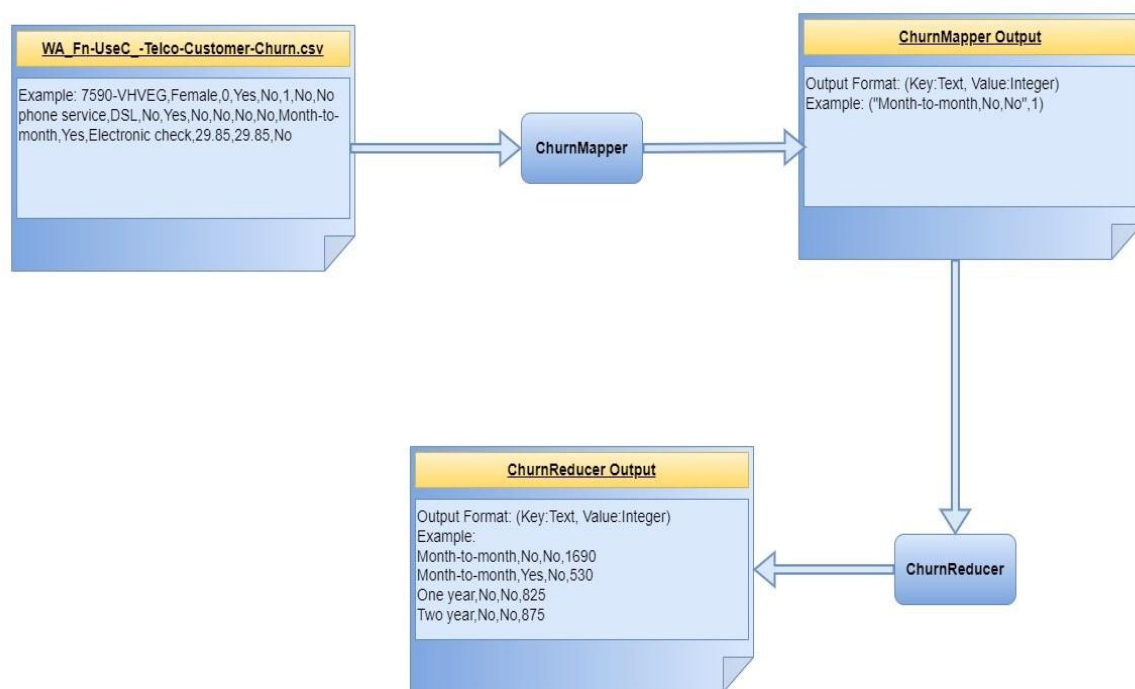**Data Flow:**



## 2.1 Implementation Steps:

1. As per the MapReduce programming model, create the below three classes and generate a .jar file for this requirement.

   a) ChurnMapper
   b) ChurnReducer
   c) ChurnAnalysis

● The complete implementation and the code are in the separate file '**Mapreduce_ChurnAnalysis**' and the file has been submitted along with this coursework.

2. Firstly, start the HDFS file system using commands 'start-dfs.sh' and 'start-yarn.sh' and confirm if all the services are running or not using a command called 'jps'



*Figure 1*

- As per figure 1, all the services are up and running.

3. Create the directory in the HDFS file system to provide our input to the HDFS and execute the MapReduce job. In this case, I have created the 'Telecom' directory and then created the 'Input' directory inside the 'Telecom' directory as shown below.



*Figure 2*

4. Copy and place the input file 'WA_Fn-UseC_-Telco- Customer-Churn.csv' in the 'Input' folder using the 'put' command as shown below.
**Note:** Here, I have renamed the original downloaded file from 'WA_Fn-UseC_-Telco- Customer-Churn' to 'Telco_Customer_Churn'.



