - **Common Crawl**
    - **NYT**
    - **Twitter**
  - **Screenshots**
- **Part3**
    - **Images**
    - **Workbook**
    - **Data folder**
- **Report**
- **readme**

**Demo Video link:** The video is made by an online software and hence has a watermark. Please ignore watermark.
Meghana Vasudeva's UB BOX:**https://buffalo.box.com/s/x6jc9nrkv4rs5gtp7dkaw68tz9gkkjcq**
Divya Srivastava's UB BOX: **https://buffalo.box.com/s/t72kq1b6zogahmzfd175kqr57028985r**

**Conclusion:**
- Automated data collection from multiple sources using the APIs offered by the businesses.
- Explained the importance of evaluating the reliability of data. Applied classical big data analytical methods: MapReduce for word count word occurrence.
- Worked on Hadoop and HDFS and processed the data using big data algorithms.
- Learnt a high level language-based data analysis by exploring Python as data processing language.
- Applied modern visualization methods and disseminate results using the web/mobile interface.

**References:**
1. Installation steps for VM: https://www.virtualbox.org/
2. Install steps for Tableau: https://www.tableau.com/academic/students
3. Steps for generating the twitter API:  https://developer.twitter.com/
4. Steps for generating the New York Times API: https://developer.nytimes.com/docs/articlesearch-product/1/overview