

CII

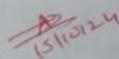
Group - 1
17

Thadomal Shahani Engineering College

Bandra (W.), Mumbai - 400 050.

© CERTIFICATE ©

Certify that Mr./Miss Harsik Bhagat
of COADS Department, Semester VII with
Roll No. 17 has completed a course of the necessary
experiments in the subject EDA under my
supervision in the **Thadomal Shahani Engineering College**
Laboratory in the year 2024 - 2025


15/01/24

Teacher In-Charge

Head of the Department

Date _____

Principal

CONTENTS

SR. NO.	EXPERIMENTS	PAGE NO.	DATE	TEACHERS SIGN.
1]	Installation of Hadoop & Experiment on HDFS commands	2	16/17	
2]	Use of Sqoop tool to transfer data between Hadoop & relational database server	10	23/17	
3]	Program execution in HBASE	18	30/17	
4]	Exp. per word counting using hadoop Map-Reduce	26	6/18	
5]	Exp. on pig	34	13/18	
6]	Create HIVE database & Descriptive analysis	42	27/18	
7]	Implement Bloomfilter using Python/R programming	50	31/19	
8]	Implement FM algorithm using Python/R programming	58	14/19	✓ (31/02/24)
9]	Data visualization using R	66	24/19	
10]	Mini Project	74	31/10	
11]	Assignment 1	82	11/10	
12]	Assignment 2	85	31/10	

EXPERIMENT NO:1

AIM: installation of Hadoop and experiment on HDFS commands.

Theory:

Cloudera is a company that provides a platform for big data processing and analytics, based on the Hadoop ecosystem. Hadoop is an open-source framework for distributed storage and processing of large datasets across clusters of computers. Hadoop Distributed File System (HDFS) is the primary storage system used in Hadoop. Here's a theoretical overview of Cloudera, Hadoop, and some common HDFS commands:

1. Cloudera: Cloudera is a software company that offers a comprehensive platform for big data management, analytics, and machine learning. Cloudera's platform is built on open-source Hadoop technologies and includes a range of tools and services to help organizations collect, store, process, and analyse vast amounts of data.
2. Hadoop: Hadoop is an open-source framework designed for distributed storage and processing of large datasets. It's widely used in various industries to handle big data challenges. The core components of Hadoop include HDFS (Hadoop Distributed File System) and the MapReduce processing engine. Hadoop is known for its scalability, fault tolerance, and flexibility.
3. Hadoop Distributed File System (HDFS): HDFS is the primary storage system in Hadoop. It's designed to store large files across a distributed cluster of commodity hardware. HDFS is based on a few key principles:
 - **Blocks:** Files are divided into fixed-size blocks (typically 128MB or 256MB) and are distributed across the cluster.
 - **Replication:** Each block is replicated across multiple nodes (usually three by default) to ensure data durability and fault tolerance.
 - **Master-Slave Architecture:** HDFS has a master node called the NameNode, which manages metadata, and multiple DataNodes that store the actual data blocks.

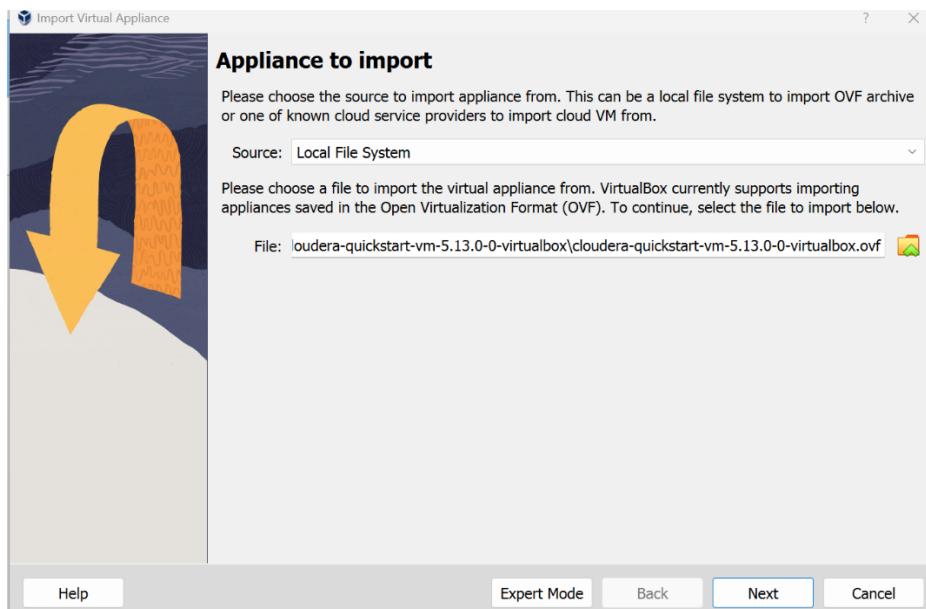
Hadoop Distributed File System (HDFS) commands are essential for managing data in a Hadoop cluster. Common commands include "ls" to list files and directories, "mkdir" to create directories, "copyFromLocal" to upload files.

Step 1:

Installing virtual box and Cloudera software.

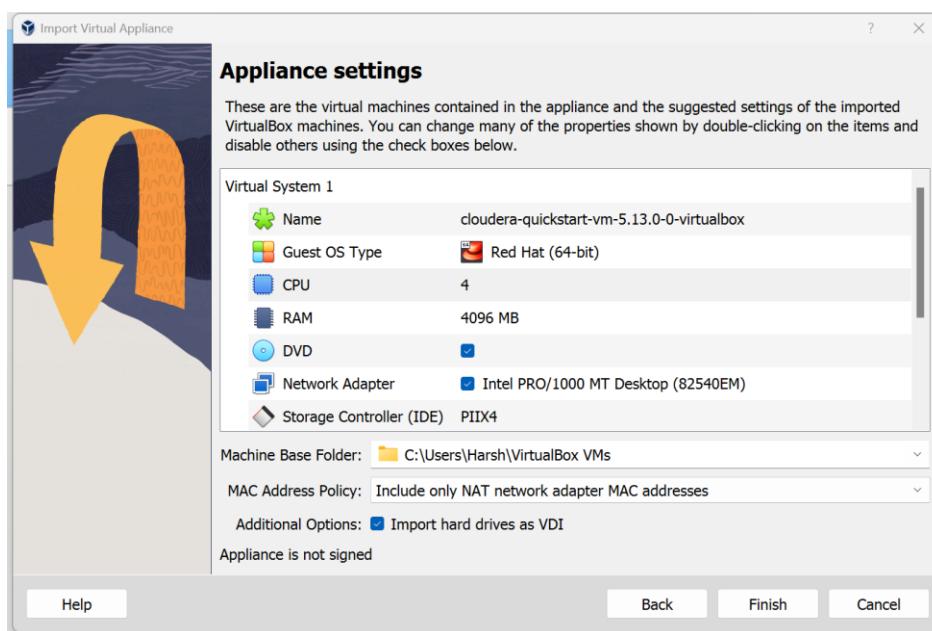
Step 2:

Unzipping the files and importing in the virtual machine. This is done by importing the link that we have downloaded.



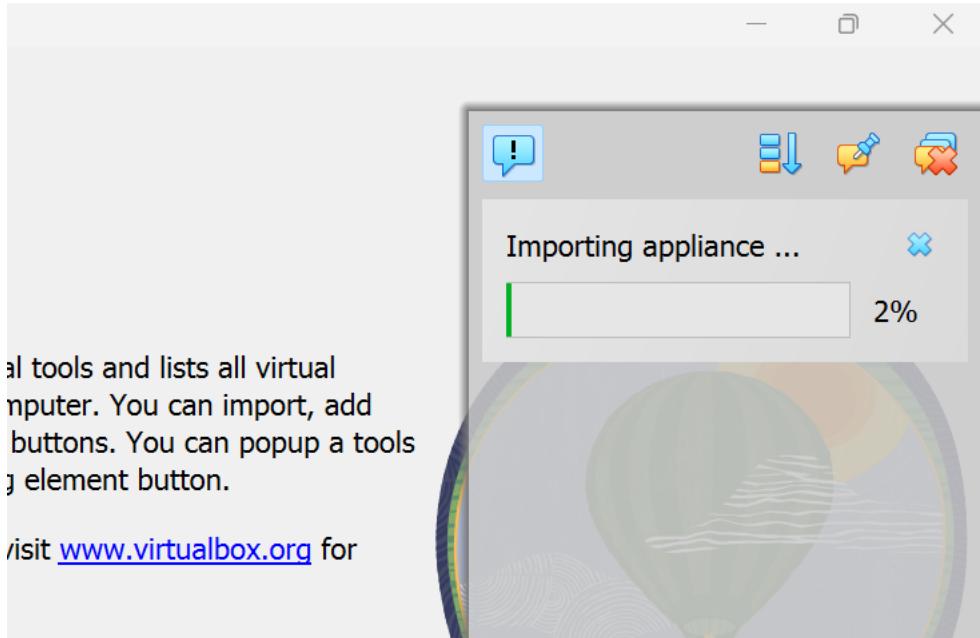
Step 3:

Selecting the import option. Also changing settings.



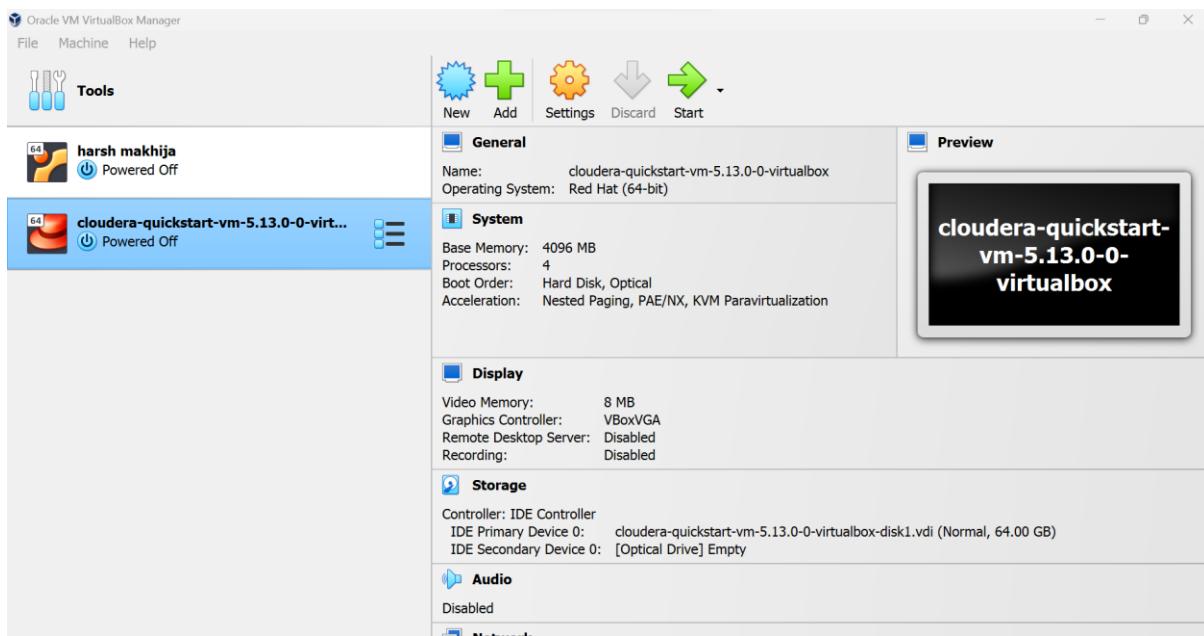
Step 4:

Importing the file and setting up.



