# Internal Assignment 2

**Student Id: 2048494**

**Student Name: Dhaniswar B.K.**

**Module Leader: Jnaneshwar Bohara**

**Tutor: Sujan Shrestha**

**Section: L6CG5**

**Submitted on : 13/04/2022**

# Table of Content

# 1. Dividing my university id by 3.

Answer: calculating the remainder of the university id using short trick

Step I => 2048494 = 2+0+4+8+4+9+4 => 31 => 3+1 = 4

Now, 4%3 = 1 So, I have to use Degree-Level_Quals.csv data set for this assessment.

# 2. Java and Hadoop

## 2.1. Copying the required files to the destination

## 2.2. Amending the Population class name to the DegreeQuals class



## 2.3. Amend the Mapper and Reducer class names to Maper_2048494 and Reducer_2048494 respectively

Reflecting these changes in the main method



## 3.   Run the code

### 3.1.   Starting the jps

## 3.2. Checking the directory in the Hadoop file system

```
dhaniswar@anonymous:~$ hdfs dfs -ls /user/dhaniswar/
Found 1 items
drwxr-xr-x   - dhaniswar supergroup          0 2022-04-13 00:35 /user/dhaniswar/spark_output_word
dhaniswar@anonymous:~$ hdfs dfs -mkdir /user/dhaniswar/input_csv
dhaniswar@anonymous:~$
```
