# Exercise Sheet 06
# Distributed Data Analytics
# Syed Khalid Ahmed

## Exercise 01: Data cleaning and text tokenization:

### Flow of the program:

The program takes a single file as input, performs operations to perform data cleaning and generates an output file. This output file will be used in Exercise # 02 as an input. I have saved the output file in directories that follow a naming convention (Book-1, Book-2 , …) because the program of Exercise # 02 first reads the URL of the file in order to identify from which file the words are coming from. By using the directory names, I can identify the file.

### Mapper:

In the mapper function, I am reading the file line by line and outputting it to the reducer for further processing. I am stripping the line and checking if the line is a blank space or not. If it is not then output it to the reducer.

### Reducer:

In the Reducer, I have made a FileCleanser() function that removes all the punctuation marks and numbers. It also removes all the common stop words that are found in the English literature. I have created a set which contains all the possible stop words. The program uses regular expressions to remove all the punctuations and numbers. It then checks whether a word appears in the stopwords's set or not. If it is not a stopword then it is merged with the string.

### Hadoop Commands Used:

Following commands were used.

### To put the files in HDFS:

hadoop fs -put /home/khalid/Documents/*.txt /Exercise6-Data/

### For Book#1:

hadoop jar /home/khalid/Downloads/hadoop-2.7.2/share/hadoop/tools/lib/hadoop-*streaming*.jar -file /home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -mapper /home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -file /home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -reducer /home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -input /Exercise6-Data/1.txt -output /Exercise6-Solution/CleansedFiles/Book1/
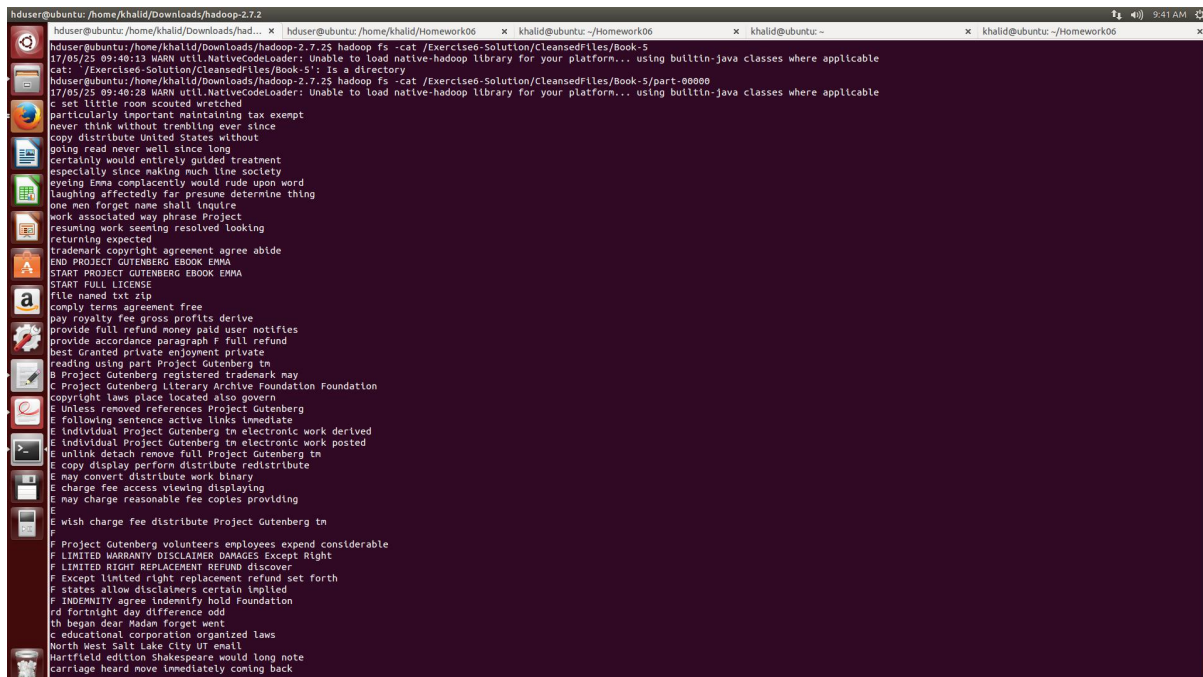
**For Book#2:**

hadoop jar /home/khalid/Downloads/hadoop-2.7.2/share/hadoop/tools/lib/hadoop-
*streaming*.jar -file /home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -mapper
/home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -file
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -reducer
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -input /Exercise6-Data/2.txt -
output /Exercise6-Solution/CleansedFiles/Book2/