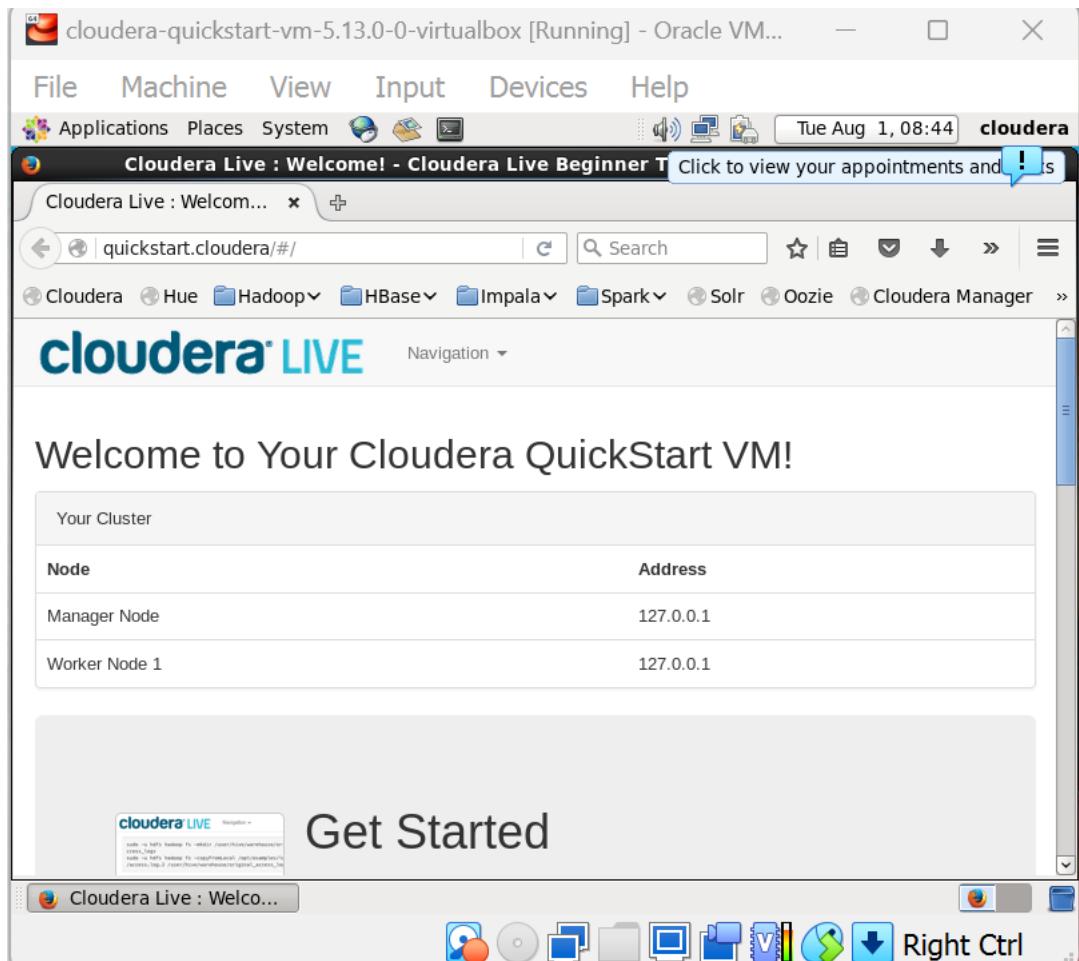
Step 5:

Import is successful, the screen will look something like this.



Step 6:

Starting the machine and loading Cloudera.



After successful completion the screen will look like this.



HDFS BASIC COMMANDS:

1. Hadoop version.

This command is used to check the installed version of the Hadoop framework on your system. It provides information about the Hadoop distribution, including its version number, build information, and other relevant details.

```
File Edit View Terminal Tabs Help
[training@localhost ~]$ # /home/training
[training@localhost ~]$ hadoop version
Hadoop 0.20.2-cdh3u2
Subversion file:///tmp/topdir/BUILD/hadoop-0.20.2-cdh3u2 -r 95a824e4005b2a
94fe1c11flef9db4c672ba43cb
Compiled by root on Thu Oct 13 21:51:41 PDT 2011
From source with checksum 644e5db6c59d45bca96cec7f220dda51
[training@localhost ~]$ hadoop fs -ls /
Found 7 items
drwxr-xr-x  - training supergroup          0 2017-02-02 02:23 /User
drwxr-xr-x  - hbase   supergroup          0 2018-01-22 21:56 /hbase
drwxr-xr-x  - training supergroup          0 2017-02-02 02:05 /home
drwxr-xr-x  - training supergroup          0 2023-07-25 01:23 /system
drwxrwxrwx  - hue    supergroup          0 2015-11-28 08:30 /tmp
drwxr-xr-x  - hue    supergroup          0 2017-03-22 23:33 /user
drwxr-xr-x  - mapred supergroup          0 2015-11-28 08:37 /var
[training@localhost ~]$ hadoop fs -df hdfs:/
Filesystem      Size  Used  Avail  Use%
hdfs:/        18611908608  95206223  12478259200  0%
[training@localhost ~]$ hadoop fs -count hdfs:/
      222       277     91843107 hdfs://localhost/
[training@localhost ~]$ hadoop fsck - /
fsck started by training from /127.0.0.1 for path / at Sun Jul 30 22:05:29
today July 30
```

2. Hadoop balancer

The balancer command is used to balance data distribution across Hadoop HDFS (Hadoop Distributed File System) data nodes.

```
[training@localhost ~]$ hadoop balancer
Time Stamp          Iteration#  Bytes Already Moved  Bytes Left To Mo
ve  Bytes Being Moved
23/07/30 22:11:59 INFO net.NetworkTopology: Adding a new node: /default-ra
ck/127.0.0.1:50010
23/07/30 22:11:59 INFO balancer.Balancer: 0 over utilized nodes:
23/07/30 22:11:59 INFO balancer.Balancer: 1 under utilized nodes:  127.0.0
.1:50010
The cluster is balanced. Exiting...
Balancing took 381.0 milliseconds
```

3. Making and removing directory.

These commands refer to basic file system operations in Hadoop HDFS. To make (create) a directory, you can use the **hadoop fs -mkdir** command followed by the directory path. To remove a directory, the **hadoop fs -rm -r** command is used, followed by the directory path.

```
[training@localhost ~]$ mkdir data/retail
[training@localhost ~]$ hadoop fs -put data/retail /user/training/hadoop
[training@localhost ~]$
```

```
[training@localhost ~]$ hadoop fs -rm hadoop/retail/customers
[training@localhost ~]$ ter: muted
[training@localhost ~]$ t remove hadoop/retail/customers: No such file or directory.
[training@localhost ~]$
```

4. Viewing directory retail

The **dus** command is used to determine the disk space usage (in bytes) of files and directories in Hadoop HDFS. It provides a summary of the storage consumed by each file and folder within the specified HDFS path, allowing administrators to monitor data storage utilization.

```
[training@localhost ~]$ hadoop fs -dus hadoop/retail
hdfs://localhost/user/training/hadoop/retail      0
[training@localhost:~]$
```

5. Creating user inside directory

These commands involve administrative actions. To create a user in Hadoop, one would typically use the underlying operating system's user management tools to add a new user. Similarly, creating a sample file would involve using Hadoop's file system commands, such as **hadoop fs -touchz** to create an empty file in HDFS.

```
[training@localhost ~]$ hadoop fs -ls /user/training/user
Found 1 items
drwxr-xr-x  - training supergroup          0 2015-09-12 02:33 /user/train
ing/user/training
[training@localhost ~]$
```

```
[training@localhost ~]$ hadoop fs -mkdir /retail/user/training/hadoop
[training@localhost ~]$ hadoop fs -ls
Found 45 items
-rw-r--r-- 1 training supergroup      1390 2015-10-02 20:20 /user/training/Books
drwxr-xr-x - training supergroup      0 2014-08-17 03:50 /user/training/WeatherData
drwxr-xr-x - training supergroup      0 2015-10-17 02:34 /user/training/_snapp
wse and run installed applications[ergroup
ing/a
drwxr-xr-x - training supergroup      0 2017-02-15 21:58 /user/training/apache_hadoop
-rw-r--r-- 1 training supergroup    2944 2015-09-26 02:37 /user/training/bookinfo
drwxr-xr-x - training supergroup      0 2015-09-20 01:38 /user/training/class2009_dir1
-rw-r--r-- 1 training supergroup     81 2015-09-18 21:05 /user/training/deptinfo
drwxr-xr-x - training supergroup      0 2014-08-17 01:59 /user/train
```

6. Making sample file

The **-put** command is a fundamental utility in Hadoop that is used to copy data from the local file system to the Hadoop Distributed File System (HDFS).

```
[training@localhost ~]$ hadoop fs -put data/sample.txt/user/training/hadoop
Usage: java FsShell [-put <localsrc> ... <dst>]
```

7. Hadoop count

Viewing a directory in Hadoop typically involves listing the contents of a directory within HDFS. The **Hadoop fs -ls** command is commonly used for this purpose. It provides a list of files and subdirectories, including their permissions, modification times, and file sizes, within the specified HDFS directory.

```
[training@localhost ~]$ hadoop fs -count hdfs:/
      228          278        91843853 hdfs://localhost/
[training@localhost ~]$
```

-rw-r--r--	1	training	supergroup	44	2017-03-22	23:31	/user/training/i
nputWC.txt							
-rw-r--r--	1	training	supergroup	69	2015-10-10	02:35	/user/training/j
oin							
-rw-r--r--	1	training	supergroup	105	2015-10-10	02:38	/user/training/j
oin2							
drwxr-xr-x	-	training	supergroup	0	2017-03-22	23:33	/user/training/o
utputWC.txt							
-rw-r--r--	1	training	supergroup	16	2015-10-04	01:23	/user/training/p
ig							
-rw-r--r--	1	training	supergroup	32	2015-10-04	02:04	/user/training/p
ig1							
-rw-r--r--	1	training	supergroup	90	2015-11-23	03:12	/user/training/p
oem							
-rw-r--r--	1	training	supergroup	90	2015-09-12	02:41	/user/training/p
oem912							
drwxr-xr-x	-	training	supergroup	0	2015-10-10	03:53	/user/training/s
tr							
drwxr-xr-x	-	training	supergroup	0	2015-10-17	01:26	/user/training/s
tud							
drwxr-xr-x	-	training	supergroup	0	2015-10-17	01:42	/user/training/s
Found 46 items							
-rw-r--r--	1	training	supergroup	1390	2015-10-02	20:20	/user/training/B
ooks							
drwxr-xr-x	-	training	supergroup	0	2014-08-17	03:50	/user/training/W
eatherData							
drwxr-xr-x	-	training	supergroup	0	2015-10-17	02:34	/user/training/_
sqoop							
-rw-r--r--	1	training	supergroup	23	2015-11-29	23:36	/user/training/a
drwxr-xr-x	-	training	supergroup	0	2017-02-15	21:58	/user/training/a
pache_hadoop							
-rw-r--r--	1	training	supergroup	2944	2015-09-26	02:37	/user/training/b
ookinfo							
drwxr-xr-x	-	training	supergroup	0	2015-09-20	01:38	/user/training/c
lass2009_dir1							
-rw-r--r--	1	training	supergroup	81	2015-09-18	21:05	/user/training/d
eptrinfo							
drwxr-xr-x	-	training	supergroup	0	2014-08-17	01:59	/user/training/d
ir1							
drwxr-xr-x	-	training	supergroup	0	2015-09-12	02:42	/user/training/d
ir10							
drwxr-xr-x	-	training	supergroup	0	2015-09-19	02:00	/user/training/d
[training@localhost ~]\$ hadoop fs -df hdfs:/							
Filesystem	Size	Used	Avail	Use%			
hdfs:/	18611908608	95203328	12473774080	0%			
[training@localhost ~]\$							

EXPERIMENT NO: 2

AIM: Use of Sqoop tool to transfer data between Hadoop and relational database servers.

Theory:

Sqoop is a versatile and powerful tool used to facilitate data transfer between Hadoop-based systems, such as the Hadoop Distributed File System (HDFS), and relational database management systems (RDBMS). Here's a theoretical overview of the use of Sqoop for transferring data between Hadoop and RDBMS:

1. **Data Integration Needs:** Many organizations need to combine and analyze data from both their traditional RDBMS and Hadoop-based systems to gain valuable insights. This necessitates the need for a tool that can efficiently transfer data between these different storage and processing environments.
2. **Sqoop Overview:** Sqoop, short for "SQL to Hadoop," is an open-source data transfer tool specifically designed for this purpose. It acts as a bridge between Hadoop and RDBMS, enabling bi-directional data transfers.
3. **Data Import:** When importing data from an RDBMS to Hadoop, Sqoop reads data from the RDBMS, converts it to Hadoop-friendly formats (e.g., Avro or Parquet), and stores it in HDFS. Users can specify various options, such as the source database connection, target HDFS directory, and data transformation methods.
4. **Data Export:** When exporting data from Hadoop to an RDBMS, Sqoop reads data from HDFS, transforms it if necessary, and loads it into the RDBMS. Users can specify the target database connection, source HDFS directory, and other export options.
5. **Use Cases:** Sqoop is valuable for various use cases, including:
 1. Loading operational data into Hadoop for analytics.
 2. Moving data from Hadoop to an RDBMS for reporting.
 3. Data migration and replication between different RDBMS systems.
 4. Continuous data integration and synchronization.

