

COMP 6730 Advanced Database Systems Project

Nikhil Nagendra Swami
nikhil_swami@student.uml.edu
01539304
December 3, 2016

1) Overview:

For this project, I have configured the Hadoop environment as well as the Eclipse with Hadoop-eclipse-plugin, then setting up the namenode and other server on local machines, start the eclipse and create a Map-Reduce project in Eclipse.

Hadoop MapReduce is a software framework for easily writing applications which process big amounts of data in-parallel on large clusters (thousands of nodes) of commodity hardware in a reliable, fault-tolerant manner.

The Map Task: This is the first task, which takes input data and converts it into a set of data, where individual elements are broken down into tuples (key/value pairs).

The Reduce Task: This task takes the output from a map task as input and combines those data tuples into a smaller set of tuples. The reduce task is always performed after the map task.

The MapReduce framework consists of a single master JobTracker and one slave TaskTracker per cluster-node.

2) Hadoop Project:

Main aim of the project is to use MapReduce programming to parse text files and extract required information performing different operations like Selection, Projection, Natural Join and Aggregation.

To implement all operations a procedure followed is:

1. Initialize Driver program to instantiate Mapper and Reducer jobs.
2. Input specific file to mapper job.
3. Parse the specified text file line by line, split the fields based on delimiter ',' and store in an array as different columns.
4. Mapper maps the required columns and send pair to reducer job.
5. Reducer receives pairs from all mapper jobs and perform calculations to output required results.
6. Results are written in to output file.

3) Selection Query:

Computing Selection by MapReduce:

Task: To find cities whose population is larger than 300,000.

Input File: city.txt

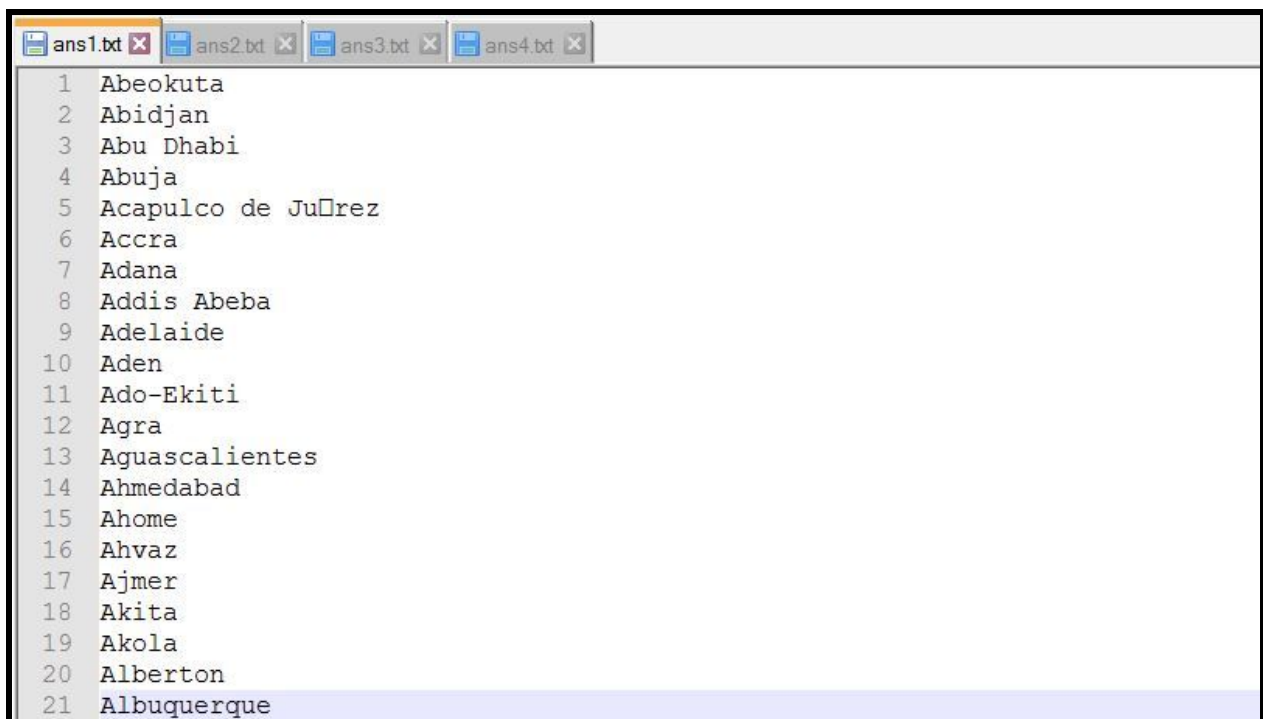
Output File: output.txt

In the map function each line of the “city.txt” are split on delimiter ‘,’ and individually stored in different columns. For this selection query we need column 2nd and 5th which are send as the key value pair(K, V) to the reducer.

