Q. Store and retrieve data in pig.

Aim :- To implement store and retrieve data in pig.

student-data.txt in the HDFS directory named (data)
with the following content.

001, Rajiv, Reddy, 9848022337, Hyderabad.
002, Siddarth, Battacharya, 9848022338, Kolkata
003, Rajesh, Khanna, 9848022359, Delhi
004, Preeti, Agarwal, 9848022330, Pune
005, Trupthi, mohanthy, 9848022336, Bhuwaneshwar

we can load data using the pig storage function as
shown below.

student = LOAD 'hdfs:// localhost:9000/pig-data/student
-data.txt' USING pigstorage (',') as (id: int,
firstname : chararray, lastname : chararray,

STORE student INTO 'hdfs:// localhost : 9000/pig
output]' USING pigstorage (',');

This will store the data into the given directory.

output

you can verify the stored data as shown below. First of
all, list out the files in the directory named pig-
output using ls command as shown below.

$ hdfs dfs-ls 'hdfs:// localhost:9000/pig-output'
found 2 items

rw-r--r-1 Hadoop supergroup 0 2020-10-05 -15:03
hdfs:// localhost:9000/pig-output/-success.

$ hdfs dfs-cat 'hdfs:// localhost:9000/pig-output'

1, Rajiv, Reddy, 9848022337, Hyderabad.
2, Siddarth, battacharya, 9848022338, Kolkata
3, Rajesh, Khanna, 9848022339, Delhi
4, Preeti, Agarwal, 9848022336, Bhuwanishwar
5, Trupthi, mohanthy, 9848022336, Pune.

Q. Perform social media analysis using cassandra.

**Aim:** To Implement social media analysis using cassandra.

Create a table for storing user posts.

```
CREATE TABLE social-media-posts (
Post-id uuid PRIMARY KEY,
user-id uuid,
post-text text,
post-time timestamp,
likes int,
shares int
);
```

Insert a post into the table

```
INSERT INTO social-media-posts (post-id, user-id,
post-text,
post-time,
likes,
shares)
VALUES (uuid(), uuid(), 'excited to be learning
about cassandra!', to timestamp (now()), 0, 0);
```

To find posts with more than 100 likes

```
SELECT * FROM social-media-posts WHERE likes > 100;
```

output

| post-id | likes | post-text |
|---------|-------|-----------|
| sddu4a87 | 101 | excited to be learning about cassandra! |

| shares | userid |
|--------|--------|
| 0 | dffbaffb-de |

Q. Buyer event analytics using cassandra on suitable product sales data.

Aim: Implement Buyer event analytics using Cassandra on suitable product sales data.

Create a table for storing product sales data events.

```
CREATE TABLE sales. product_events (
  event_id uuid PRIMARY KEY,
  product_id uuid,
  buyer_id uuid,
  event_time timestamp,
  event_type text,
  quantity int,
  price decimal
);
```

Insert a sales event into the table.

```
INSERT INTO sales. product_events ( event_id,
product_id,
buyer_id,
event_time,
event_type,
quantity, price)
VALUES ( uuid (), uuid(), to Timestamp ()), 'purchase'
1, 19.99)
```

output

```
SELECT * FROM sales. product_events
WHERE product_id = <specific_product_id>
AND event_type = "purchase";
```

| event_id | buyer_id | event_time | event_type |
|----------|----------|------------|------------|
| 781a003d | 64b96de1 | 2024-06-11 09:33:15051000 | Purchase |

| price | product_id | quantity |
|-------|------------|----------|
| 19.99 | 0b95658e-d15a | |

MGIT

8. Use R-project to carry out statistical analysis of big data.

Sample data for big-data.csv

product_ID, product_category, sales_Amount, Date

1. Electronics, 150, 2022-03-15
2. Clothing, 80, 2022-04-22
3. Books, 120, 2022-05-10
4. Home decor, 90, 2022-08-05
5. Electronics, 200, 2022-01-28
6. Clothing, 50, 2022-11-14
7. Books, 110, 2022-09-19
8. Home Decor, 70, 2022-04-03
9. Electronics, 180, 2022-06-30
10. Clothing, 70, 2022-10-17

```r
install.packages("dplyr")
install.packages("ggplot2")

library(dplyr)
library(ggplot2)

big_data <- read.csv("big-data.csv")
str(big_data)
total_sales <- big_data .1.>.1.
group_by(product_category) .1.>.1.
summarise(total_sales = sum(sales))
total_sales

ggplot(total_sales, aes(x=product_category,
y=total_sales)) +
geom_bar(stat="identity", fill="skyblue") +
labs(title = "Total sales by product category",
     x="product category",
     y = "Total sales") + theme_minimal()
```

output

| | product_category | Total_sales |
|---|---|---|
| 1 | Books | 230 |
| 2 | Clothing | 200 |
| 3 | Electronics | 530 |
| 4 | Home Decor | 160 |