Step 1: open cluders and connect to mysql and view databases.

```
[training@localhost ~]$ mysql -u root -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 4
Server version: 5.0.77 Source distribution

Type 'help;' or '\h' for help. Type '\c' to clear the buffer.

mysql> create database bank;
ERROR 1007 (HY000): Can't create database 'bank'; database exists
mysql> create database bank1;
Query OK, 1 row affected (0.00 sec)

mysql> show databases;
+-----+
| Database      |
+-----+
| information_schema |
| bank          |
| bank1         |
| db1           |
| db2           |
| example        |
| hivemetastore |
| movielens      |
| mysql          |
+-----+
10 rows in set (0.00 sec)

ning@localhost:~
```

Step 2: create database:

```
mysql> use bank1;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
mysql> create table registercopy(accno varchar(20), name varchar(20), number va
Query OK, 0 rows affected (0.04 sec)

mysql> describe registercopy;
+-----+-----+-----+-----+-----+
| Field    | Type     | Null | Key | Default | Extra |
+-----+-----+-----+-----+-----+
| accno    | varchar(20) | YES |   | NULL    |       |
| name     | varchar(20) | YES |   | NULL    |       |
| number   | varchar(20) | YES |   | NULL    |       |
| email    | varchar(20) | YES |   | NULL    |       |
| password | varchar(20) | YES |   | NULL    |       |
| age      | varchar(20) | YES |   | NULL    |       |
+-----+-----+-----+-----+-----+
6 rows in set (0.00 sec)
```

Step 3: insert values

```
+-----+-----+-----+-----+
| accno | varchar(20) | YES |      | NULL   |
| name  | varchar(20)  | YES |      | NULL   |
| number | varchar(20) | YES |      | NULL   |
| email  | varchar(20)  | YES |      | NULL   |
| password | varchar(20) | YES |      | NULL   |
| age    | varchar(3)   | YES |      | NULL   |
+-----+-----+-----+-----+
6 rows in set (0.00 sec)

mysql> insert into register values("1","harsh","21","harsh@email","harsh"
,"21");
Query OK, 1 row affected (0.00 sec)

mysql> insert into register values("2","disha","22","dishah@email","disha","20");
Query OK, 1 row affected (0.00 sec)

mysql> insert into register values("3","sahil","23","sahil@email","sahil","20");
Query OK, 1 row affected (0.00 sec)

mysql> insert into register values("4","ishwar","24","ishwar@email","ishwar","20");
Query OK, 1 row affected (0.00 sec)

mysql> insert into register values("5","pankaj","25","pankaj@email","pankaj","20");
Query OK, 1 row affected (0.00 sec)
```

Step 4: view data

```
mysql> select * from registercopy;
+-----+-----+-----+-----+-----+
| accno | name   | number | email     | password | age   |
+-----+-----+-----+-----+-----+
| 1     | harsh  | 21    | harsh@email | harsh    | 21   |
| 2     | disha  | 22    | dishah@email | disha   | 20   |
| 3     | sahil  | 23    | sahil@email | sahil   | 20   |
| 4     | ishwar | 24    | ishwar@email | ishwar  | 20   |
| 5     | pankaj | 25    | pankaj@email | pankaj  | 20   |
+-----+-----+-----+-----+-----+
5 rows in set (0.00 sec)
```

Step 5: exporting the data from hdfs to mysql

Syntax:

```
Sqoop export --connect jdbc:mysql://localhost/db --username root --
table<table_name> --export-dir<directory>
```

```
[cloudera@quickstart ~]$ sqoop export --connect jdbc:mysql://localhost/bank --username root --password cloudera --table registercopy --export-dir /home/cloudera/myfirstdata
Warning: /usr/lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
23/09/13 05:25:52 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.13.0
23/09/13 05:25:52 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
23/09/13 05:25:52 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
23/09/13 05:25:52 INFO tool.CodeGenTool: Beginning code generation
23/09/13 05:25:53 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `registercopy` AS t LIMIT 1
23/09/13 05:25:53 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `registercopy` AS t LIMIT 1
23/09/13 05:25:53 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/lib/hadoop-mapreduce
Note: /tmp/sqoop-cloudera/compile/1651b0eae1f9ed0cb04ac1839851920a/registercopy.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
23/09/13 05:26:00 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-cloudera/compile/1651b0eae1f9ed0cb04ac1839851920a/registercopy.jar
23/09/13 05:26:00 INFO mapreduce.ExportJobBase: Beginning export of registercopy
23/09/13 05:26:00 INFO Configuration.deprecation: mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
23/09/13 05:26:01 INFO Configuration.deprecation: mapred.jar is deprecated. Instead, use mapreduce.job.jar
23/09/13 05:26:03 INFO Configuration.deprecation: mapred.reduce.tasks.speculative.execution is deprecated. Instead, use mapreduce.reduce.speculative
23/09/13 05:26:03 INFO Configuration.deprecation: mapred.map.tasks.speculative.execution is deprecated. Instead, use mapreduce.map.speculative
23/09/13 05:26:03 INFO Configuration.deprecation: mapred.map.tasks is deprecated. Instead, use mapreduce.job.maps
23/09/13 05:26:03 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
23/09/13 05:26:08 INFO input.FileInputFormat: Total input paths to process : 1
23/09/13 05:26:08 INFO input.FileInputFormat: Total input paths to process : 1
23/09/13 05:26:08 INFO mapreduce.JobSubmitter: number of splits:4
23/09/13 05:26:08 INFO Configuration.deprecation: mapred.map.tasks.speculative.execution is deprecated. Instead, use mapreduce.map.speculative
23/09/13 05:26:09 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1694606951119_0001
23/09/13 05:26:11 INFO impl.YarnClientImpl: Submitted application application_1694606951119_0001
23/09/13 05:26:11 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1694606951119_0001/
23/09/13 05:26:11 INFO mapreduce.Job: Running job: job_1694606951119_0001
23/09/13 05:26:42 INFO mapreduce.Job: Job job_1694606951119_0001 running in uber mode : false
23/09/13 05:26:42 INFO mapreduce.Job: map 0% reduce 0%
23/09/13 05:27:08 INFO mapreduce.Job: map 100% reduce 0%
23/09/13 05:27:09 INFO mapreduce.Job: Job job_1694606951119_0001 completed successfully
23/09/13 05:27:09 INFO mapreduce.Job: Counters: 30
File System Counters
```

```

4606951119_0001/
23/09/13 05:26:11 INFO mapreduce.Job: Running job: job_1694606951119_0001
23/09/13 05:26:42 INFO mapreduce.Job: Job job_1694606951119_0001 running in uber mode : false
23/09/13 05:26:42 INFO mapreduce.Job: map 0% reduce 0%
23/09/13 05:27:08 INFO mapreduce.Job: map 100% reduce 0%
23/09/13 05:27:09 INFO mapreduce.Job: Job job_1694606951119_0001 completed successfully
23/09/13 05:27:09 INFO mapreduce.Job: Counters: 30
    File System Counters
        FILE: Number of bytes read=0
        FILE: Number of bytes written=683776
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=1209
        HDFS: Number of bytes written=0
        HDFS: Number of read operations=19
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=0
    Job Counters
        Launched map tasks=4
        Data-local map tasks=4
        Total time spent by all maps in occupied slots (ms)=85742
        Total time spent by all reduces in occupied slots (ms)=0
        Total time spent by all map tasks (ms)=85742
        Total vcore-milliseconds taken by all map tasks=85742
        Total megabyte-milliseconds taken by all map tasks=87799808
    Map-Reduce Framework
        Map input records=5
        Map output records=5
        Input split bytes=691
        Spilled Records=0
        Failed Shuffles=0
        Merged Map outputs=0
        GC time elapsed (ms)=504
        CPU time spent (ms)=4650
        Physical memory (bytes) snapshot=691548160
        Virtual memory (bytes) snapshot=6285340672
        Total committed heap usage (bytes)=695730176
    File Input Format Counters
        Bytes Read=0
    File Output Format Counters
        Bytes Written=0
23/09/13 05:27:09 INFO mapreduce.ExportJobBase: Transferred 1.1807 KB in 66.0434 seconds (18.3062 bytes/sec)
23/09/13 05:27:09 INFO mapreduce.ExportJobBase: Exported 5 records.
[clooudera@quickstart ~]$ █

```

Step 6: we have successfully exported the data

```

mysql> select * from register;
+-----+-----+-----+-----+-----+
| accmo | name   | number | email      | password | age   |
+-----+-----+-----+-----+-----+
| 1     | harsh  | 21    | harsh@email | harsh    | 21   |
| 2     | disha  | 22    | dishah@email | disha    | 20   |
| 3     | sahil  | 23    | sahil@email | sahil    | 20   |
| 4     | ishwar | 24    | ishwar@email | ishwar   | 20   |
| 5     | pankaj | 25    | pankaj@email | pankaj   | 20   |
+-----+-----+-----+-----+-----+
5 rows in set (0.00 sec)

```

EXPERIMENT NO: 3

AIM: Programming exercises on HBASE.

Theory:

HBase, short for Hadoop Database, is an open-source, distributed, and scalable NoSQL database that is designed to handle large volumes of sparse data. It is part of the Hadoop ecosystem and is often used for real-time, read and write access to big data. Here's a comprehensive overview of HBase:

Key Characteristics and Concepts:

1. **NoSQL Database:** HBase is a NoSQL database, which means it is not based on traditional relational database models. Instead, it uses a schema-less design, allowing for flexible data storage.
2. **Distributed and Scalable:** HBase is built on the Hadoop Distributed File System (HDFS) and is designed to scale horizontally, making it capable of handling petabytes of data spread across a cluster of commodity hardware.
3. **Columnar Storage:** HBase stores data in tables, similar to traditional databases. However, the data within tables is stored in a columnar format, which enables efficient data retrieval and compression.
4. **Sparse Data Model:** HBase is optimized for sparse data, meaning that it is well-suited for situations where data is missing or not present in every record. This makes it particularly useful for handling semi-structured or rapidly changing data.
5. **Strong Consistency:** HBase provides strong consistency guarantees, which means that data is consistent and up-to-date across the entire cluster.
6. **High Write Throughput:** It excels in write-heavy workloads, making it ideal for applications with frequent data inserts and updates.

HBase is a powerful NoSQL database designed for handling large-scale, sparse data and real-time access. It plays a significant role in the Hadoop ecosystem, complementing other tools for big data processing and analytics. Its distributed architecture, scalability, and flexibility make it a valuable choice for various applications, particularly those requiring high write throughput and real-time data access.