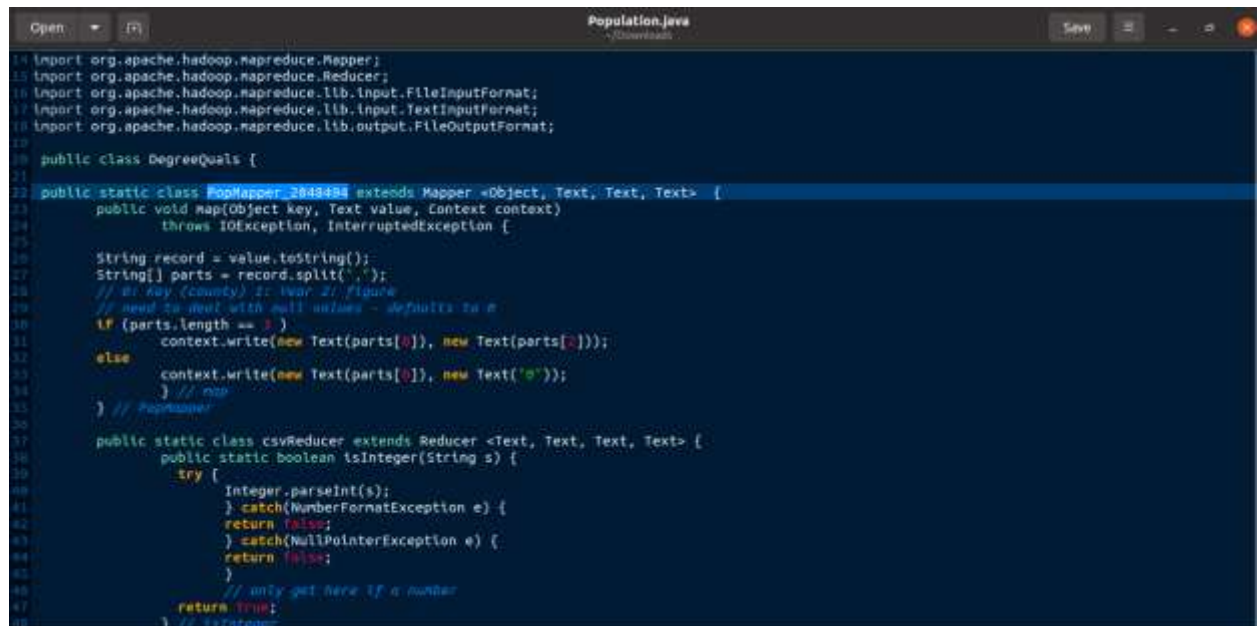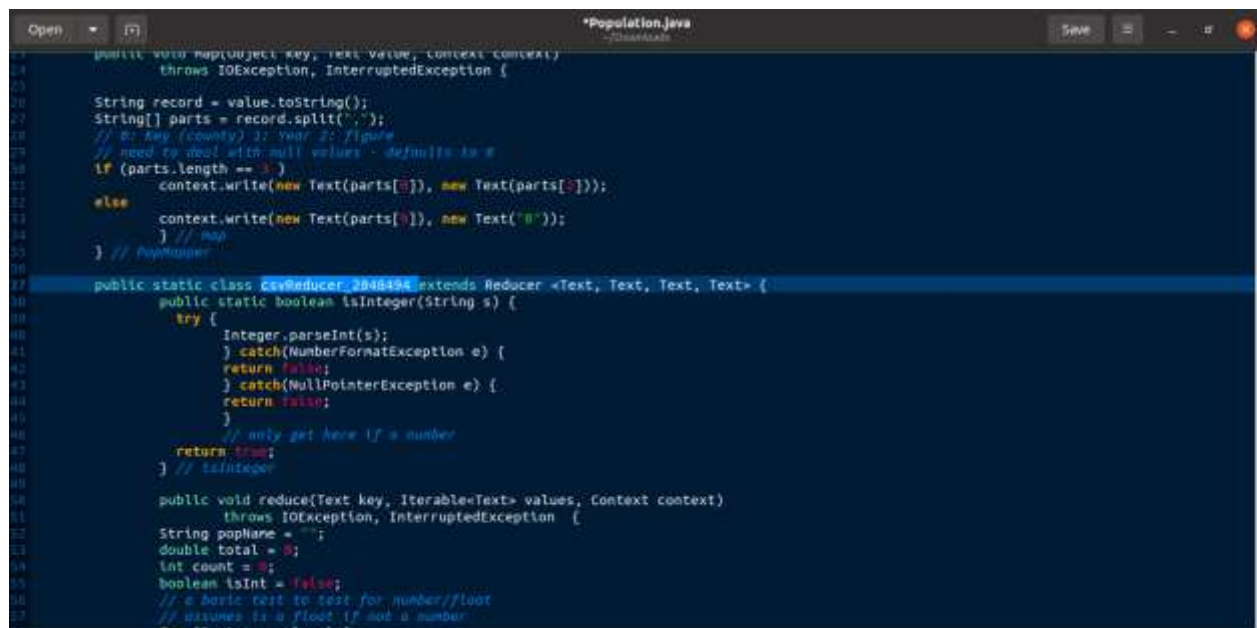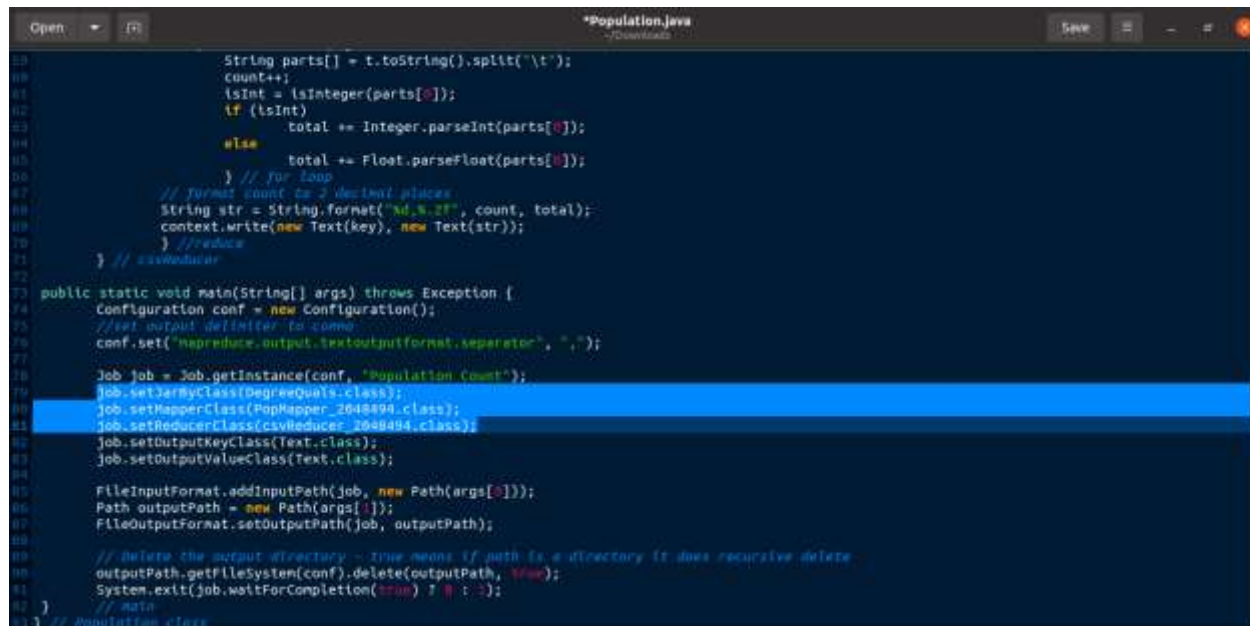
## 3.3. Compiling the java file and creating the jar file

```
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ ls
Degree-Level_Quals.csv  Population.java
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ mv Population.java DegreeQuals.java
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ ls
Degree-Level_Quals.csv  DegreeQuals.java
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ javac -classpath $(hadoop classpath) DegreeQuals.java
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ ls
 Degree-Level_Quals.csv  'DegreeQuals$csvReducer_2048494.class'  'DegreeQuals$PopMapper_2048494.class'   DegreeQuals.class   DegreeQuals.java
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ jar cf DegreeQuals.jar DegreeQ*.class
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ ls
 Degree-Level_Quals.csv                  'DegreeQuals$PopMapper_2048494.class'   DegreeQuals.jar
'DegreeQuals$csvReducer_2048494.class'   DegreeQuals.class                       DegreeQuals.java
dhaniswar@anonymous:~/Desktop/Internal_Assignment$
```