*Figure 3*

5. Now, check whether the input file 'Telco_Customer_Churn.csv' is placed inside the /Telecom/Input directory or not. To do that, login to the interface, by default, it is available at http://localhost:9870. Next, go to Utilities and browse the file system **'/Telecom/Input/'.** Notice that the file 'Telco_Customer_Churn' is inserted in **/Telecom/Input/** directory.
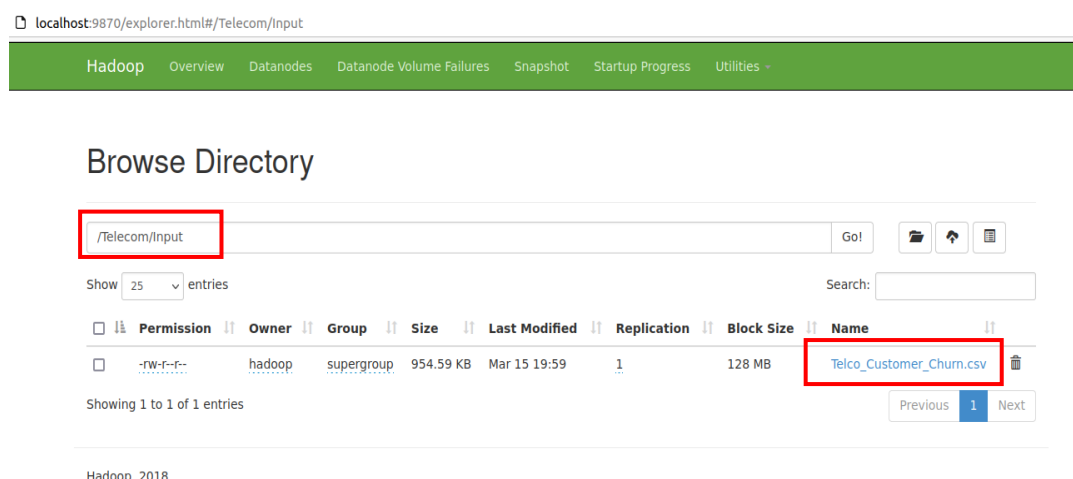


*Figure 4*

6. Now, In the Eclipse IDE, create the project and the required Java classes and export the .jar file. (Steps to create the java classes and exporting the .jar file mentioned in section 2.3 JAVA Classes Implementation and 2.4 Export the Jar file respectively)

7. Finally, run the exported jar file using the below command.

**hadoop jar <jar file name> <driver class name> <input path> <output path>**

Ex:<hadoop jar Downloads/CustomerChurn.jar org.myorg.ChurnAnalysis /Telecom/Input /Telecom/Output>



*Figure 5*

8. When the job runs successfully, no error messages will be displayed after the completion of MapReduce as shown below.



*Figure 6*

*Figure 7*

9. Check the output of the job which is created in the following location. <mark>/Telecom/Output/part-r-00000</mark>



*Figure 8*

## 2.2 JAVA Classes Implementation

1. Create the project as shown in figure 9. Also, four additional jars were used in the development of this project.

    a. hadoop-common.jar – available at
        <mark></usr/local/hadoop/share/hadoop/common></mark>

    b. hadoop-mapreduct-client-core.jar – available at
        <mark></usr/local/hadoop/share/hadoop/mapreduce></mark>

    c. hadoop-mapreduct-client-jobclient.jar – available at
        <mark></usr/local/hadoop/share/hadoop/mapreduce></mark>

    d. log4j.jar – available at

`</home/hadoop/HPCI_Jars >`

2. To obtain the log4j jar file needed for the project, you can download it from the official Apache log4j website at https://logging.apache.org/log4j/1.2/download.html. A separate file named "Mapreduce_ChurnAnalysis" has been provided with clear instructions for downloading the jar file. Log4j is a Java logging framework developed by the Apache Software Foundation and is used to track and manage program logs.



*Figure 9*

## ChurnMapper (Mapper Class):

The Mapper class will take a line from the CSV file one by one as input. For each line, first, it will convert the line into a string and store it in a line variable. If the line is the header of data, then, it will avoid that line else it will split the line by comma to separate every field of a row. We need only contracts, dependents, and churn fields for this research question. So, these fields are accessed by using their respective indexes. The result of the mapper class consists of a key and a value. Mapper class will return contract, dependents and churn as the key in the form of Text and 1 as the value in the form of IntWritable. So, for each customer, we will get its contract, dependents, and churn as keys and 1 as values by using the mapper class.

*Figure 10*

Inheritance from a mapper class

Input key value

Output key and Value datatypes

if statement is used to avoid headers in input data.

```java
1 package org.myorg;
2 //Importing the libraries that contain all methods needed for mapper class
3   import java.io.IOException;
8 //creating a subclass called ChurnReducer from parent class using keyword extends
9     public class ChurnMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
10
11        private final static IntWritable one = new IntWritable(1);
12        private Text KEY = new Text();
13
14 /*To handle the objects in Hadoop manner, hadoop uses Text instead of Java's String.Hadoop's test class is
15  comparable to a java string. Yet, the text implements interfaces like LongWritable and IntWritable,
16  both of which are required for MapReduce. As the reducer sorts the keys, it uses the similar interface
17  for comparison, and Writable can output the results to the local disc. Hadoop objects can be lightly serialised
18  using IntWritable and LongWritable.
19        */
20        @Override
21        public void map(LongWritable key, Text value, Context context)
22            throws IOException, InterruptedException {
23
24        String line = value.toString();
25        if (!line.startsWith("customerID")){
26            String[] fields = line.split(",");
27            String result = fields[15].concat(",").concat(fields[4]).concat(",").concat(fields[20]);
28            KEY.set(result);
29            context.write(KEY,one);
30        }
31
32      }
33    }
34
```