9. Use R-project for data visualization of social media data.

social - media - data.csv

Text, Hashtags

" Excited to announce the launch of our new product!
" Check out our latest blog post about sustainability!
" We're hosting a webinar next week on digital marketing startegies join us!
" Happy friday everyone! Have a great weekend!
" Throwback to our team outing last summer.
" Excited to attend the conference next month!

library (ggplot2)
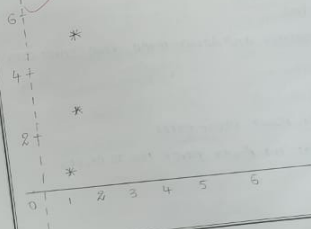
social_media_data<-read.csv("social-media-data.csv"

particular - hashtag <- "your - hashtag"
hashtag - data <- subset (social-media - data, grepl
(pasteo (" llb", particular-hashtag, "llb"),
hashtags, ignore.case = TRUE))

hashtag - frequency <- nrow (hashtag-data)

ggplot() +

geom - bar (data = NULL, aes(x = " ", y,hashtag -
frequency), fill = "skyblue", stat = " identity") +
theme - minimalt )

output

6. Using power point (excel) perform the following on any dataset.

a) Big data analytics.

b) Big data charting.

a) **Aim:** using power point to perform big data analytics.

procedure.

i) consider sample data

→ we use two datasets one is customer info table and the other is the order info table both have customerID as the common field.

ii) Getting excel power pivot Add In

→ open excel

→ select files > options

→ select add-ins

→ select the manage drop down menu, then select com Add-ins

→ select Go

→ select microsoft power pivot excel

→ select OK. it adds the power pivot tab to excel.

3) Adding data into the data model

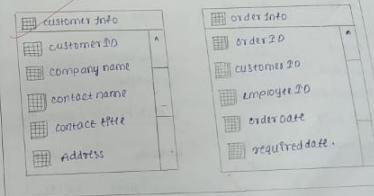→ select the range of customer into table. Then, Add to data model from power pivot tab.

→ you will notice that new pop up window appears. This is power pivot window.
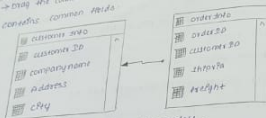
4) creating realationship between tables.

→ select the power pivot window Go to home then diagram view.

→ The imported table will appear as separate boxes in diagram view.



| Reset layout Display : ☑ columns ☑ calculated fields ☑ KPI |
|---|

| customer info | order info |
|---|---|
| customerID | orderID |
| company name | customer ID |
| contact name | employee ID |
| contact title | order date |
| Address | required date |

→ Drag the column heading from one table to another that contains common fields.



5) create pivot tables using data models.
→ In power pivot window go to home → pivot table.
→ create pivot table using dialog box will appear. select new worksheet and then select ok.
→ we can then calculate and various things.

| Row tables | 🔽 | sum of freight | Pivot table fields | -x |
|---|---|---|---|---|
| Alfreds | | 285·58 | Active /All | |
| Ana | | 77·42 | choose fields | |
| Blavis | | 60·3 | ☐ ship via | |
| Bon app | | 135·7·87 | ☑ freight | |
| Bottom dollas | | 231·31 | | |
| chop 3evy | | 267·11 | ▽ FILTERS | ▥ COLUMN |
| ble mond | | 63·7 | | |
| | | 832 | ▤ ROWS | Σ VALUES |
| | | | companyname | sum of freight |

---

6) Aim:- using power point to do key data charting.

Procedure:

1. Accessing power pivot
file → menu → options → Add·ins → microsoft power pivot for excel.

2. Importing data
click on power pivot tab to button → manage data → get external.
→ This are alot of options in the data source list. In this we will use data from another excel file.
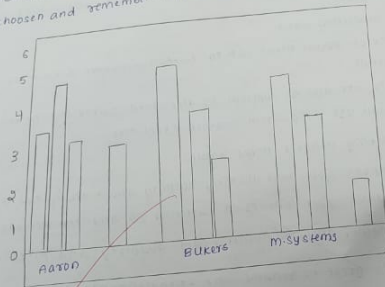
3. making a basic pivot table
→ suppose you have looking keeping about that has a item deals. each line is a customer ID and the sections are name, invoice number, cate, quantity and price.

4. It is great to rename the friendly name. hedec to a title that the informational index for this situation the title has been changed to invoices click finish.

5. making a pivot check a power pivot tables.
→ To make this from power pivot click the pivot table on the excel cheap and click pivot chart

→ Another exercise manual will open utilize the fields on the option to chosen fields.

→ So customer name, Date, and quantity have been choosen and remembered for pivot chart.