**For Book#3:**
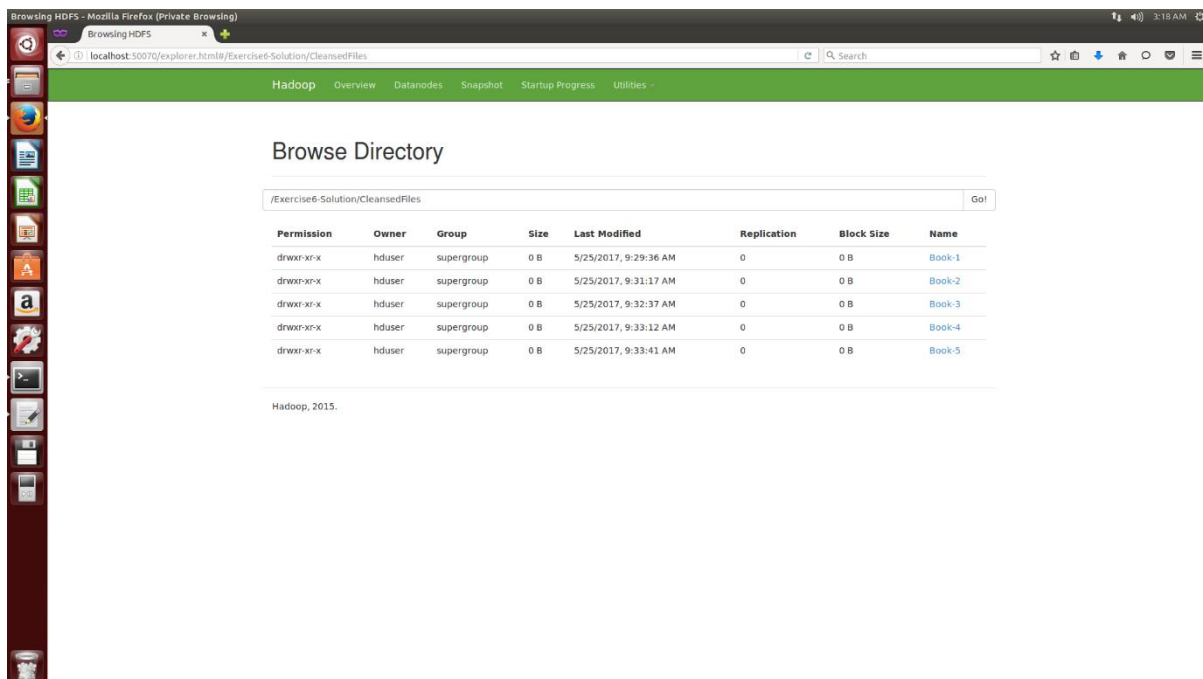
hadoop jar /home/khalid/Downloads/hadoop-2.7.2/share/hadoop/tools/lib/hadoop-
*streaming*.jar -file /home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -mapper
/home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -file
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -reducer
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -input /Exercise6-Data/3.txt -
output /Exercise6-Solution/CleansedFiles/Book3/

**For Book#4:**

hadoop jar /home/khalid/Downloads/hadoop-2.7.2/share/hadoop/tools/lib/hadoop-
*streaming*.jar -file /home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -mapper
/home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -file
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -reducer
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -input /Exercise6-Data/4.txt -
output /Exercise6-Solution/CleansedFiles/Book4/

**For Book#5:**

hadoop jar /home/khalid/Downloads/hadoop-2.7.2/share/hadoop/tools/lib/hadoop-
*streaming*.jar -file /home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -mapper
/home/khalid/Homework06/Exercise01/FileCleansing-mapper.py -file
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -reducer
/home/khalid/Homework06/Exercise01/FileCleansing-reducer.py -input /Exercise6-Data/5.txt -
output /Exercise6-Solution/CleansedFiles/Book5/

## Images:

This is the output of the Hadoop cat command.



Directory for each book

Each Book directory has part-00000 as the output file.



## Exercise 02: TFIDF scores of words:

### Flow of the program:

The program takes the all of the input files from the previous program, performs operations to compute TF , IDF and TF*IDF scores and generates an output file.

### Mapper:

In the mapper function, I am reading the files line by line and outputting it to the reducer for further processing. Since I have to find the id of the file from which the word is coming, so I have used:

input_file = os.environ['map_input_file']

to get the URL of the file that is currently being used by the mapper. It the returns a URL string from which I extract the directory name, perform a split and then use the number to identify the file. For example, my directory name is Book-2. I can extract 2 from this directory name which will tell me that this is file # 2.

The program then reads line from the file and generates and generates a key/value pair in which key is the word and value is the file id.

## Reducer:

      In the Reducer, the first function Compute() function counts the number of words in each file and the word counts for a specific file. The second function TFIDF() uses the information obtained in the first function to calculate the TF, IDF and TF*IDF scores for each word of a particular file.