In the reducer function all the key value pairs received from mapper are filtered. And based on the filter condition of population > 300,000 are written to the output.

Resulting rows are taken as output.

OUTPUT:

A screenshot of a text editor window with four tabs labeled 'ans1.txt', 'ans2.txt', 'ans3.txt', and 'ans4.txt'. The 'ans1.txt' tab is active and displays a list of 21 cities, each preceded by a line number from 1 to 21. The cities are: Abeokuta, Abidjan, Abu Dhabi, Abuja, Acapulco de Juarez, Accra, Adana, Addis Abeba, Adelaide, Aden, Ado-Ekiti, Agra, Aguascalientes, Ahmedabad, Ahome, Ahvaz, Ajmer, Akita, Akola, Alberton, and Albuquerque. The last line, '21 Albuquerque', is highlighted in light blue.

```
1 Abeokuta
2 Abidjan
3 Abu Dhabi
4 Abuja
5 Acapulco de Juarez
6 Accra
7 Adana
8 Addis Abeba
9 Adelaide
10 Aden
11 Ado-Ekiti
12 Agra
13 Aguascalientes
14 Ahmedabad
15 Ahome
16 Ahvaz
17 Ajmer
18 Akita
19 Akola
20 Alberton
21 Albuquerque
```

4) Projection

Computing Projection by MapReduce:

Task: To find all the name of the cities and corresponding district.

Input File: city.txt

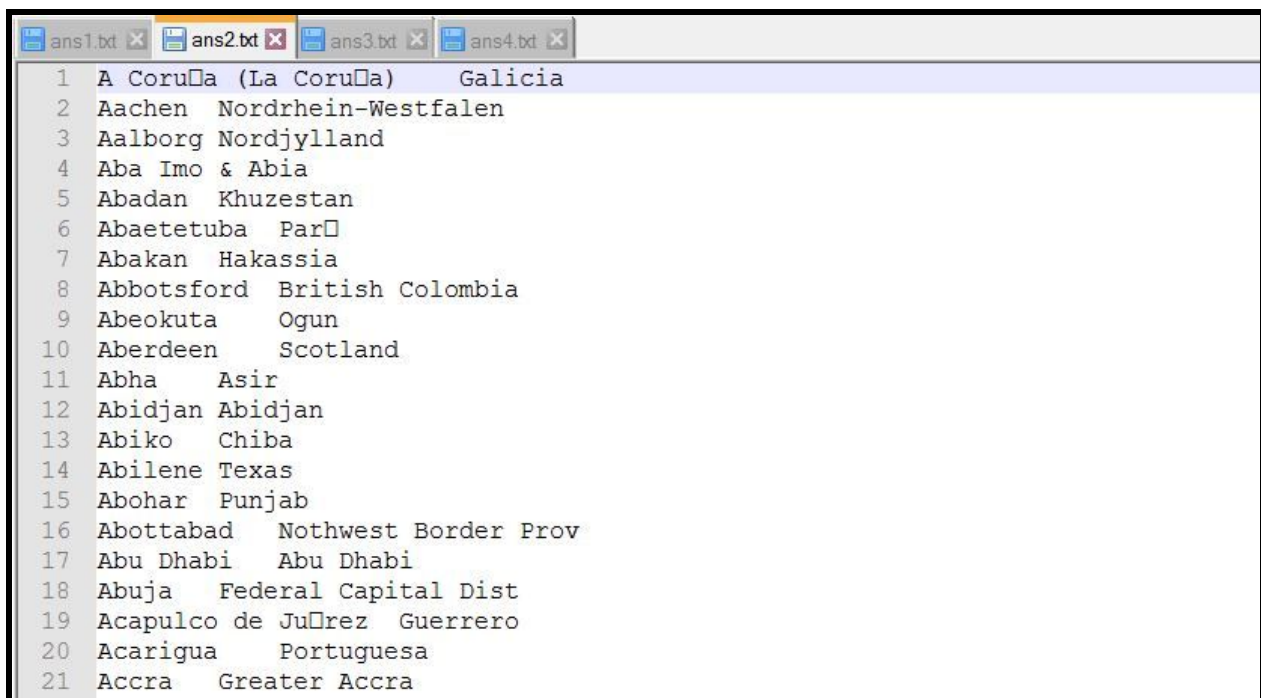
Output File: output.txt

In the mapping function, the mapper reads each line of the city.txt which has text separated by delimiter ',' and stored in separate columns.