## ChurnReducer (Reducer Class):

The Reducer class will take the results of the mapper class as input. It will group the values of the mapper class result by key and will aggregate those values. In this reducer class, the sum is applied as an aggregation function. For each key, it will take the 'sum' of all of its values. Reducer class will return contract, dependents, and churn as key in form of Text and 'sum of all values of that key' as the value in the form of IntWritable. So, the output would be in the following form:

**Contract, dependent, churn, Sumofvalues**

**Example: Month-to-Month, No, Yes, 1396**

In this example, 'Month-to-month' is a contract, 'No' is dependents, and 'Yes' is churn. Thus, 1396 is the total number of customers who have a Month-to-month contract, have 'No' dependents, and shave churned the subscription (which means the customer has unsubscribed/left the contract).

*Figure 11*

## ChurnAnalysis (Driver Class):

This is the driver class of the program. It is implemented to run the whole project. Make sure both the input and output paths are passed in the command. It will create the job instance and is responsible for passing inputs to the mapper and reducer classes. It will delete the output folder from hdfs if that folder already exists. It returns a success message after the job gets completed successfully.

- **Job Name:** name of the job, job
- **Executable (Jar) Class:** the main executable class. Here it's, ChurnAnalysis.
- **Mapper Class:** class that overrides the "map" function. Here it's, ChurnMapper.
- **Reducer Class:** class that overrides the "reduce" function. Here it's, ChurnReducer.
- **Output Key:** type of output key. It's, Text.
- **Output Value:** type of output value. It's IntWritable.
- **File Input Path:** /Telecom/Input
- **File Output Path:** /Telecom/Output

```
 1 package org.myorg;
 2 //Import all the libraries that contain all methods needed for
 3⊖ import org.apache.hadoop.conf.Configuration;
11
12 public class ChurnAnalysis {
13
14⊖   public static void main(String[] args) throws Exception {