**Aim :** To implement a simple word count using map Reduce.

**Program :** 1) BDA Word Count Sample Driver. Java

```java
import java.io.IOException;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.InputFile
                                        InputFormat;
import org.apache.hadoop.mapreduce.lib.Input.Text
                                        InputFormat;
import org.apache.hadoop.mapreduce.lib.output.File
                                        OutputFormat;
import org.apache.hadoop.mapreduce.lib.output.Text
                                        OutputFormat;

public class BDAWordCountSampleDriver {
    public static void main (String []args) throws
                                    IOException,
        InterruptedException. ClassNotFoundException
    Job job =new Job();
    job.setJobName (" word counter');
```

```java
job.setJarByClass(BDA_WordCount_Sample_Driver.class);
job.setMapperClass(BDA_WordCount_Sample_Mapper.class);
job.setReducerClass(BDA_WordCount_Sample_Reducer.class);
job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);
FileInputFormat.addInputPath(Job.newPath(C(sample
                                     word text)));
FileOutputFormat.setOutputPath(C(Job.newpath(C(sample
                                     word count)));
System.exit(jb.waitForCompletion(true)?0:1);
}
}

// BDA Word Count Sample_Mapper.java
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
public class BDA_WordCount_Sample_Mapper extends
Mapper {
    @Override
    public void map(
```

```java
    protected void map(LongWritable key, Text value, Context
                                                       context)
    throws IOException, InterruptedException {
        String[] words = value.toString().split(",");
        for(String word : words){
            context.write(new Text(word), new IntWritable(1));
        }
    }
}

// BDA_WordCount_Sample_Reducer.java
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;
public class BDA_WordCount_Sample_Reducer extends
Reducer
<Text, IntWritable, Text, IntWritable>{
    @Override
    protected void reduce(Text word, Iterable<IntWritable>
    values, Context context)
    throws IOException, InterruptedException {
        integer count = 0;
        for(Int Writable val : values){
```

```
count += value.get();
}
context.write(word, new IntWritable(count));
}
}
```

**Aim:** To implement a simple map-reduce job that builds and inverted index on the set of input docs (Hadoop)

## Program

```
// mapper class
import java.io.IOException;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
public class InvertedIndexmapper extends mapper<long
                    Writable, Text, Text, Text> {
    private Text word = new Text();
    private Text docid = new Text();
    @Override
    protected void map(Longwritable key, text value, context
                        context)
    throws IOException, InterruptedException {
        string line = value.toString();
        string [] parts = line.split("\t");
        if (parts.length >= 2) {
            string DocId = parts[0];
            docid.set(DocId);
            string [] words = parts[1].split("");
```

```java
    for (string w: words) {
      word.set(w);
      context.write(word, docid);
    }
  }
}

// Reducer class
import java.io.IOException;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
public class InvertedIndexReducer extends Reducer<
   Text, Text, Text, Text> {
   @Override
   protected void reduce(Text key, Iterable<Text> values,
      context context)
   throws IOException, InterruptedException {
      string builder docList = new string builder();
      for (Text docid: values) {
         if (docList.length() > 0)
            docList.append(",");
```

```java
         docList.append(docid.to string());
      }
      context.write(key, new Text(docList.to string()));
   }
}

// main class
import org.apache.hadoop.fs.path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.File
   Input Format;
import org.apache.hadoop.mapreduce.lib.output.File
   Output Format;
public class InvertedIndex {
   public static void main(String[] args) throws IOException, Exception {
      if (args.length != 2) {
         System.err.println("usage: Inverted Index <in> <out>");
         System.exit(-1);
      }
      Job job = new Job();
      Job.setJarByClass(InvertedIndex.class);
      Job.setJobName("InvertedIndex");
```

```
File Input Format . add Input path (Job. new path (args(0)));
File output Format . set output path (Job. new path (args(1)));

Job . set mapper class (Inverted Index mapper . class);
Job . set Reducer class (Inverted Index Reducer . class);

Job . set output key class (Text . class);
Job . set output value class (Text . class);

system . exit (Job . wait For completion (true) ? 0 : 1);
}
}

compile your code and create a jar file . Then you can
run your map - reduce Job using hadoop with the following
command .

hadoop jar Inverted Index . jar Inverted Index < Input path>

<output path>

DOC 1 . txt
hello world
hello hadoop

DOC 2 . txt

hadoop . is a framework

hello map reduce .
```

| output | | |
|---|---|---|
| framework | DOC 2 . txt | |
| hello | DOC 1 . txt , DOC 2 . txt | |
| hadoop | DOC 1 . txt , DOC 2 . txt | |
| is | DOC 2 . txt | |
| mapreduce | DOC 2 . txt | |
| world | DOC 1 . txt . | |