## 3.4. Executing the java programme

```
dhaniswar@anonymous:~/Desktop/Internal_Assignment$ hadoop jar DegreeQuals.jar DegreeQuals input_csv/Degree-Level_Quals.csv output_csv
2022-04-13 15:55:02,707 INFO client.RMProxy: Connecting to ResourceManager at /127.0.0.1:8032
2022-04-13 15:55:03,102 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute
 your application with ToolRunner to remedy this.
2022-04-13 15:55:03,124 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/dhaniswar/.staging/job_1649843
133371_0002
2022-04-13 15:55:03,420 INFO input.FileInputFormat: Total input files to process : 1
2022-04-13 15:55:03,520 INFO mapreduce.JobSubmitter: number of splits:1
2022-04-13 15:55:03,884 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1649843133371_0002
2022-04-13 15:55:03,886 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-04-13 15:55:04,024 INFO conf.Configuration: resource-types.xml not found
2022-04-13 15:55:04,024 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2022-04-13 15:55:04,490 INFO impl.YarnClientImpl: Submitted application application_1649843133371_0002
2022-04-13 15:55:04,536 INFO mapreduce.Job: The url to track the job: http://anonymous:8088/proxy/application_1649843133371_0002/
2022-04-13 15:55:04,537 INFO mapreduce.Job: Running job: job_1649843133371_0002
2022-04-13 15:55:13,695 INFO mapreduce.Job: Job job_1649843133371_0002 running in uber mode : false
2022-04-13 15:55:13,698 INFO mapreduce.Job:  map 0% reduce 0%
2022-04-13 15:55:18,773 INFO mapreduce.Job:  map 100% reduce 0%
2022-04-13 15:55:24,818 INFO mapreduce.Job:  map 100% reduce 100%
2022-04-13 15:55:24,829 INFO mapreduce.Job: Job job_1649843133371_0002 completed successfully
2022-04-13 15:55:24,955 INFO mapreduce.Job: Counters: 54
        File System Counters
                FILE: Number of bytes read=84903
                FILE: Number of bytes written=638575
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=157831
                HDFS: Number of bytes written=7182
                HDFS: Number of read operations=8
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Launched map tasks=1
                Launched reduce tasks=1
```

```
                Map input records=4550
                Map output records=4550
                Map output bytes=75797
                Map output materialized bytes=84903
                Input split bytes=134
                Combine input records=0
                Combine output records=0
                Reduce input groups=325
                Reduce shuffle bytes=84903
                Reduce input records=4550
                Reduce output records=325
                Spilled Records=9100
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=61
                CPU time spent (ms)=2640
                Physical memory (bytes) snapshot=439874816
                Virtual memory (bytes) snapshot=5476130816
                Total committed heap usage (bytes)=370147328
                Peak Map Physical memory (bytes)=261586944
                Peak Map Virtual memory (bytes)=2733181056
                Peak Reduce Physical memory (bytes)=177487872
                Peak Reduce Virtual memory (bytes)=2743029760
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=157697
        File Output Format Counters
                Bytes Written=7182
dhaniswar@anonymous:~/Desktop/Internal_Assignment$
```