15 /*create the conf object from the configuration class which provides ... configuration parameters
16 necessary for hadoop job*/
17     Configuration conf = new Configuration();
18     conf.set("mapreduce.output.textoutputformat.separator", ",");
19 /*if the number of arguments are 2, then print the message that two arguments are needed, else exit*/
20     if (args.length != 2) {
21       System.out.printf("Usage: ChurnAnalysis <input dir> <output dir>\n");
22       System.exit(-1);
23     }
24 /*create a new job from the job class.The job represents the task that needs to be executed
25 The Job class is the most important class in the mapreduce API. The user is able to set up the job,
26 submit it, manage its execution, and check its status. Set parameters for the newly created job.*/
27     Job job;
28     job = Job.getInstance(conf,"Telecommunication Churn Analysis");
29 /* setting the driver class in the JAR file */
30     job.setJarByClass(ChurnAnalysis.class);
31 /* setting the input and output paths for the job */
32     FileInputFormat.addInputPath(job, new Path(args[0]));
33     FileOutputFormat.setOutputPath(job, new Path(args[1]));
34 /* setting the mapper and reducer for the job */
35     job.setMapperClass(ChurnMapper.class);
36     job.setReducerClass(ChurnReducer.class);
37     job.setCombinerClass(ChurnReducer.class);
38 /* Setting the key class (Text) and the class for the job output data*/
39     job.setOutputKeyClass(Text.class);
40     job.setOutputValueClass(IntWritable.class);
41 /*Delete Output file If already exist*/
42     FileSystem hdfs = FileSystem.get(conf);
43     Path outputDir = new Path(args[1]);
44     if (hdfs.exists(outputDir)){
45       hdfs.delete(outputDir, true);
46     }
47
48     boolean success = job.waitForCompletion(true);
49 /* check the status of the job (success or not) and exit when its done*/
50     System.exit(success ? 0 : 1);
51   }
52 }
53
```
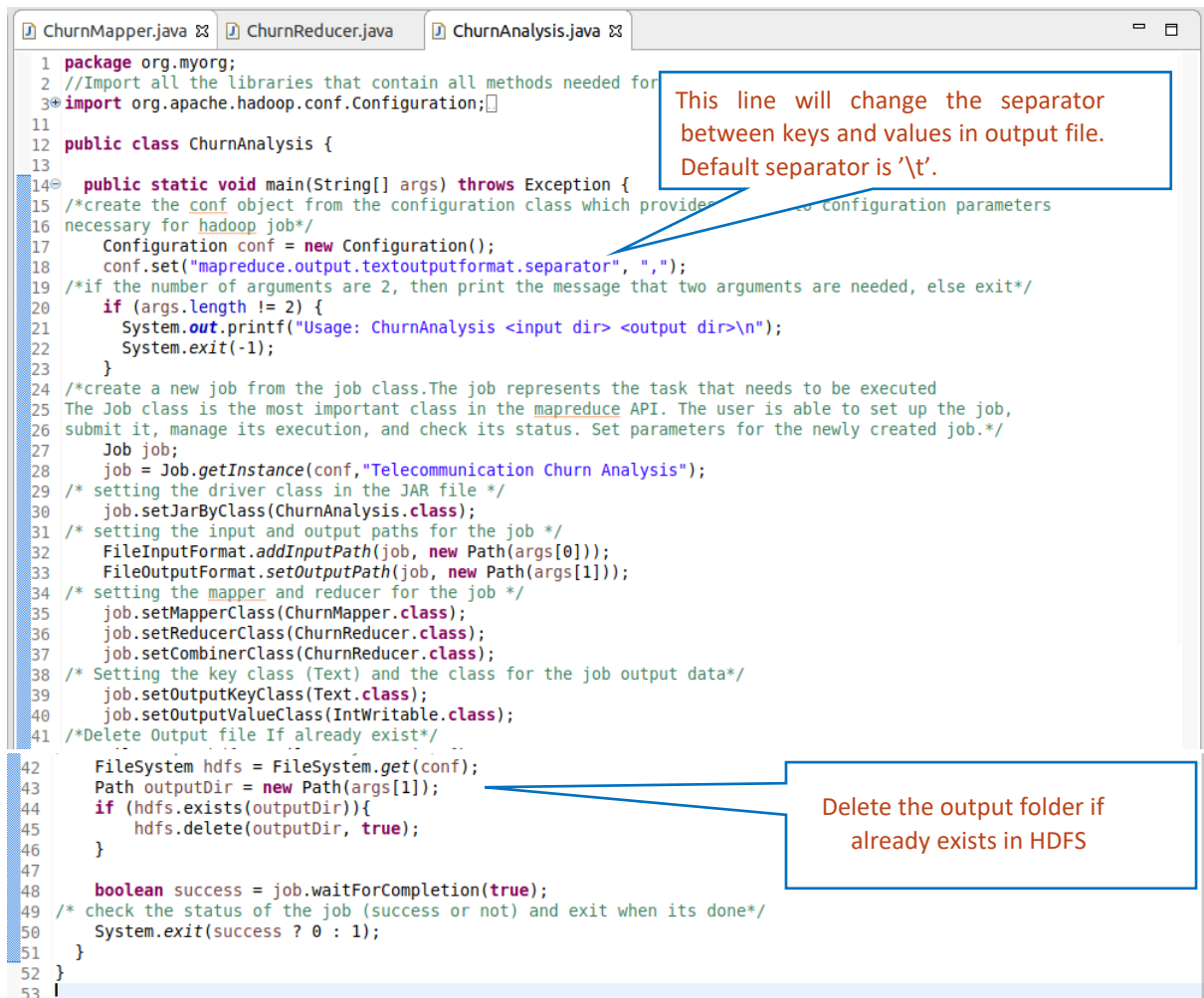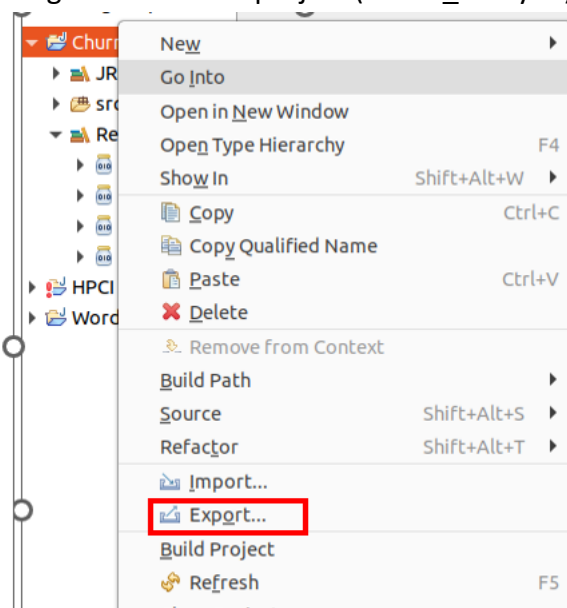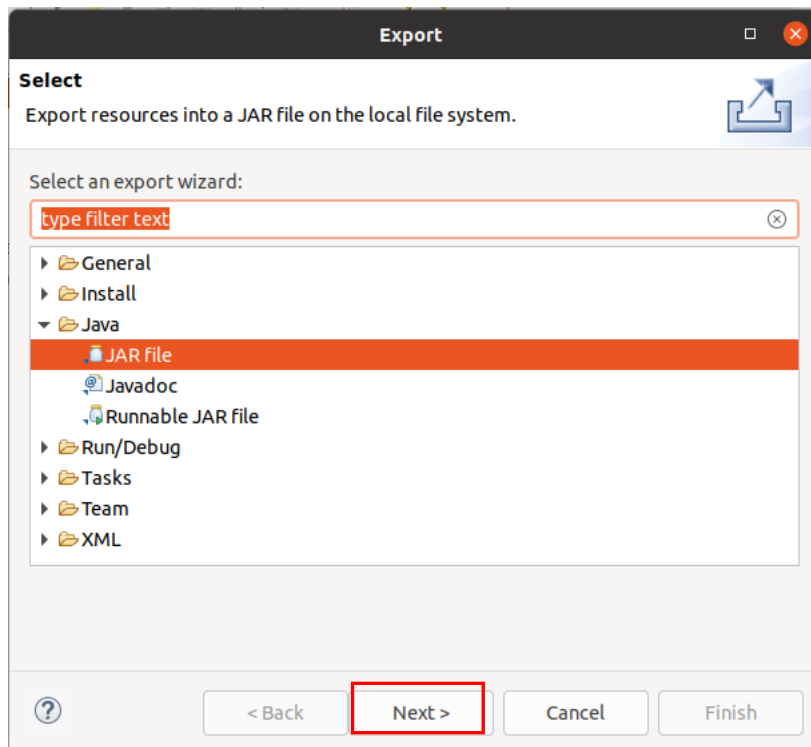
This line will change the separator between keys and values in output file. Default separator is '\t'.