STEP 1:

Starting with HBase shell, creating ‘register’ and learning how to disable and enable the table.

syntax

Create ‘table’ ‘columns’, ‘columns’

```
[training@localhost ~]$ hbase shell
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 0.90.4-cdh3u2, r, Thu Oct 13 20:32:26 PDT 2011

hbase(main):001:0> create 'table', 'name', 'age', 'gender'
0 row(s) in 0.2490 seconds

hbase(main):002:0> create 'register', 'accno', 'name', 'password' , 'email' , 'age'
0 row(s) in 1.0760 seconds

hbase(main):003:0> list
TABLE
emp
etash_register
hemp
hr_class_1810
register
table
6 row(s) in 0.0330 seconds

hbase(main):004:0> disable 'register'
0 row(s) in 2.0330 seconds

hbase(main):005:0> is_disabled 'register'
true
0 row(s) in 0.0110 seconds

hbase(main):006:0> enable 'register'
0 row(s) in 2.0310 seconds

hbase(main):007:0> describe 'register'
```

Step 2: accessing various columns and describing table:

Syntax:

Describe 'register'

```
hbase(main):006:0> enable 'register'
0 row(s) in 2.0310 seconds

hbase(main):007:0> describe 'register'
ERROR: Failed to find table named register

Here is some help for this command:
Describe the named table. For example:
  hbase> describe 't1'

hbase(main):008:0> describe 'register'
DESCRIPTION                                                 ENABLED
{NAME => 'register', FAMILIES => [{NAME => 'accno', BLOOMFILTER => 'NONE', REPLICATION_SCOPE => '0' true
, COMPRESSION => 'NONE', VERSIONS => '3', TTL => '2147483647', BLOCKSIZE => '65536', IN_MEMORY => '
false', BLOCKCACHE => 'true'}, {NAME => 'age', BLOOMFILTER => 'NONE', REPLICATION_SCOPE => '0', COM
PRESSION => 'NONE', VERSIONS => '3', TTL => '2147483647', BLOCKSIZE => '65536', IN_MEMORY => 'false
', BLOCKCACHE => 'true'}, {NAME => 'email', BLOOMFILTER => 'NONE', REPLICATION_SCOPE => '0', COMPRE
SSION => 'NONE', VERSIONS => '3', TTL => '2147483647', BLOCKSIZE => '65536', IN_MEMORY => 'false',
BLOCKCACHE => 'true'}, {NAME => 'name', BLOOMFILTER => 'NONE', REPLICATION_SCOPE => '0', COMPRESSIO
N => 'NONE', VERSIONS => '3', TTL => '2147483647', BLOCKSIZE => '65536', IN_MEMORY => 'false', BLOC
KCACHE => 'true'}, {NAME => 'password', BLOOMFILTER => 'NONE', REPLICATION_SCOPE => '0', COMPRESSIO
N => 'NONE', VERSIONS => '3', TTL => '2147483647', BLOCKSIZE => '65536', IN_MEMORY => 'false', BLOC
KCACHE => 'true'}]}
1 row(s) in 0.0260 seconds

hbase(main):009:0> alter 'register' ,NAME => 'name', VERSION => 5
ERROR: Table register is enabled. Disable it first before altering.

Here is some help for this command:
Alter column family schema; pass table name and a dictionary
```

Step 3: working on various commands

Alter' table_name'

Drop"column_name"

Put "table_name' , 'data'

```
hbase(main):013:0> disable 'register'
0 row(s) in 2.0420 seconds

hbase(main):014:0> alter 'register',NAME => 'name',VERSION => 5
0 row(s) in 0.0290 seconds

hbase(main):015:0> exists 'register'
Table register does exist
0 row(s) in 0.0080 seconds

hbase(main):016:0> drop 'register'
0 row(s) in 1.0380 seconds

hbase(main):017:0> create 'register' , 'personal data' , 'account data'
0 row(s) in 1.0850 seconds

hbase(main):018:0> put 'register' , '1','personal nme:name','raj'

ERROR: org.apache.hadoop.hbase.client.RetryExhaustedException: Failed 1 action
host.localdomain:50462,

Here is some help for this command:
Put a cell 'value' at specified table/row/column and optionally
timestamp coordinates. To put a cell value into table 't1' at
row 'r1' under column 'c1' marked with the time 'ts1', do:

    hbase> put 't1', 'r1', 'c1', 'value', ts1

hbase(main):019:0> put 'register' , '1','personal data:name','harsh'
0 row(s) in 0.0180 seconds
```

Step 4: truncating table and other commands :

```
Here is some help for this command:  
Put a cell 'value' at specified table/row/column and optionally  
timestamp coordinates. To put a cell value into table 't1' at  
row 'r1' under column 'c1' marked with the time 'ts1', do:  
  
hbase> put 't1', 'r1', 'c1', 'value', ts1  
  
hbase(main):019:0> put 'register' , '1','personal data:name','harsh'  
0 row(s) in 0.0180 seconds  
  
hbase(main):020:0> put 'register' , '1','personal data:age','21'  
0 row(s) in 0.0060 seconds  
  
hbase(main):021:0> put 'register' , '1','personal data:email','harsh@21'  
0 row(s) in 0.0050 seconds  
  
hbase(main):022:0> put 'register' , '1','personal data:accno','21'  
0 row(s) in 0.0040 seconds  
  
hbase(main):023:0> scan 'register'  
ROW                                     COLUMN+CELL  
1                                         column=personal data:accno, timestamp=1697268915750, value=21  
1                                         column=personal data:age, timestamp=1697268881268, value=21  
1                                         column=personal data:email, timestamp=1697268893044, value=harsh@21  
1                                         column=personal data:name, timestamp=1697268818994, value=harsh  
1 row(s) in 0.0200 seconds  
  
hbase(main):024:0> get 'register' '1'  
ERROR: wrong number of arguments (1 for 2)
```

```
hbase(main):025:0> get 'register', '1'  
COLUMN                                     CELL  
personal data:accno                         timestamp=1697268915750, value=21  
personal data:age                           timestamp=1697268881268, value=21  
personal data:email                          timestamp=1697268893044, value=harsh@21  
personal data:name                          timestamp=1697268818994, value=harsh  
4 row(s) in 0.0260 seconds  
  
hbase(main):026:0> delete 'register' , '1','personal data:name',1661584013135  
0 row(s) in 0.0190 seconds  
  
hbase(main):027:0> count 'register'  
1 row(s) in 0.0260 seconds  
  
hbase(main):028:0> truncate 'register'  
Truncating 'register' table (it may take a while):  
- Disabling table...  
- Dropping table...  
- Creating table...  
0 row(s) in 3.2210 seconds  
  
hbase(main):029:0>
```

EXPERIMENT NO: 4

AIM: Experiment on word counting using Hadoop map-reduce

Theory:

Word counting using Hadoop MapReduce is a fundamental example of how distributed computing can process vast amounts of data efficiently. This experiment demonstrates how to leverage the power of Hadoop, a popular big data framework, to analyse and count the occurrences of words within large datasets. Below is a comprehensive explanation of this experiment:

Introduction to Word Counting with Hadoop MapReduce:

Hadoop MapReduce is a programming model and processing framework that allows users to perform distributed data processing tasks. One of the most common introductory exercises in Hadoop MapReduce is word counting, which involves determining the frequency of each word in a given set of documents or text data.

Experiment Setup:

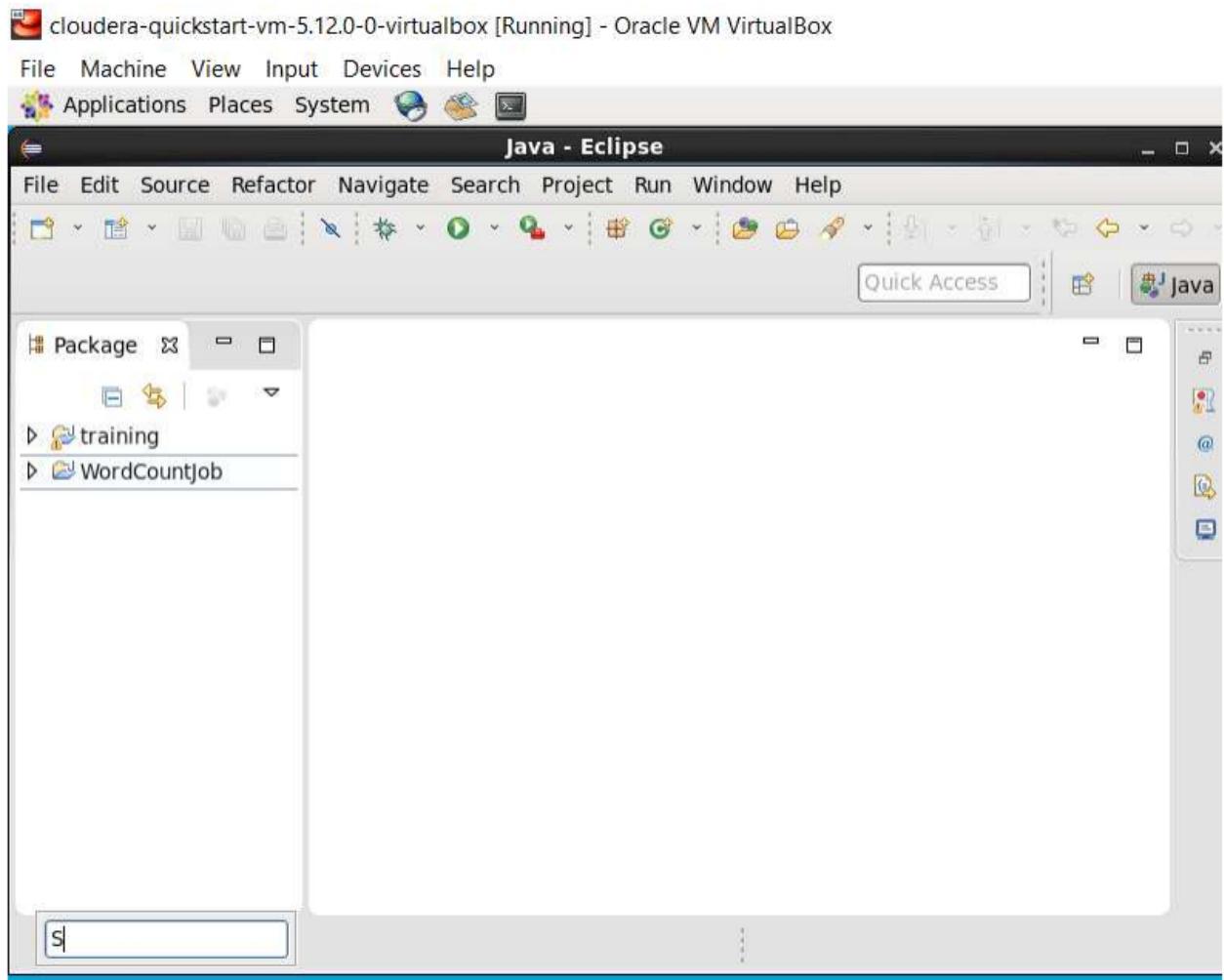
1. **Data Input:** The first step in this experiment is to have a set of input text documents. These documents could be stored in Hadoop HDFS (Hadoop Distributed File System), and they are divided into smaller chunks called blocks.
2. **Map Phase:** The MapReduce process starts with the "Map" phase, where each Mapper task processes a portion of the input data. Mappers read the text, tokenize it into words, and emit key-value pairs where the key is the word and the value is usually 1, signifying the occurrence of that word.
3. **Shuffle and Sort:** After the Map phase, the framework groups all the emitted key-value pairs by keys and sorts them. This phase ensures that all occurrences of the same word are grouped together.
4. **Reduce Phase:** The "Reduce" phase takes the grouped and sorted key-value pairs and applies a "Reducer" function to aggregate the counts. In the context of word counting, the Reducer sums up the values for each word key, giving the final count of occurrences.

STEP 1 :

Create New Java Project in Eclipse

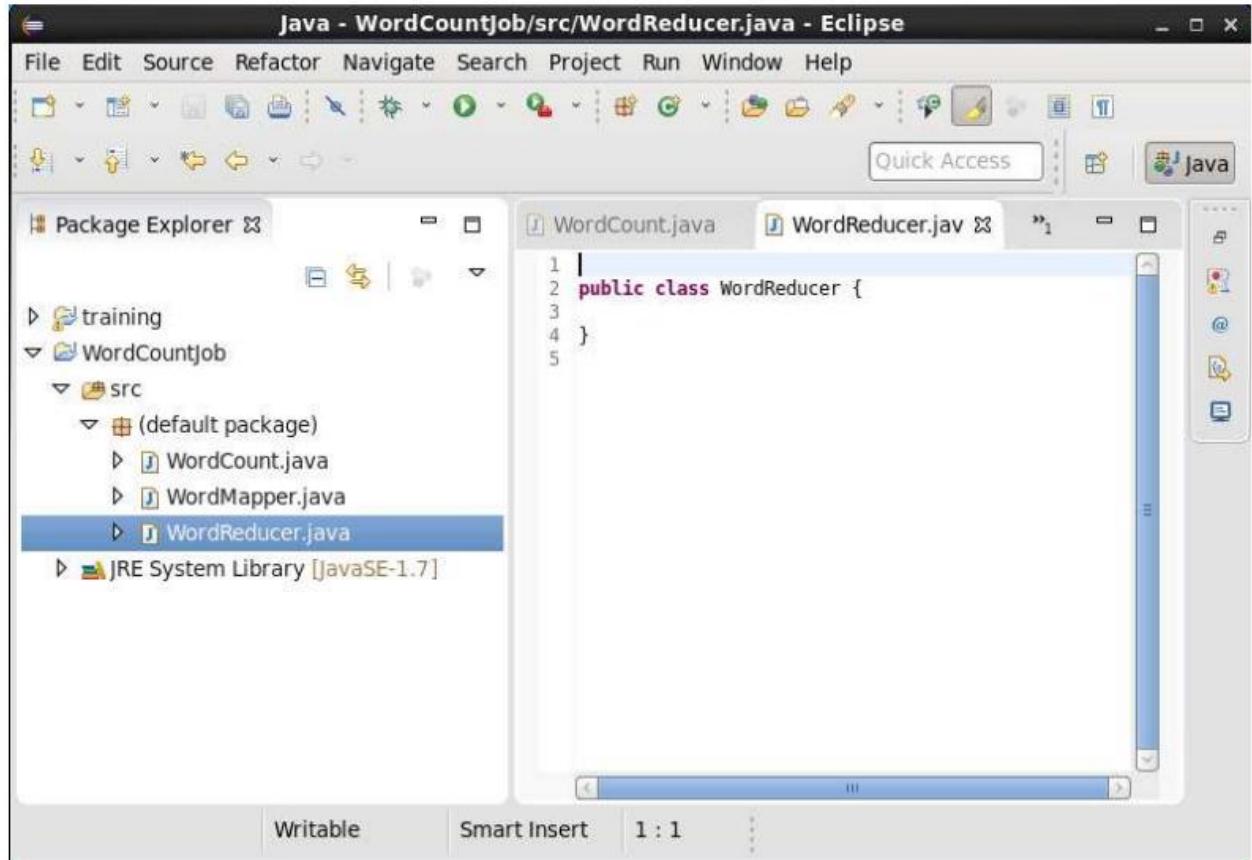
Name project :

