# CSCI 5408 Data Management, Warehousing, and Analytics

## Assignment – 2

Banner ID: - B00843607                                    Darpan Patel

---

### A. Cluster Setup: -

I have created Apache spark cluster in windows [1] (local machine).

Steps:
1. Install JDK 8, set JAVA_HOME = path to jdk file (ex. 'C:\Program files\Java\jdk') and PATH = "%JAVA_HOME%\bin"
2. Download "spark-2.4.4-bin-hadoop2.7.tgz" unpack and put the folder at "C:" drive (Or any location of your choice)
3. Set SPARK_HOME = "path where spark folder is located" (ex. 'C:\spark') and PATH = "%SPARK_HOME%\bin"
4. Check "pyspark" command in cmd.
5. Go inside "bin" folder in spark.
6. Run 'spark-class org.apache.spark.deploy.master.Master' (start Master and note IP and port address)
7. Run 'spark-class org.apache.spark.deploy.worker.Worker spark://IP:PORT'

### B. Twitter and News API Data Extraction and Cleaning: -

`       For extracting tweets:
- Created developer account in Twitter
- Create application to get twitter credentials to access tweeter APIs.
- Using "api.search" [2] to get tweets [tw.Cursor(api.search, q=query, tweet_mode="extended").items(numOfTweets)]

For extracting news:
- Create account to use NewsAPI, get API key.
- I have used api.get_everything()[3] to get articles.
  Example
       api.get_everything(q='University',
                          from_param='2019-10-05',
                          to='2019-11-02',
                          language='en',
                          sort_by='relevancy',
                          page_size = 100)

I have extracted total **3572** tweets and **500** news articles.

To clean data I have removed special characters, links, emoticons and converted everything to lower case using regex. After cleaning check for empty values in tweets and news, if present remove that row.

## C. Import data to MongoDb: -

- mongoimport --db data --collection tweets --type csv --file C:\Users\darpa\Desktop\Tweepy_tutorial\final_codes\news.csv --headerline
- mongoimport --db data --collection news --type csv --file C:\Users\darpa\Desktop\Tweepy_tutorial\final_codes\tweets.csv --headerline
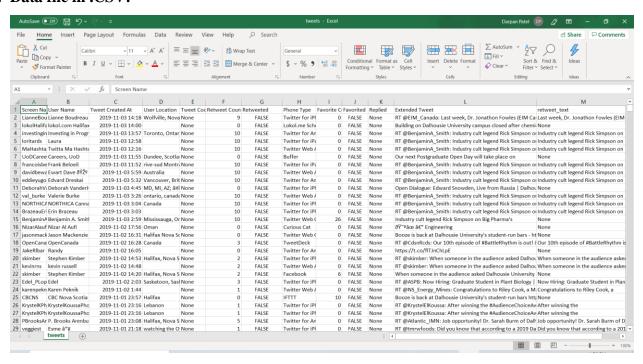
## D. Data file in .CSV: -



*Figure 1: Tweets in csv format*

## E. Data Processing: -

To process data, first I have merged 'tweets.csv' and 'news.csv' with axis = 0 in dataframe. Now, using flatMap and reduceByKey in spark I have calcuted frequency for 1-Grams (Canada, Dalhousie, university, graduate, faculty, education, expensive) and Bi-Grams (good school, bad schools, good schools, bad school, computer science). Whole program was run on cluster with 1 slave node (14.7 Gb RAM assigned)

```
freq_lst = freq_bigrams.collect()
```

```
freq_lst
```

```
[(5, ('computer', 'science'))]
```

```
sc.stop()
```

```
f = open('output.txt',"a+")
```

```
for item in word_count:
    f.write('{}:{}'.format(item[0],item[1]))
    f.write("\n")
```

```
for item in freq_lst:
    f.write('{}:{}'.format(item[1][0] + " " +item[1][1],item[0]))
    f.write("\n")
```

```
f.close()
```

output - Notepad

File  Edit  Format  View  Help

canada:1449
dalhousie:348
university:1267
graduate:16
faculty:13
education:818
expensive:5
computer science:5

*Figure 2 - Frequency Count Output*

# References

[1] "Installing Apache Spark(PySpark)" *Meduim*. [Online]. Available: https://medium.com/@loldja/installing-apache-spark-pyspark-the-missing-quick-start-guide-for-windows-ad81702ba62d.[Accessed: 27-Oct-2019].

[2] "tweepy.api – Tweeter API wrapper" Available: http://docs.tweepy.org/en/latest/api.html [Accessed: 01-Oct-2019].

[3] "Python client library," *News API*. [Online]. Available: https://newsapi.org/docs/client-libraries/python. [Accessed: 01-Oct-2019].

[4] *PySpark – Word Count Example*. [Online]. Available: https://pythonexamples.org/pyspark-word-count-example/. [Accessed: 03-Nov-2019].