Delete the output folder if already exists in HDFS

*Figure 12*
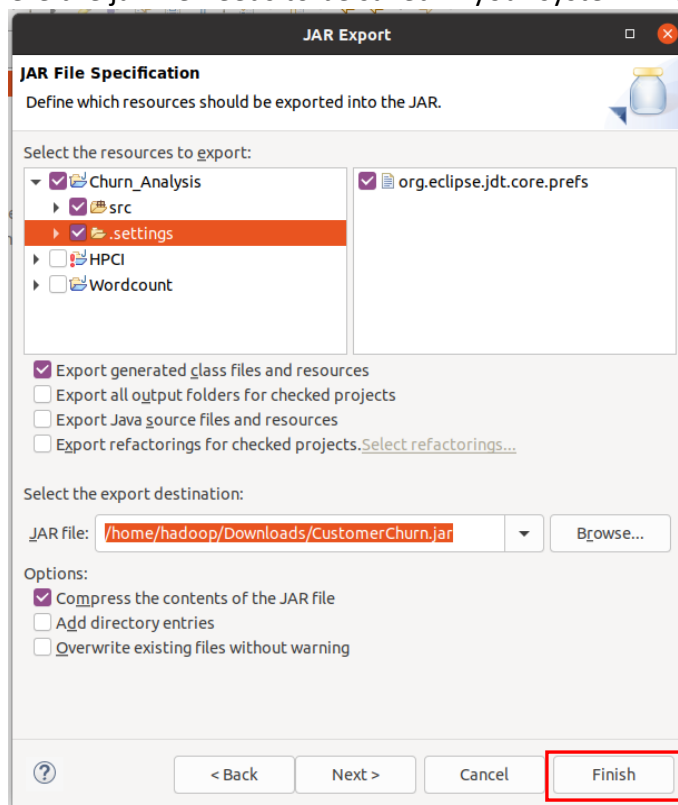
## 2.3 Export the Jar file

1. Right-click on the project (Churn_Analysis) and click on Export.



2. Click on the Next button

3. Select the resources to export (Churn_Analysis) and select the export destination where the jar file needs to be saved in your system. Hit the next button.