For this selection query we need column 2nd and 4th which are send as the key value pair(K, V) to the reducer.

In the reducer, it receives data from the mapper as key value pair and the columns are further filtered down based on the condition of finding all names of cities and corresponding district. Resulting rows are taken as output.

OUTPUT:



1	A Coruña (La Coruña)	Galicia
2	Aachen	Nordrhein-Westfalen
3	Aalborg	Nordjylland
4	Aba Imo & Abia	
5	Abadan	Khuzestan
6	Abaetetuba	Pará
7	Abakan	Hakassia
8	Abbotsford	British Columbia
9	Abeokuta	Ogun
10	Aberdeen	Scotland
11	Abha	Asir
12	Abidjan	Abidjan
13	Abiko	Chiba
14	Abilene	Texas
15	Abohar	Punjab
16	Abottabad	Nothwest Border Prov
17	Abu Dhabi	Abu Dhabi
18	Abuja	Federal Capital Dist
19	Acapulco de Juárez	Guerrero
20	Acarigua	Portuguesa
21	Accra	Greater Accra

5) Natural Join:

Computing Natural Join by MapReduce: Task: To find all countries whose official language is English.

Input Files: country.txt and countrylanguage.txt

Output File: output.txt

In the map function each line of the “countrylanguage.txt” are split on delimiter ‘,’ and individually stored in different columns. For this selection query we need column 1st and 2nd columns which has the country code and respective language are send as the key value pair (K, V) to the reducer.

In the reducer, the key value pair of the country code and respective language send by the mapper are captured and stored. They are further used to find the name of the country where language is english using natural join between the 2 columns.

Resulting rows are taken as output.

OUTPUT:

A screenshot of a text editor window with four tabs labeled 'ans1.txt', 'ans2.txt', 'ans3.txt', and 'ans4.txt'. The 'ans3.txt' tab is active and displays a list of 21 countries, each preceded by a number from 1 to 21. The list is as follows:

```
1 Aruba
2 Anguilla
3 Netherlands Antilles
4 American Samoa
5 Antigua and Barbuda
6 Australia
7 Bahrain
8 Belize
9 Bermuda
10 Barbados
11 Brunei
12 Canada
13 Cocos (Keeling) Islands
14 Cook Islands
15 Christmas Island
16 Cayman Islands
17 Denmark
18 Falkland Islands
19 United Kingdom
20 Gibraltar
21 Guam
```

6) Aggregation:

Computing Aggregation by MapReduce:

Task: To find how many cities each district has.