WordCounJob click : finish



STEP 2:

Right Click on project -> New -> Class create class files as Name : WordCount
Name : WordMapper Name : WordReducer



Code for all the files is :

Mapper:

```
// Importing libraries
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.mapred.MapReduceBase;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reporter;

public class WCMapper extends MapReduceBase implements Mapper<LongWritable,
Text, IntWritable> {
```

Text,

```

// Map function

public void map(LongWritable key, Text value, OutputCollector<Text,
                IntWritable> output, Reporter rep) throws IOException

{

    String line = value.toString();

    // Splitting the line on spaces

    for (String word : line.split(" "))

    {

        if (word.length() > 0)

        {

            output.collect(new Text(word), new IntWritable(1));
        }
    }
}

```

reducer:

```

// Importing libraries

import java.io.IOException;

import java.util.Iterator;

import org.apache.hadoop.io.IntWritable;

import org.apache.hadoop.io.Text;

import org.apache.hadoop.mapred.MapReduceBase;

import org.apache.hadoop.mapred.OutputCollector;

import org.apache.hadoop.mapred.Reducer;

import org.apache.hadoop.mapred.Reporter;

public class WCReducer extends MapReduceBase implements Reducer<Text,
                                         IntWritable, Text, IntWritable> {

    // Reduce function

    public void reduce(Text key, Iterator<IntWritable> value,
                      OutputCollector<Text, IntWritable> output,
                      Reporter rep) throws IOException

    {

        int count = 0;

```

```

    // Counting the frequency of each words
    while (value.hasNext())
    {
        IntWritable i = value.next();
        count += i.get();
    }

    output.collect(key, new IntWritable(count));
}

}

```

Word count:

```

// Importing libraries

import java.io.IOException;
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobClient;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;

public class WCDriver extends Configured implements Tool {

    public int run(String args[]) throws IOException
    {
        if (args.length < 2)
        {
            System.out.println("Please give valid inputs");
            return -1;
        }

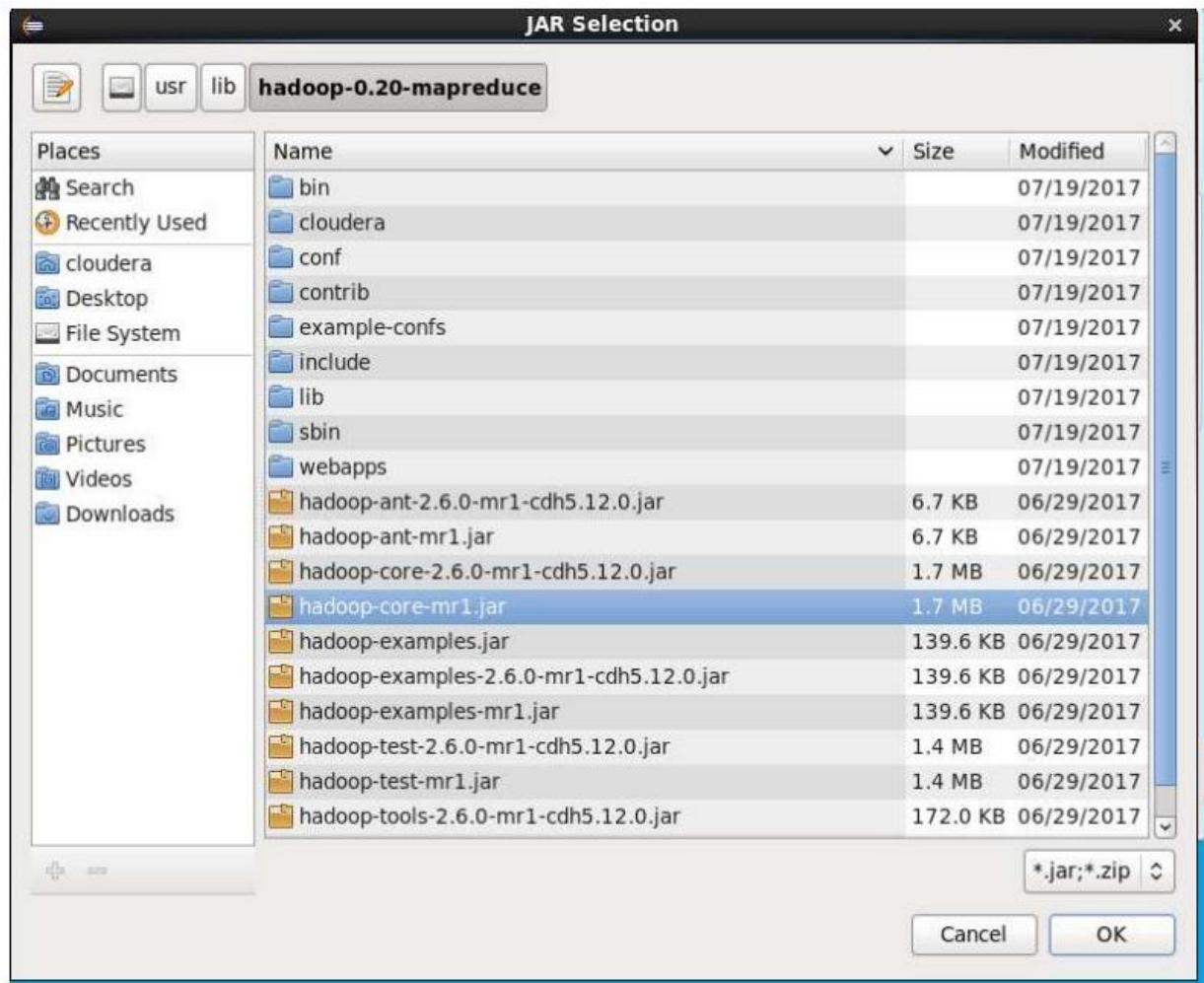
        JobConf conf = new JobConf(WCDriver.class);

```

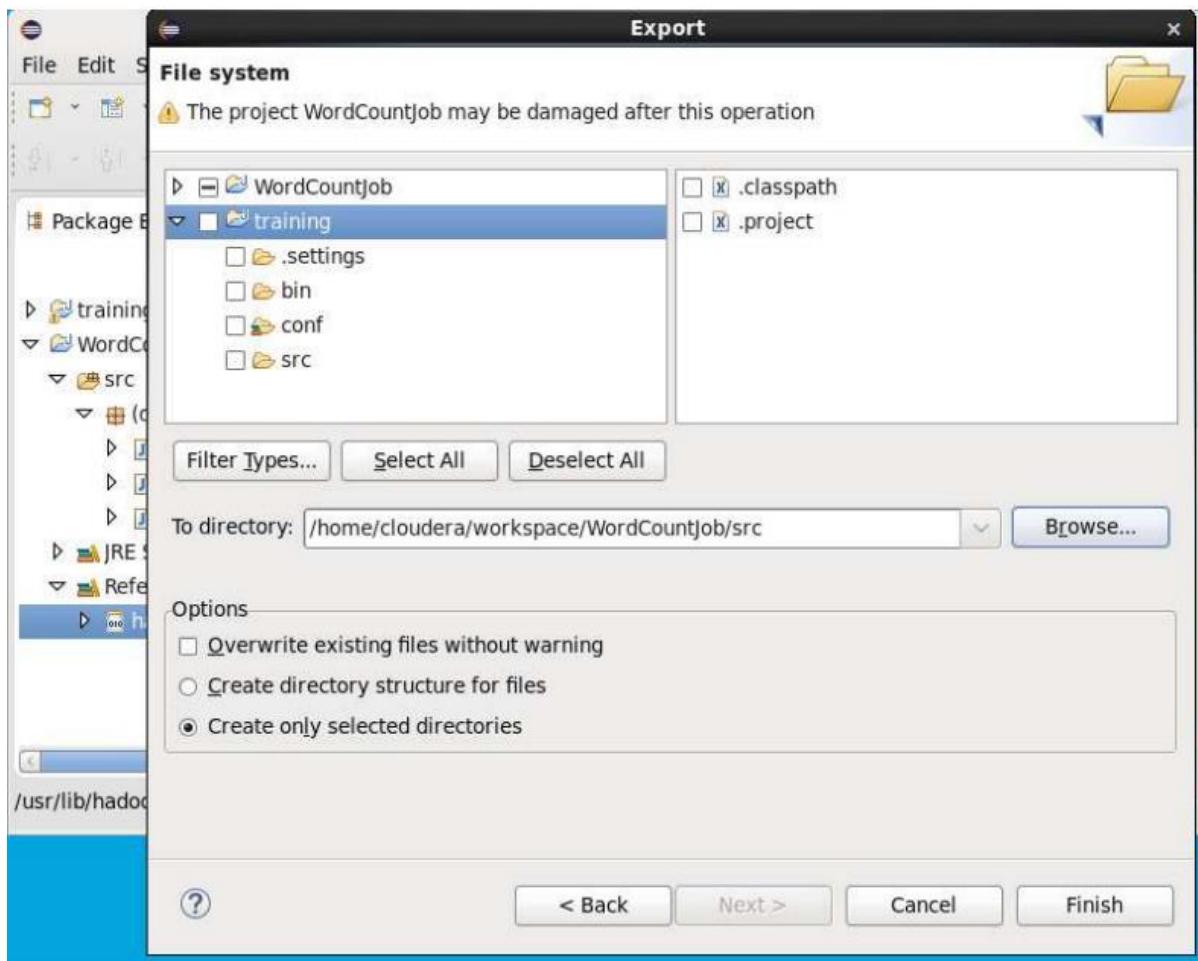
```
    FileInputFormat.setInputPaths(conf, new Path(args[0]));
    FileOutputFormat.setOutputPath(conf, new Path(args[1]));
    conf.setMapperClass(WCMapper.class);
    conf.setReducerClass(WCReducer.class);
    conf.setMapOutputKeyClass(Text.class);
    conf.setMapOutputValueClass(IntWritable.class);
    conf.setOutputKeyClass(Text.class);
    conf.setOutputValueClass(IntWritable.class);
    JobClient.runJob(conf);
    return 0;
}

// Main Method
public static void main(String args[]) throws Exception
{
    int exitCode = ToolRunner.run(new WCDriver(), args);
    System.out.println(exitCode);
}
```

STEP 3 : Right Click Project -> build path -> add external archievs -> filesystem
-> usr -> lib -> hadoop-0.20 -> hadoop-core.jar click add



STEP 4 : export the jar file in same folder
(training/workspace/WordCountJob/src)