4. Select the main class and click on finish. It will export the jar in the given destination folder.
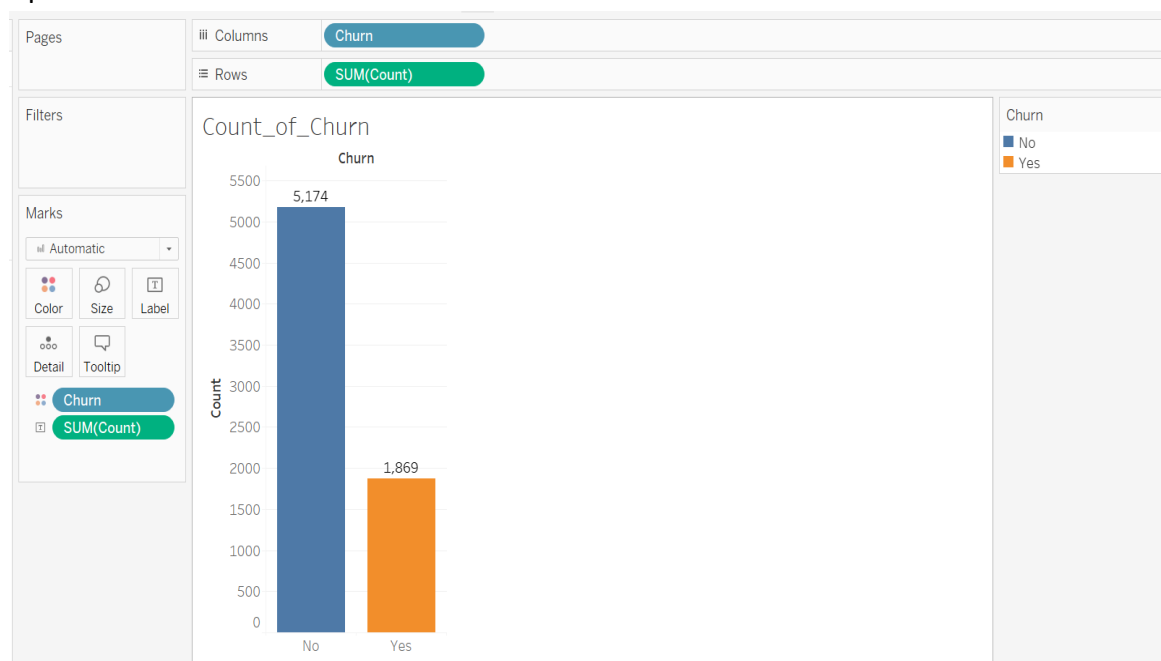
### 3. Results & Evaluation
1. The obtained output file 'part-r-00000' is in the format shown below.



```
(base) hadoop@hadoop-VirtualBox:~$ hdfs dfs -cat /Telecom/Output/part-r-00000
Month-to-month,No,No,1690
Month-to-month,No,Yes,1396
Month-to-month,Yes,No,530
Month-to-month,Yes,Yes,259
One year,No,No,825
One year,No,Yes,117
One year,Yes,No,482
One year,Yes,Yes,49
Two year,No,No,875
Two year,No,Yes,30
Two year,Yes,No,772
Two year,Yes,Yes,18
(base) hadoop@hadoop-VirtualBox:~$
```
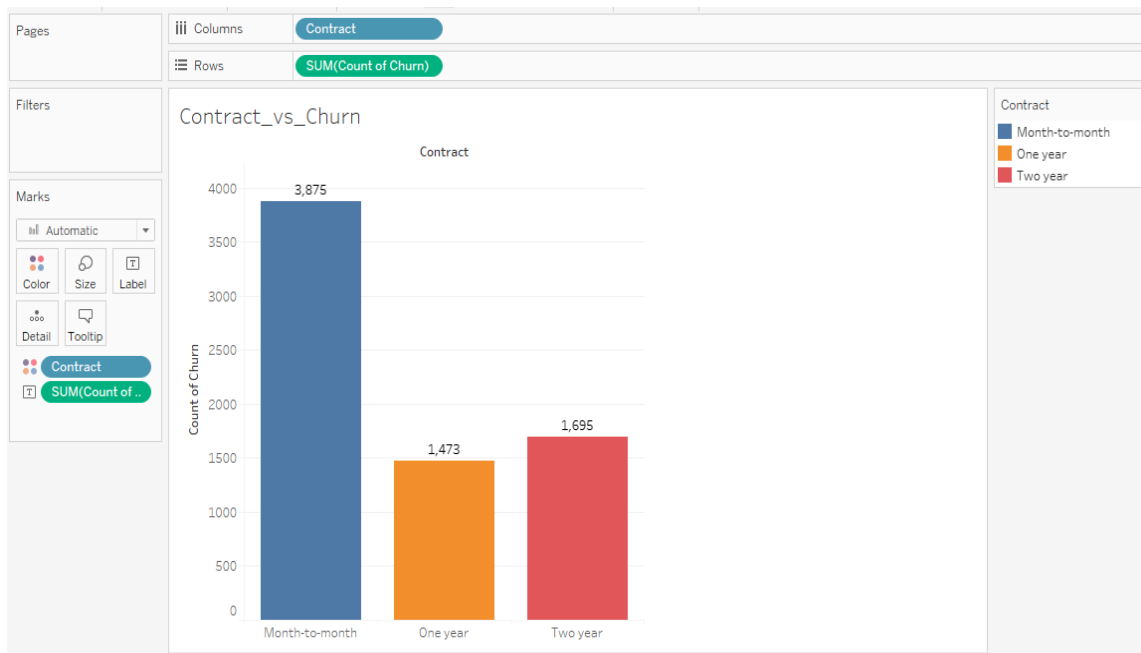
2. The output file has been converted to CSV and is imported into the tableau tool to visualize and identify the customer churned vary from contract period and people having dependents.