Input File: city.txt

Output File: output.txt

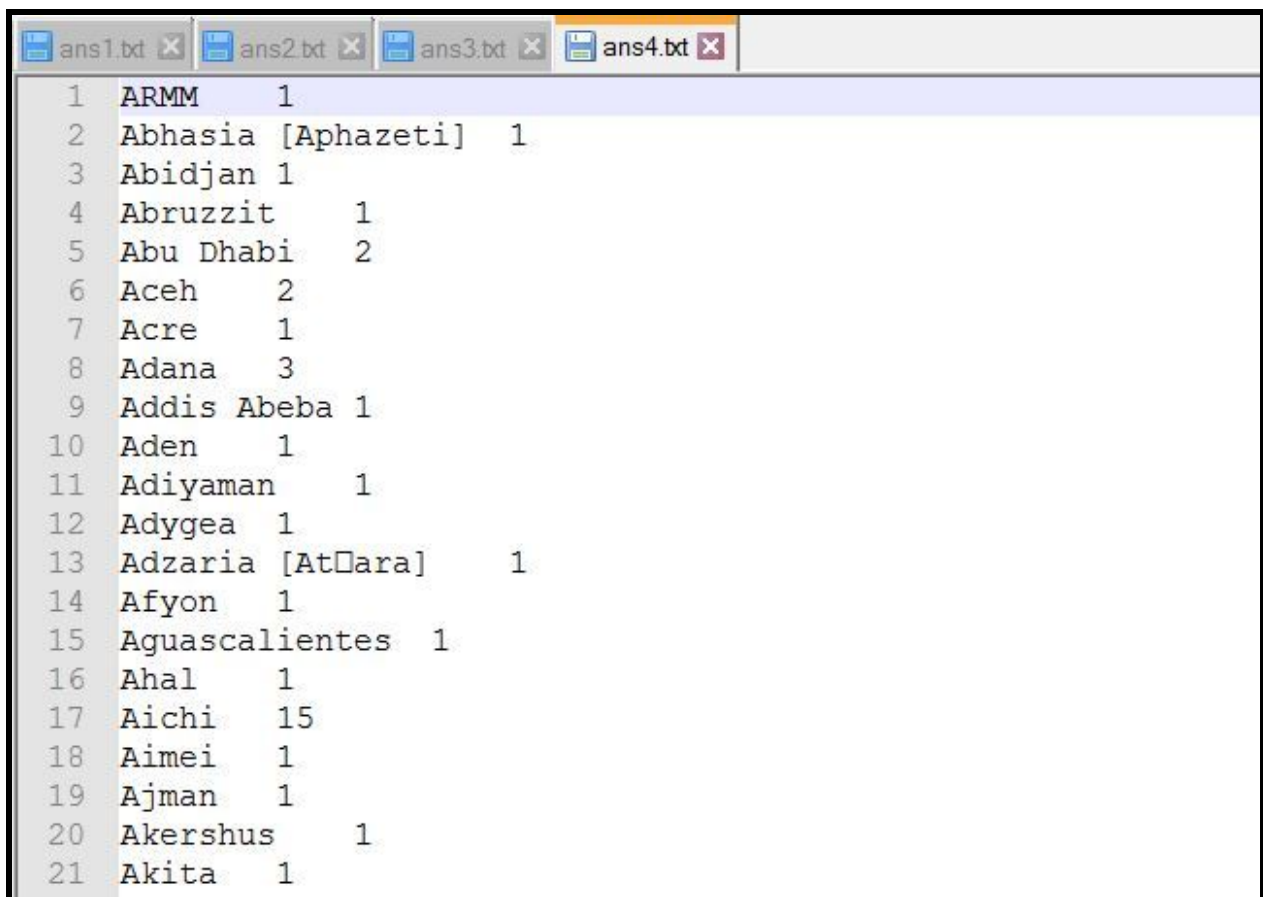
Map:

In the map function each line of the "city.txt" are split on delimiter ',' and individually stored in different columns. We get the city-district key value pair (K, V).

Reduce:

All the <key, value> pairs are received from all mapper jobs are used to count of number of cities is calculated based on received pairs.

OUTPUT:



1	ARMM	1
2	Abhasia [Aphazeti]	1
3	Abidjan	1
4	Abruzzit	1
5	Abu Dhabi	2
6	Aceh	2
7	Acre	1
8	Adana	3
9	Addis Abeba	1
10	Aden	1
11	Adiyaman	1
12	Adygea	1
13	Adzaria [Atāara]	1
14	Afyon	1
15	Aguascalientes	1
16	Ahal	1
17	Aichi	15
18	Aimei	1
19	Ajman	1
20	Akershus	1
21	Akita	1

7) Problems Faced:

One of the problems faced in this project was the installation of Hadoop on a Linux system and using it to get the jar out file for processing of the map-reduce function.

Selection of the columns and getting the values by the delimiters gave some problem but once I was able to get for the first program, it was easier to get the values for processing.

8) Conclusion:

The main conclusion from this project was to get a sense of Hadoop working and to process huge amount of data and selecting the required data to find according to the query processed. The assignments were pretty clear and understandable. Processing on the key value pair helped in processing the queries and getting the desired output.

9) Reference:

1. "Welcome To Apache™ Hadoop®!". Hadoop.apache.org. N.p., 2016. Web. 3 Dec. 2016.
2. "Hadoop - Command Reference". www.tutorialspoint.com. N.p., 2016. Web. 3 Dec. 2016.
3. "Setting Up A Apache Hadoop 2.7 Single Node On Ubuntu 14.04". the power of data. N.p., 2016. Web. 3 Dec. 2016.
4. "Word Count Job Implementation In Hadoop". YouTube. N.p., 2016. Web. 3 Dec. 2016.