STEP 5: open terminal [training@localhost ~]\$ ls -l

```
[cloudera@quickstart ~]$ ls -l
total 252
drwxrwxr-x 2 cloudera cloudera 4096 Aug 26 03:29 BDApracs
-rwxrwxr-x 1 cloudera cloudera 5387 Jul 19 2017 cloudera-manager
-rwxrwxr-x 1 cloudera cloudera 9964 Jul 19 2017 cm_api.py
drwxrwxr-x 2 cloudera cloudera 4096 Aug 28 01:24 data
-rw-rw-r-- 1 cloudera cloudera 67 Aug 27 00:31 demo.txt
drwxrwxr-x 3 cloudera cloudera 4096 Aug 26 03:44 Desktop
drwxrwxr-x 4 cloudera cloudera 4096 Jul 19 2017 Documents
drwxr-xr-x 2 cloudera cloudera 4096 Aug 26 03:21 Downloads
drwxrwsr-x 9 cloudera cloudera 4096 Aug 29 04:02 eclipse
-rw-rw-r-- 1 cloudera cloudera 53655 Jul 19 2017 enterprise-deployment.json
-rw-rw-r-- 1 cloudera cloudera 50515 Jul 19 2017 express-deployment.json
-rwxrwxr-x 1 cloudera cloudera 5007 Jul 19 2017 kerberos
drwxrwxr-x 2 cloudera cloudera 4096 Jul 19 2017 lib
drwxr-xr-x 2 cloudera cloudera 4096 Aug 26 03:21 Music
drwxrwxr-x 2 cloudera cloudera 4096 Aug 26 23:29 myfirstdata
-rwxrwxr-x 1 cloudera cloudera 4228 Jul 19 2017 parcels
drwxr-xr-x 2 cloudera cloudera 4096 Aug 26 03:21 Pictures
drwxr-xr-x 2 cloudera cloudera 4096 Aug 26 03:21 Public
-rw-rw-r-- 1 cloudera cloudera 18353 Aug 26 23:46 registercopy.java
-rw-rw-r-- 1 cloudera cloudera 18281 Aug 26 23:36 register.java
drwxr-xr-x 2 cloudera cloudera 4096 Aug 26 03:21 Templates
drwxrwxr-x 2 cloudera cloudera 4096 Aug 29 00:11 test
drwxrwxr-x 2 cloudera cloudera 4096 Aug 27 00:33 training
drwxr-xr-x 2 cloudera cloudera 4096 Aug 26 03:21 Videos
drwxrwxr-x 6 cloudera cloudera 4096 Aug 29 05:26 workspace
[cloudera@quickstart ~]$ █
```

STEP 6: create sample file [training@localhost ~]\$ cat WCFfile.txt

Hello we are doing BDA pracs Hello I'm a student Put the file in Hadoop

```
[cloudera@quickstart workspace]$ cat WCFfile.txt
Hello we are doing BDA pracs
Hello I'm a student
[cloudera@quickstart workspace]$ █
```

STEP 7 : Put the content of sample.txt to samplehadoop.txt

```
[training@localhost ~]$ hadoop fs -put WCFfile.txt WCFfile.txt  
[training@localhost ~]$ hadoop fs -ls
```

```
[cloudera@quickstart workspace]$ hadoop fs -put WCFfile.txt WCFfile.txt  
[cloudera@quickstart workspace]$ hadoop fs -ls  
Found 4 items  
-rw-r--r-- 1 cloudera cloudera 49 2022-08-29 05:32 WCFfile.txt  
drwxr-xr-x - cloudera cloudera 0 2022-08-26 03:41 hadoop  
drwxr-xr-x - cloudera cloudera 0 2022-08-26 04:08 raunak  
drwxr-xr-x - cloudera cloudera 0 2022-08-26 04:11 test.txt  
[cloudera@quickstart workspace]$ █
```

STEP 8 : change the directory to [training@localhost ~]\$ cd /home/training/workspace/WordCount/src

```
[cloudera@quickstart ~]$ cd workspace  
[cloudera@quickstart workspace]$ cd WordCount  
[cloudera@quickstart WordCount]$ cd src  
[cloudera@quickstart src]$ ls  
WCDriver.java WCMapper.java WCReducer.java  
[cloudera@quickstart src]$ █
```

STEP 9: Listing the contents of that directory [cloudera@quickstart src]\$ ls WCDriver.java WCMapper.java WCReducer.java

STEP 10 : Run the wordcount program and collect the output in WCOutput [cloudera@quickstart workspace]\$ hadoop jar WordCount.jar WCDriver WCFfile.txt WCOutput

```
[cloudera@quickstart workspace]$ hadoop jar WordCount.jar WCDriver WCFfile.txt WCOutput
22/08/29 05:37:49 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
22/08/29 05:37:50 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
22/08/29 05:37:50 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
22/08/29 05:37:51 INFO mapred.FileInputFormat: Total input paths to process : 1
22/08/29 05:37:51 INFO mapreduce.JobSubmitter: number of splits:2
22/08/29 05:37:51 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1661770860486_0001
22/08/29 05:37:52 INFO impl.YarnClientImpl: Submitted application application_1661770860486_0001
22/08/29 05:37:53 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1661770860486_0001/
22/08/29 05:37:53 INFO mapreduce.Job: Running job: job_1661770860486_0001
22/08/29 05:38:08 INFO mapreduce.Job: Job job_1661770860486_0001 running in uber mode : false
22/08/29 05:38:08 INFO mapreduce.Job: map 0% reduce 0%
22/08/29 05:38:21 INFO mapreduce.Job: map 50% reduce 0%
22/08/29 05:38:22 INFO mapreduce.Job: map 100% reduce 0%
22/08/29 05:38:28 INFO mapreduce.Job: map 100% reduce 100%
22/08/29 05:38:29 INFO mapreduce.Job: Job job_1661770860486_0001 completed successfully
22/08/29 05:38:30 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=115
    FILE: Number of bytes written=375712
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=290
```

STEP 11 : chk the output file [cloudera@quickstart workspace]\$
hadoop fs -ls WCOutput

STEP 12: Displaying the Output [cloudera@quickstart workspace]\$
hadoop fs -cat WCOutput/part-00000

```
[cloudera@quickstart workspace]$ hadoop fs -cat WCOutput/part-00000
BDA      1
Hello    2
I'm     1
a       1
are     1
doing   1
pracs   1
student 1
we      1
[cloudera@quickstart workspace]t ■
```

EXPERIMENT NO: 6

AIM: Create HIVE database and descriptive analytics (basic statistics)

Theory:

Hive is a powerful data warehousing and SQL-like query language tool in the Hadoop ecosystem. Developed by Facebook and later open-sourced, Hive is designed to make it easier for users to query and analyze large datasets stored in Hadoop's distributed file system (HDFS) without requiring extensive programming skills. In this 400-word explanation, we'll explore the key features and use cases of Hive.

Key Features and Concepts:

1. **SQL-like Query Language:** Hive Query Language (HQL) resembles SQL, making it familiar to analysts and data scientists. Users can write queries to extract, transform, and analyze data in a manner similar to traditional relational databases.
2. **Schema-on-Read:** Unlike traditional databases with a strict schema-on-write, Hive follows a schema-on-read approach. This means that data is read from HDFS as is, and you can apply the schema when querying the data. This flexibility is particularly useful for handling semi-structured and unstructured data.
3. **Data Storage:** Hive organizes data into tables, databases, and partitions. It can store data in various formats, such as Text, Parquet, ORC, and more. This flexibility allows you to choose the best storage format for your data.
4. **Extensibility:** Hive supports custom user-defined functions (UDFs) and user-defined aggregates (UDAs) in multiple programming languages, such as Java, Python, and more, allowing you to perform complex data processing within your queries.
5. **Optimization:** Hive includes a query optimizer that can optimize query execution plans, improving query performance.
6. **Integration with Hadoop Ecosystem:** Hive seamlessly integrates with other Hadoop ecosystem components like HBase, Pig, and Spark, making it a valuable tool for a wide range of data processing tasks.

Hive Query Language (HQL):

HQL is a SQL-like language used for querying and manipulating data in Hive. It includes familiar SQL commands like SELECT, JOIN, GROUP BY, and ORDER BY. However, it also incorporates some Hive-specific extensions to work with distributed data and the schema-on-read approach.

Step 1: starting with HIVE, viewing database and tables

```
[training@localhost ~]$ hive
Hive history file=/tmp/training/hive_job_log_training_202309132302_441038358.tx
hive> show databases;
OK
books
data1
default
employee
example
hivecl276
mobilo
test
Time taken: 1.025 seconds
hive> show tables;
OK
message
messagel
posts
Time taken: 0.131 seconds
hive> use bank;
FAILED: Error in metadata: ERROR: The database bank does not exist.
FAILED: Execution Error, return code 1 from org.apache.hadoop.hive.ql.exec.DDLT
er: muted late database bank.
```

Step 2: creating database and bank

```
hive> create database bank;
OK
Time taken: 0.11 seconds
hive> show databases;
OK
bank
books
data1
default
employee
example
hivecl276
mobilo
test
Time taken: 0.058 seconds
hive> use bank;
OK
Time taken: 0.024 seconds
hive> create table emp(id INT,name STRING ,sal DOUBLE)row format delimited
      fields terminated by ',' stored as textfile;
FAILED: Parse Error: line 1:69 mismatched input 'fis' expecting EOF near 'delim

hive> create table emp(id INT,name STRING ,sal DOUBLE)row format delimited
      fields terminated by ',' stored as textfile;
OK
```

```

Time taken: 0.058 seconds
hive> use bank;
OK
Time taken: 0.024 seconds
hive> create table emp(id INT,name STRING ,sal DOUBLE)row format delimited field
s terminated by ',' stored as textfile;
FAILED: Parse Error: line 1:69 mismatched input 'fis' expecting EOF near 'delimited'

hive> create table emp(id INT,name STRING ,sal DOUBLE)row format delimited fields terminated by ',' stored as
OK
Time taken: 0.081 seconds
hive> show tables;
hive and run installed applications
emp
Time taken: 0.096 seconds
hive> describe emp;
OK
id      int
name    string
sal     double
Time taken: 0.07 seconds
hive> load data local inpath '/home/training/demo.txt' into table emp;
Copying data from file:/home/training/demo.txt
Copying file: file:/home/training/demo.txt
Loading data to table bank.emp
OK
Time taken: 0.155 seconds

```

Step 3: performing queries on the table

Various statistical queries like AVG () , MAX() atc

```

4
Time taken: 11.344 seconds
hive> select AVG(sal) as avg_salary from emp_sal;
Total MapReduce jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapred.reduce.tasks=<number>
Starting Job = job_202309132232_0005, Tracking URL = http://localhost:50030/jobdetails.jsp?jobid=job_2023091
Kill Command = /usr/lib/hadoop/bin/hadoop job -Dmapred.job.tracker=localhost:8021 -kill job_202309132232_00
2023-09-13 23:27:41,550 Stage-1 map = 0%,  reduce = 0%
2023-09-13 23:27:42,554 Stage-1 map = 100%,  reduce = 0%
2023-09-13 23:27:49,607 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_202309132232_0005
OK
4375.0
Time taken: 10.296 seconds
hive> select MAX(sal) as max_salary from emp_sal;
Total MapReduce jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
click here to hide all windows and show the desktop. set=<number>
```

```

Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_202309132232_0003, Tracking URL = http://localhost:50030/jobdetails.jsp?jobid=job_2023091
Kill Command = /usr/lib/hadoop/bin/hadoop job -Dmapred.job.tracker=localhost:8021 -kill job_202309132232_00
2023-09-13 23:22:24,723 Stage-1 map = 0%,  reduce = 0%
2023-09-13 23:22:25,749 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_202309132232_0003
OK
1      ABC      4000.0
1      abc      4500.0
Time taken: 3.569 seconds
hive> select count(*) from emp_sal;
Total MapReduce jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapred.reduce.tasks=<number>
Starting Job = job_202309132232_0004, Tracking URL = http://localhost:50030/jobdetails.jsp?jobid=job_2023091
Kill Command = /usr/lib/hadoop/bin/hadoop job -Dmapred.job.tracker=localhost:8021 -kill job_202309132232_00
2023-09-13 23:23:08,691 Stage-1 map = 0%,  reduce = 0%
2023-09-13 23:23:10,700 Stage-1 map = 100%,  reduce = 0%
2023-09-13 23:23:17,756 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_202309132232_0004

```

Wednesday September 13

```

emp
Time taken: 0.096 seconds
hive> describe emp;
OK
id      int
name    string
sal     double
Time taken: 0.07 seconds
hive> load data local inpath '/home/training/demo.txt' into table emp;
Copying data from file:/home/training/demo.txt
Copying file: file:/home/training/demo.txt
Loading data to table bank.emp
OK
Time taken: 0.155 seconds
hive> select * from emp;
OK
1      ABC      4000.0
2      PQR      4500.0
1      abc      4500.0
2      pqr      4500.0
Time taken: 0.122 seconds
hive> alter table emp rename to emp_sal;
OK
Time taken: 0.118 seconds
hive> select * from emp_sal where id=1;
Total MapReduce jobs = 1

```

Click here to hide all windows and show the desktop.