**3.5.** **Navigating the input and output directories of Hadoop file system and found success message in output_csv directory and storing the successfully created file from Hadoop file system to local directory.**



**3.6.** **Content of the output directory**

## 4. Apache Spark

### 4.1. Starting the apace spark using command pyspark



### 4.2. Loading the same CSV file into Apache Spark

## 4.3. Selecting the first columns of the dataset with 25 rows using data frame

```
>>> df.select("_c0").show(25)
+--------+
|     _c0|
+--------+
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
| Babergh|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
|Basildon|
+--------+
only showing top 25 rows
```

## 4.4. Filtering specific records using data frame

```
>>> df.filter(df["_c0"]=="Basildon").show()
+--------+-------------------+----+
|     _c0|                _c1| _c2|
+--------+-------------------+----+
|Basildon|Jan 2004-Dec 2004|18.7|
|Basildon|Jan 2005-Dec 2005|16.7|
|Basildon|Jan 2006-Dec 2006|19.7|
|Basildon|Jan 2007-Dec 2007|17.6|
|Basildon|Jan 2008-Dec 2008|19.4|
|Basildon|Jan 2009-Dec 2009|18.6|
|Basildon|Jan 2010-Dec 2010|16.8|
|Basildon|Jan 2011-Dec 2011|17.1|
|Basildon|Jan 2012-Dec 2012|  25|
|Basildon|Jan 2013-Dec 2013|24.0|
|Basildon|Jan 2014-Dec 2014|  25|
|Basildon|Jan 2015-Dec 2015|28.9|
|Basildon|Jan 2016-Dec 2016|33.5|
|Basildon|Jan 2017-Dec 2017|27.8|
+--------+-------------------+----+

>>>
```

## 4.5. Selecting one columns using SQL query in spark

```
>>> sqlDF = spark.sql("SELECT _c2 FROM Degree")
>>> sqlDF.show()
+----+
| _c2|
+----+
|24.1|
|23.2|
|25.3|
|27.1|
|23.8|
|23.2|
|23.5|
|31.4|
|27.8|
|  18|
|24.7|
|23.4|
|32.9|
|40.7|
|18.7|
|16.7|
|19.7|
|17.6|
|19.4|
|18.6|
+----+
only showing top 20 rows

>>>
```

## 4.6. Selecting two columns on the basis of first column value equals to Babergh using SQL query in spark

```
>>> df.createOrReplaceTempView("Degree")
>>> sqlDF = spark.sql("SELECT _c1, _c2 FROM Degree WHERE _c0=='Babergh'")
>>> sqlDF.show()
+-------------------+----+
|                _c1| _c2|
+-------------------+----+
|Jan 2004-Dec 2004|24.1|
|Jan 2005-Dec 2005|23.2|
|Jan 2006-Dec 2006|25.3|
|Jan 2007-Dec 2007|27.1|
|Jan 2008-Dec 2008|23.8|
|Jan 2009-Dec 2009|23.2|
|Jan 2010-Dec 2010|23.5|
|Jan 2011-Dec 2011|31.4|
|Jan 2012-Dec 2012|27.8|
|Jan 2013-Dec 2013|  18|
|Jan 2014-Dec 2014|24.7|
|Jan 2015-Dec 2015|23.4|
|Jan 2016-Dec 2016|32.9|
|Jan 2017-Dec 2017|40.7|
+-------------------+----+
```

## 5. Hadoop is fast (advantage)

Data processing tools are frequently housed on the same servers as the data, resulting in the substantially quicker data processing. Hadoop can easily handle terabytes of data in minutes or petabytes in hours if you're working with enormous amounts of unstructured data. Hadoop's one-of-a-kind storage system is built on a distributed file system that essentially "maps" data to any location on a cluster. Data processing tools are frequently housed on the same servers as the data, resulting in the substantially quicker data processing. Hadoop can easily handle terabytes of data in minutes or petabytes in hours if you're working with enormous amounts of unstructured data.

## 6. Issue with Small Files (disadvantage)

Hadoop is not suitable for handling modest amounts of data. (HDFS) Because of its high-capacity architecture, the Hadoop distributed file system is unable to efficiently support the random reading of tiny files. In HDFS, small files are the most common issue. The size of a tiny file is much smaller than the HDFS block size (default 128MB). HDFS can't manage this many little files since it's designed to deal with a limited number of large files for storing massive data sets rather than a large number of small files. Because it holds the HDFS namespace, the Name Node will get overloaded if there are too many tiny files.