      The output consists of four columns for word, TF, IDF and TF*IDF scores respectively.

## Hadoop Commands Used:

      Following command was used

hadoop jar /home/khalid/Downloads/hadoop-2.7.2/share/hadoop/tools/lib/hadoop-*streaming*.jar -file /home/khalid/Homework06/Exercise02/TF-IDF-mapper.py -mapper /home/khalid/Homework06/Exercise02/TF-IDF-mapper.py -file /home/khalid/Homework06/Exercise02/TF-IDF-reducer.py -reducer /home/khalid/Homework06/Exercise02/TF-IDF-reducer.py -input /Exercise6-Solution/CleansedFiles/Book-*/part-00000 -output /Exercise6-Solutions/Output/

In the above command, Book-* means that it will get all the directories that match this wildcard and output their part-00000 files.

## Images:

Output of the program

hduser@ubuntu: /home/khalid/Downloads/hadoop-2.7.2

hduser@ubuntu: /home/khalid/Downloads/had... ✕ | hduser@ubuntu: /home/khalid/Homework06 ✕ | khalid@ubuntu: ~/Homework06 ✕ | khalid@ubuntu: ~ ✕ | khalid@ubuntu: ~/Homework06 ✕

```
branch           2.73356614205e-05          0.0969100130081          2.64909930384e-06
disproportion             6.83391535512e-06          0.698970004336          4.7767018454e-06
disobliging               6.83391535512e-06          0.698970004336          4.7767018454e-06
er               7.51730689064e-05          0.221848749616           1.66770513417e-05
jiggedy          1.36678307102e-05          0.698970004336           9.55340369081e-06
incendiated               6.83391535512e-06          0.698970004336          4.7767018454e-06
malefactor                6.83391535512e-06          0.698970004336          4.7767018454e-06
album            1.36678307102e-05          0.698970004336           9.55340369081e-06
junk             6.83391535512e-06          0.698970004336           4.7767018454e-06
servest          6.83391535512e-06          0.698970004336           4.7767018454e-06
goodlooking               4.10034921307e-05          0.698970004336          2.86602110724e-05
quietus          6.83391535512e-06          0.698970004336           4.7767018454e-06
earthly          3.41695767756e-05          0.397940008672           1.35974416784e-05
geysers          6.83391535512e-06          0.698970004336           4.7767018454e-06
Musical          2.05017460654e-05          0.397940008672           8.15846500705e-06
shadowing                 1.36678307102e-05          0.698970004336          9.55340369081e-06
serpentplants             6.83391535512e-06          0.698970004336          4.7767018454e-06
Larch            6.83391535512e-06          0.698970004336           4.7767018454e-06
rotting          6.83391535512e-06          0.221848749616           1.51609557652e-06
weatherwise               6.83391535512e-06          0.698970004336          4.7767018454e-06
chinchopper               2.05017460654e-05          0.698970004336          1.43301055362e-05
space            0.000184515714588          0.0              0.0
jewel            2.73356614205e-05          0.397940008672           1.08779533427e-05
Voisin           6.83391535512e-06          0.698970004336           4.7767018454e-06
intentionally             6.83391535512e-06          0.397940008672          2.71948833568e-06
Lyster           1.36678307102e-05          0.698970004336           9.55340369081e-06
vane             6.83391535512e-06          0.698970004336           4.7767018454e-06
Boosed           6.83391535512e-06          0.698970004336           4.7767018454e-06
silverbuckled             6.83391535512e-06          0.698970004336          4.7767018454e-06
urchin           6.83391535512e-06          0.397940008672           2.71948833568e-06
patriarchs                6.83391535512e-06          0.698970004336          4.7767018454e-06


                 For Book # 5


Word             TF-score                   IDF-score                  TF*IDF-score


secondly                  1.33209004929e-05          0.0969100130081          1.29092864004e-06
pardon           0.000106567203943          0.0              0.0
similarity                1.33209004929e-05          0.221848749616           2.95522511811e-06
sashed           2.66418009857e-05          0.397940008672          1.06018385153e-05
yellow           3.99627014786e-05          0.0              0.0
four             0.000386306114293          0.0              0.0
protest          0.000106567203943          0.221848749616          2.36418009449e-05
woods            1.33209004929e-05          0.0              0.0
abide            2.66418009857e-05          0.0              0.0
hanging          3.99627014786e-05          0.0              0.0
aggression                1.33209004929e-05          0.698970004336          9.31090987526e-06
briskness                 1.33209004929e-05          0.397940008672          5.30091925765e-06
spoiled          5.32836019715e-05          0.0969100130081          5.16371456017e-06
Foundation                0.000319701611829          0.0              0.0
unfeeling                 3.99627014786e-05          0.397940008672           1.5902757773e-05
```