Step 4:

EXIT from hive

```
OK
4375.0
Time taken: 10.296 seconds
hive> select MAX(sal) as max_salary from emp_sal;
Total MapReduce jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapred.reduce.tasks=<number>
Starting Job = job_202309132232_0006, Tracking URL = http://localhost:50030/jobdetails.jsp?jobid=job_2023091
Kill Command = /usr/lib/hadoop/bin/hadoop job -Dmapred.job.tracker=localhost:8021 -kill job_202309132232_00
2023-09-13 23:28:17,824 Stage-1 map = 0%,  reduce = 0%
2023-09-13 23:28:18,827 Stage-1 map = 100%,  reduce = 0%
2023-09-13 23:28:26,011 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_202309132232_0006
OK
4500.0
Time taken: 9.427 seconds
hive> drop table emp_sal;
OK
Time taken: 0.129 seconds
hive> exit;
```

EXPERIMENT NO: 7

AIM: implement Bloom Filter using python / R

Theory:

A Bloom filter is a space-efficient probabilistic data structure that is used to test whether an element is a member of a set. For example, checking availability of username is set membership problem, where the set is the list of all registered username. A Bloom filter is a probabilistic data structure that plays a significant role in big data analytics by providing an efficient way to test whether a given element is a member of a set. It is a memory-efficient and scalable tool commonly used in various applications within the realm of distributed and big data systems. The price we pay for efficiency is that it is probabilistic in nature that means, there might be some False Positive results. False positive means, it might tell that given username is already taken but actually it's not.

Interesting Properties of Bloom Filters

- Unlike a standard hash table, a Bloom filter of a fixed size can represent a set with an arbitrarily large number of elements.
- Adding an element never fails. However, the false positive rate increases steadily as elements are added until all bits in the filter are set to 1, at which point all queries yield a positive result.
- Bloom filters never generate **false negative** result, i.e., telling you that a username doesn't exist when it actually exists.
- Deleting elements from filter is not possible because, if we delete a single element by clearing bits at indices generated by k hash functions, it might cause deletion of few other elements. Example – if we delete “geeks” (in given example below) by clearing bit at 1, 4 and 7, we might end up deleting “nerd” also Because bit at index 4 becomes 0 and bloom filter claims that “nerd” is not present.

CODE:

```
import java.util.Scanner;

public class StreamData {
    public static void main(String[] args) {
        Scanner scanner = new Scanner(System.in);

        System.out.print("Enter the length of the stream data (m): ");
        int m = scanner.nextInt();
        int[] dataStream = new int[m];

        System.out.print("Enter the number of inputs to train the array (n): ");
        int n = scanner.nextInt();

        // Training the array
        for (int i = 0; i < n; i++) {
            System.out.print("Enter a number to train the array: ");
            int input = scanner.nextInt();
            int hash1 = (input % 5) % m;
            int hash2 = ((2 * input + 3) % 5) % m;

            if (dataStream[hash1] != 1) {
                dataStream[hash1] = 1;
            }
            if (dataStream[hash2] != 1) {
                dataStream[hash2] = 1;
            }
        }

        // Print the training array
        System.out.print("Training Array: ");
        for (int i = 0; i < m; i++) {
            System.out.print(dataStream[i] + " ");
        }
        System.out.println();

        System.out.print("Enter the number of times to test (x): ");
        int x = scanner.nextInt();

        // Testing the array
        for (int i = 0; i < x; i++) {
            System.out.print("Enter a number to test: ");
            int testValue = scanner.nextInt();
            int hash1 = (testValue % 5) % m;
            int hash2 = ((2 * testValue + 3) % 5) % m;
```

```
        if (dataStream[hash1] == 0 && dataStream[hash2] == 0) {
            System.out.println("Number is not in the stream.");
        } else if (dataStream[hash1] == 1 && dataStream[hash2] == 1) {
            System.out.println("Number is in the stream.");
        } else {
            System.out.println("Value not present in the stream.");
        }
    }

    scanner.close();
}
}
```

OUTPUT:

```
"C:\Program Files\Java\jdk-20\bin\java.exe" "-javaagent:C:\Prog
Enter the length of the stream data (m): 5
Enter the number of inputs to train the array (n): 2
Enter a number to train the array: 9
Enter a number to train the array: 11
Training Array: 1 1 0 0 1
Enter the number of times to test (x): 2
Enter a number to test: 16
Number is in the stream.
Enter a number to test: 15
Value not present in the stream.

Process finished with exit code 0
```

EXPERIMENT NO: 8

AIM: implement FM algorithm using python / R

Theory:

The Flajolet-Martin (FM) algorithm is a probabilistic algorithm used to estimate the number of distinct elements in a stream of data without storing or revisiting each data item. This algorithm is particularly useful in big data analytics and data stream processing, where it's often infeasible to store and process every data point individually.

Key Features and Concepts:

- Probabilistic Counting:** The FM algorithm uses hash functions and probability theory to provide an approximate count of distinct elements. It leverages the principle that different hash functions produce distinct values with a certain probability.
- Bit Patterns:** Instead of counting distinct elements directly, the FM algorithm focuses on the leading zeros in the binary representation of the hash values of the data elements. These patterns are used to estimate the number of distinct elements.
- Randomization:** The algorithm introduces randomization by using multiple hash functions. Each hash function is applied to each data element independently, and the maximum number of leading zeros across all hash values is tracked.
- Estimation Process:** The FM algorithm calculates the estimate of the number of distinct elements using the formula $2^{(\text{maximum leading zeros})}$. This estimation provides a probabilistic count, which may have some error but is highly memory-efficient.

Use Cases:

- Large Data Streams:** FM is particularly useful in scenarios where the data stream is too large to store and process in memory. It provides a memory-efficient method for estimating distinct element counts.
- Web Traffic Analysis:** In web analytics, the FM algorithm can estimate the number of distinct visitors or IP addresses accessing a website in real-time.

CODE:

```
import java.util.Scanner;

public class Main {

    static int hash(int x) {
        return (6 * x + 1) % 5;
    }

    static String toThreeBitBinary(int num) {
        String binary = Integer.toBinaryString(num);
        while (binary.length() < 3) {
            binary = "0" + binary;
        }
        return binary;
    }

    static void countTrailingZeros(String[] arr, int[] result) {
        for (int i = 0; i < arr.length; i++) {
            String binary = arr[i];
            int count = 0;
            boolean encounteredOne = false;
            for (int j = binary.length() - 1; j >= 0; j--) {
                if (binary.charAt(j) == '0' && !encounteredOne) {
                    count++;
                } else if (binary.charAt(j) == '1') {
                    encounteredOne = true;
                }
            }
            if (!encounteredOne) {
                count = 0;
            }
            result[i] = count;
        }
    }

    public static void main(String[] args) {
        Scanner scanner = new Scanner(System.in);

        System.out.print("Enter the size of the array: ");
        int size = scanner.nextInt();
        int[] inputArray = new int[size];

        System.out.println("Enter elements of the array:");
        for (int i = 0; i < size; i++) {
            System.out.print("Element " + (i + 1) + ": ");
            inputArray[i] = scanner.nextInt();
        }

        int[] hashedArray = new int[size];
        for (int i = 0; i < size; i++) {
            hashedArray[i] = hash(inputArray[i]);
        }
    }
}
```

```

String[] binaryArray = new String[size];
for (int i = 0; i < size; i++) {
    binaryArray[i] = toThreeBitBinary(hashedArray[i]);
}
hint[] trailingZerosArray = new int[size];
countTrailingZeros(binaryArray, trailingZerosArray);

int maxTrailingZeros = Integer.MIN_VALUE;
for (int count : trailingZerosArray) {
    if (count > maxTrailingZeros) {
        maxTrailingZeros = count;
    }
}

System.out.println("Maximum trailing zeros: " + maxTrailingZeros);
double powerOfTwo = Math.pow(2, maxTrailingZeros);
System.out.println("the number of distinct elements " +
powerOfTwo);

scanner.close();
}
}

```

OUTPUT:

```

"C:\Program Files\Java\jdk-20\bin\java.exe" "-ja
Enter the size of the array: 12
Enter elements of the array:
Element 1: 1
Element 2: 3
Element 3: 2
Element 4: 1
Element 5: 2
Element 6: 3
Element 7: 4
Element 8: 3
Element 9: 1
Element 10: 2
Element 11: 3
Element 12: 1
Maximum trailing zeros: 2
the number of distinct elements 4.0

```

EXPERIMENT NO: 9

AIM: Data visualization using R

Theory:

Data visualization is a crucial aspect of data analysis and interpretation. It involves representing data graphically to identify trends, patterns, and insights that might be difficult to discern from raw data alone. R is a powerful and versatile tool for data visualization, with numerous packages and libraries that make it easy to create a wide range of visualizations.

Key Concepts in Data Visualization:

1. **Data Visualization Types:** Data can be visualized in various forms, including charts, graphs, plots, maps, and more. The choice of visualization type depends on the nature of the data and the specific insights you aim to convey.
2. **Visualization Libraries in R:** R provides several packages for data visualization, with **ggplot2**, **base**, and **lattice** being some of the most commonly used. These libraries offer different approaches to creating visualizations.
3. **Aesthetics and Mapping:** In R, aesthetics refer to the visual properties of the plot, such as color, size, and shape. Mapping aesthetics to variables in the dataset allows you to create dynamic and informative visualizations.
4. **Grammar of Graphics:** **ggplot2** follows the "Grammar of Graphics" philosophy, making it easy to build up a plot by adding layers, such as data, geometries (e.g., points, lines, bars), scales, and labels, in a structured and coherent way.

Creating Data Visualizations in R:

1. **Loading Data:** Start by importing your data into R, typically in the form of a data frame. You can load data from various sources, including CSV files, databases, or web APIs.
2. **Selecting Visualization Library:** Choose the R visualization library that best suits your needs. **ggplot2** is popular for its flexibility and the ability to create complex and customized visualizations.

3. **Data Preparation:** Before creating visualizations, preprocess the data as needed. This may include filtering, aggregating, and transforming the data to make it suitable for visualization.
4. **Creating Basic Plots:** Use basic functions like `plot()` or `hist()` from the base graphics package to create simple visualizations. For example, `plot(x, y)` generates a scatterplot, and `hist(data)` creates a histogram.
5. **Advanced Visualizations with ggplot2:** For more advanced and customized visualizations, `ggplot2` is a powerful choice. Start by creating a `ggplot` object and adding layers with the `+` operator. For example, to create a scatterplot, you can use `ggplot(data, aes(x, y)) + geom_point()`.
6. **Customization:** Customize your visualization by modifying aesthetics, adjusting labels, adding titles, and changing the theme. `ggplot2` allows for extensive customization to make your plots visually appealing and informative.
7. **Exporting Plots:** Save your visualizations as image files (e.g., PNG, JPEG) or as vector graphics (e.g., PDF, SVG) using functions like `ggsave()` or `pdf()`. You can also display them directly in your R environment.