### 4. Observations:
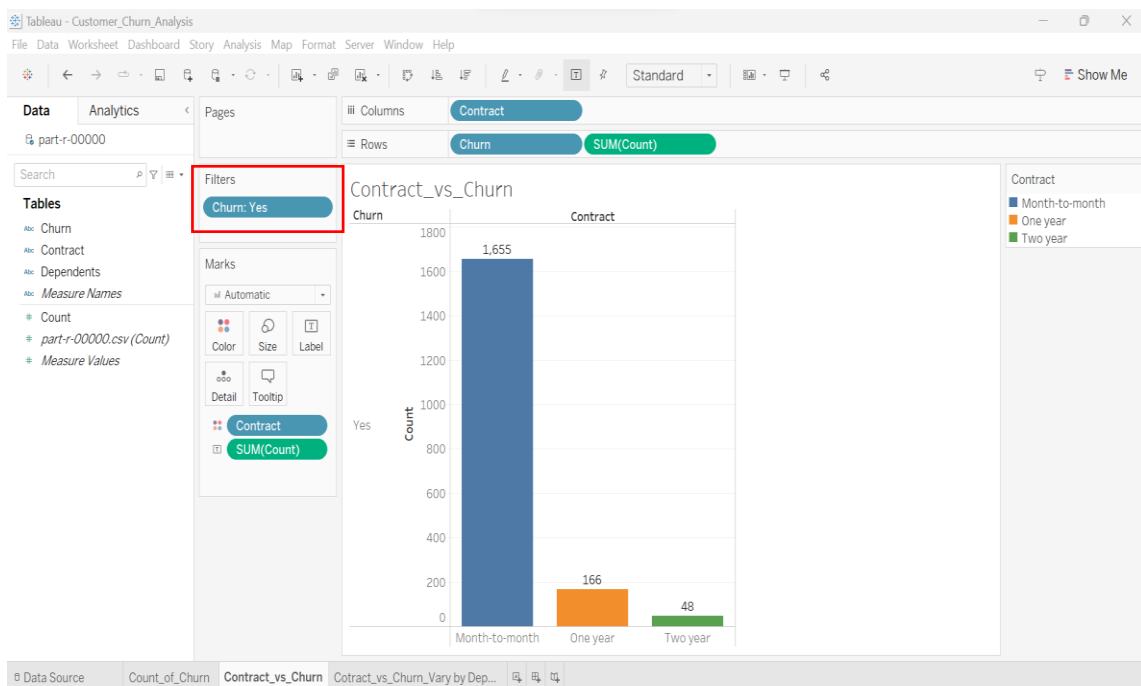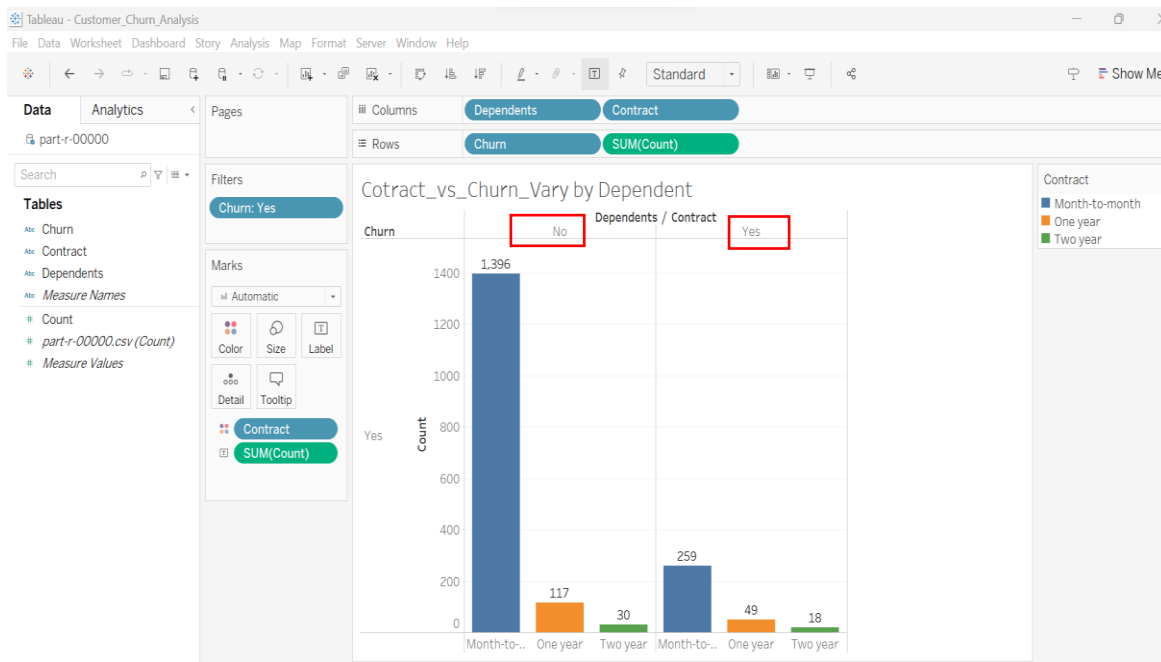a. It is observed that the company has lost around 25% of subscribers as shown in the plot below.



b. It is observed that most of the subscribers like to choose month-to-month contracts as shown in the plot below.

c. It is observed that customers having month-to-month contract periods are more likely to churn compared to customers having one-year or two-year contract periods.



d. It is observed that the customers who do not have dependents and have subscribed for less contract period (month-to-month) are more likely to churn compared to the others as shown in the plot below.

## Conclusion

With these observations, our research question is answered. It can be concluded that customers who subscribe for a month-to-month subscription, having no dependents, tend to cancel their subscription more frequently than those who opt for longer contract periods and have dependents. This difference is significant enough to affect the overall conclusion. Additionally, telecom companies can develop multiple strategies to improve customer retention and reduce churn. For example, companies can focus on offering long-term contract periods to customers, especially those without dependents. It not only develops their own business but also improves overall customer satisfaction.

## Appendix

As part of the coursework, The Mapreduce_ChurnAnalysis Word document has been attached containing all supporting documents as an appendix to the coursework. It includes all the contents, such as the java classes code, the original dataset used for this project, the renamed dataset, the Log4j file, the output file, and a visualization created in Tableau (.twb file) to evaluate the results. This comprehensive list of contents in a single word document is intended to support and enhance the project's overall quality, allowing for a thorough analysis and evaluation of the project's findings.

## References

hadoop.apache.org. (n.d.). Apache Hadoop 3.3.1 – Hadoop: Setting up a Single Node Cluster. [online] Available at: https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html.

hadoop.apache.org. (n.d.). MapReduce Tutorial. [online] Available at: https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html.

Apache.org. (2015). Apache log4j 1.2 - Download Apache log4j 1.2. [online] Available at: https://logging.apache.org/log4j/1.2/download.html [Accessed 17 Mar. 2023].

help.tableau.com. (n.d.). Aggregate Functions in Tableau. [online] Available at:

https://help.tableau.com/current/pro/desktop/enus/calculations_calculatedfields_aggregate_create.htm [Accessed on 09 Mar, 2023].

logging.apache.org. (n.d.). LogManager (Apache Log4j API). [online] Available at: https://logging.apache.org/log4j/2.x/log4j-api/apidocs/org/apache/logging/log4j/LogManager.html [Accessed on 08 Mar, 2023].

Apache.org. (2015). Apache log4j 1.2 - Download Apache log4j 1.2. [online] Available at: https://logging.apache.org/log4j/1.2/download.html.

www.tutorialspoint.com. (n.d.). Eclipse - Create Jar Files. [online] Available at: https://www.tutorialspoint.com/eclipse/eclipse_create_jar_files.htm# [Accessed 17 Apr. 2023].

GeeksforGeeks. (2019). HDFS Commands. [online] Available at: https://www.geeksforgeeks.org/hdfs-commands/ [Accessed on 08 Mar, 2023].