Types of Visualizations in R:

R supports a wide variety of visualizations, including:

- **Scatter Plots:** Use `geom_point()` in `ggplot2` to create scatter plots, which are useful for displaying the relationship between two numeric variables.
- **Bar Charts and Histograms:** Represent categorical or discrete data with bar charts and continuous data distributions with histograms using `geom_bar()` and `geom_histogram()`.
- **Line Charts:** Show trends and patterns over time or a continuous axis with line charts using `geom_line()`.
- **Box Plots:** Display the distribution, median, and outliers of a dataset using `geom_boxplot()`.
- **Heatmaps:** Visualize patterns in large matrices or tables with colors indicating values using the `heatmap()` function or packages like `pheatmap`.

```
In [24]: # Read the CSV file with semicolon as the delimiter
data <- read.csv("e-shop clothing 2008.csv", sep = ";", header = TRUE)

head(data)
```

A data.frame: 6 × 14

	year	month	day	order	country	session.ID	page.1..main.category.	page.2..clothing.model.	...
	<int>	<int>	<int>	<int>	<int>	<int>	<int>	<int>	<chr>
1	2008	4	1	1	29	1	1	1	A13
2	2008	4	1	2	29	1	1	1	A16
3	2008	4	1	3	29	1	2	2	B4
4	2008	4	1	4	29	1	2	2	B17
5	2008	4	1	5	29	1	2	2	B8
6	2008	4	1	6	29	1	3	3	C56

```
In [32]: # Load the dplyr package
library(dplyr)

# Create a new DataFrame with relevant columns
csdf <- data %>%
  select(month, day, `page.1..main.category.`, price)

# Rename the columns
csdf <- csdf %>%
  rename(Month = month, Day = day, Type = `page.1..main.category.`, Price = price)

# Replace values in the "Type" column
csdf$Type <- recode(csdf$Type, "1" = "Trousers", "2" = "Skirts", "3" = "Blouses", "4" = "Jackets")

# Replace values in the "Month" column
csdf$Month <- recode(csdf$Month, "4" = "April", "5" = "May", "6" = "June", "7" = "July")

# Print the first few rows of the new DataFrame
head(csdf)
```

A data.frame: 6 × 4

	Month	Day	Type	Price
	<chr>	<int>	<chr>	<int>
1	April	1	Trousers	28
2	April	1	Trousers	33
3	April	1	Skirts	52
4	April	1	Skirts	38
5	April	1	Skirts	52
6	April	1	Blouses	57

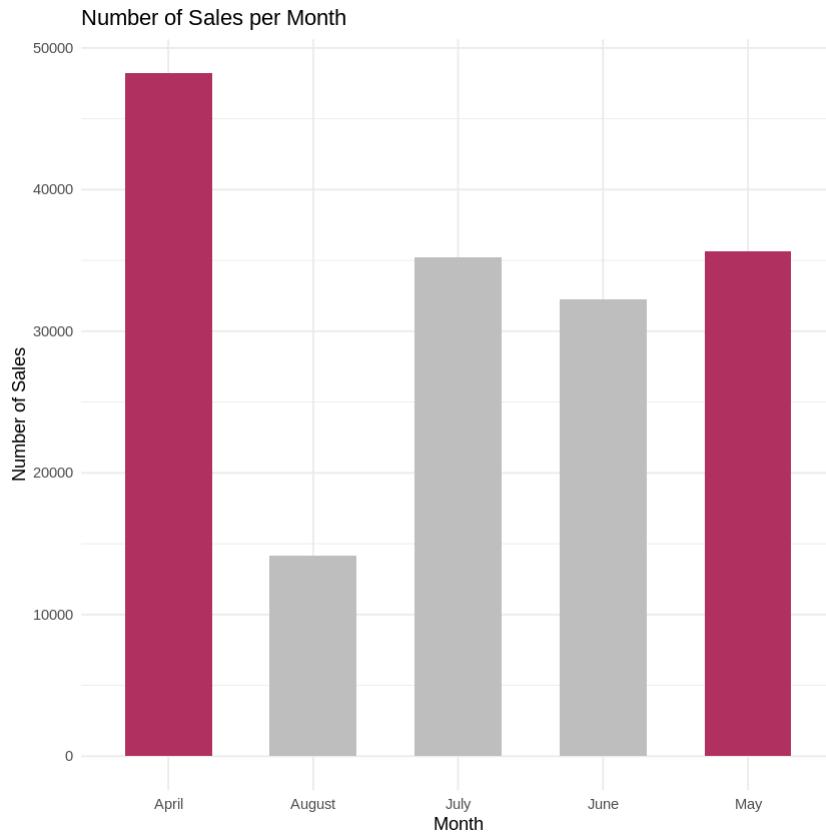
```
In [33]: # Assuming 'csdf' is your data frame
# Load the ggplot2 package
library(ggplot2)

# Count the number of goods sold each month
csmsm <- table(csdf$Month)

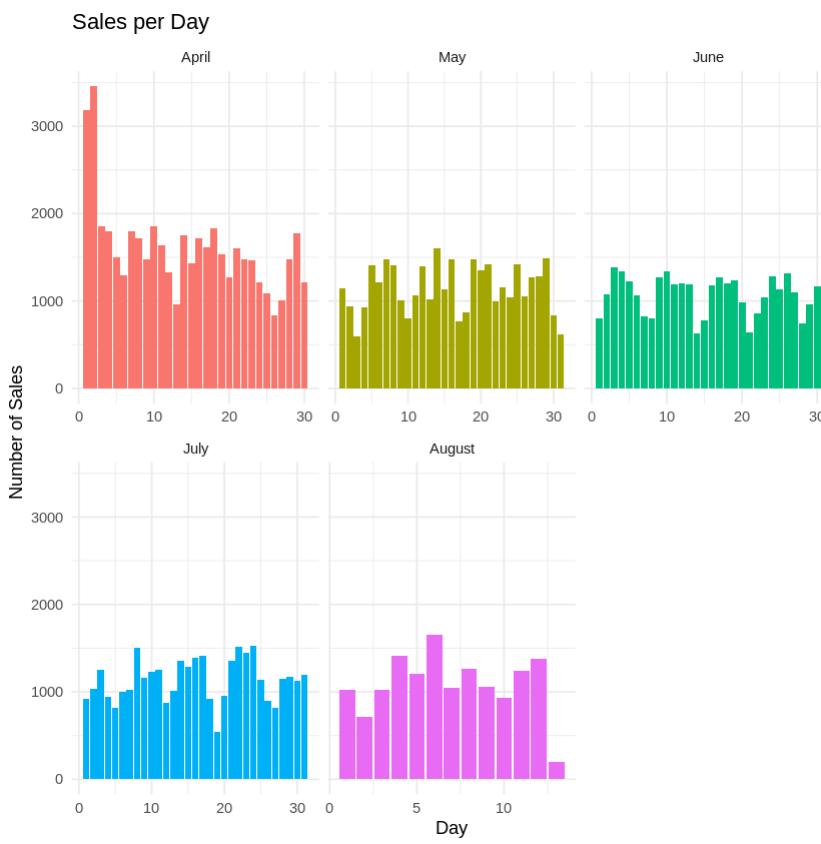
# Create a data frame for plotting
data <- data.frame(Month = names(csmsm), Sales = as.vector(csmsm))

# Create the bar plot
ggplot(data, aes(x = Month, y = Sales, fill = Month)) +
  geom_bar(stat = "identity", width = 0.6) +
  labs(title = "Number of Sales per Month",
       y = "Number of Sales") +
  scale_fill_manual(values = c("maroon", "gray", "gray", "gray", "maroon")) +
```

```
theme_minimal() +  
theme(legend.position = "none")
```



```
In [36]: # Assuming 'csdf' is your data frame  
# Load the ggplot2 package  
library(ggplot2)  
  
# Create data frames for each month  
csap <- subset(csdf, Month == "April")  
csmay <- subset(csdf, Month == "May")  
csjn <- subset(csdf, Month == "June")  
csjl <- subset(csdf, Month == "July")  
csau <- subset(csdf, Month == "August")  
  
# Combine all data frames into one  
combined_df <- rbind(csap, csmay, csjn, csjl, csau)  
combined_df$Month <- factor(combined_df$Month, levels = c("April", "May", "June", "July", "August"))  
  
# Create subplots with facets  
p <- ggplot(combined_df, aes(x = Day)) +  
  geom_bar(aes(fill = Month), position = "dodge") +  
  labs(title = "Sales per Day", y = "Number of Sales") +  
  theme_minimal() +  
  theme(legend.position = "none") +  
  facet_wrap(~Month, scales = "free_x")  
  
print(p)
```

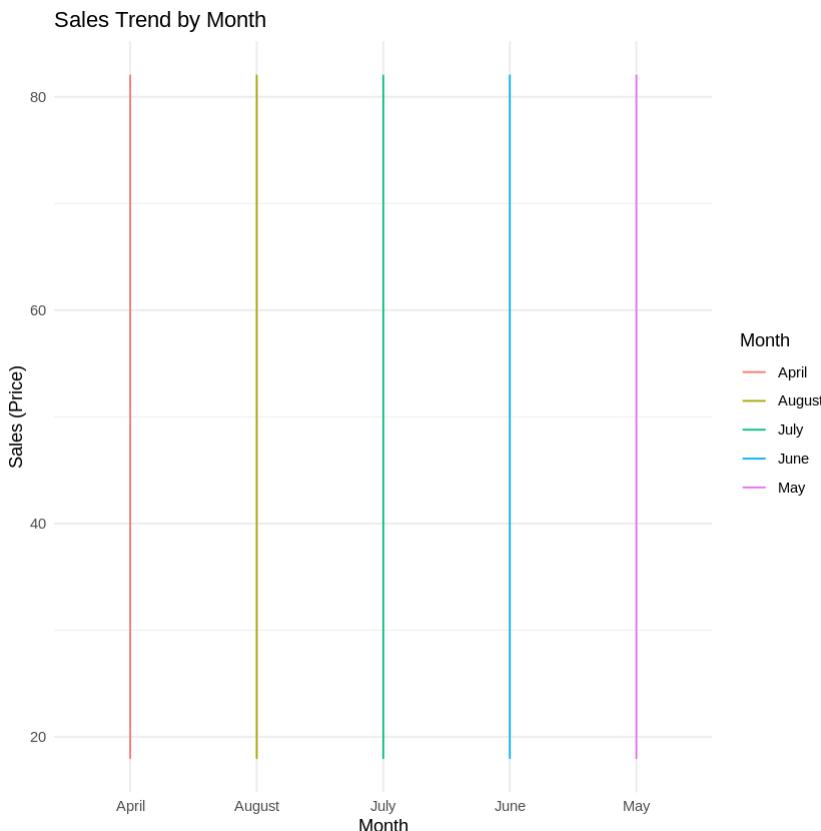


```
In [72]: # Load the ggplot2 package
library(ggplot2)

# Assuming you have a data frame named "csdf" with columns "Month" and "Price"

# Create a line plot
line_plot <- ggplot(csdf, aes(x = Month, y = Price, color = Month, group = Month)) +
  geom_line() +
  labs(title = "Sales Trend by Month", x = "Month", y = "Sales (Price)") +
  theme_minimal()

# Display the line plot
print(line_plot)
```



```
In [67]: # Load the ggplot2 package
library(ggplot2)

# Assuming you have a data frame named "csdf" with columns "Type" and "Price"

# Aggregate total sales for each clothing type
total_sales <- aggregate(Price ~ Type, data = csdf, FUN = sum)

# Create a pie chart
pie_chart <- ggplot(total_sales, aes(x = "", y = Price, fill = Type)) +
```

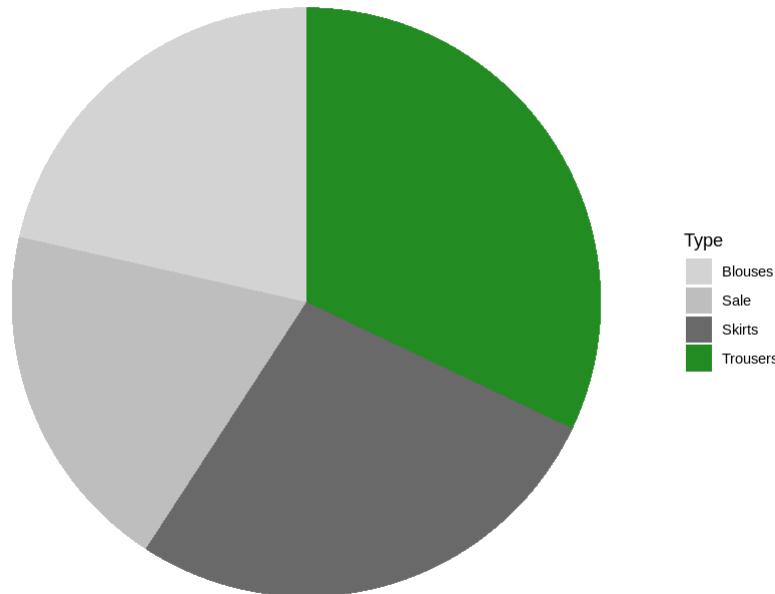
```

geom_bar(stat = "identity", width = 1) +
coord_polar(theta = "y") +
labs(title = "Share of Total Sales in Dollars per Type") +
scale_fill_manual(values = c('lightgray', 'gray', 'dimgray', 'forestgreen')) +
theme_void()

# Display the pie chart
print(pie_chart)

```

Share of Total Sales in Dollars per Type



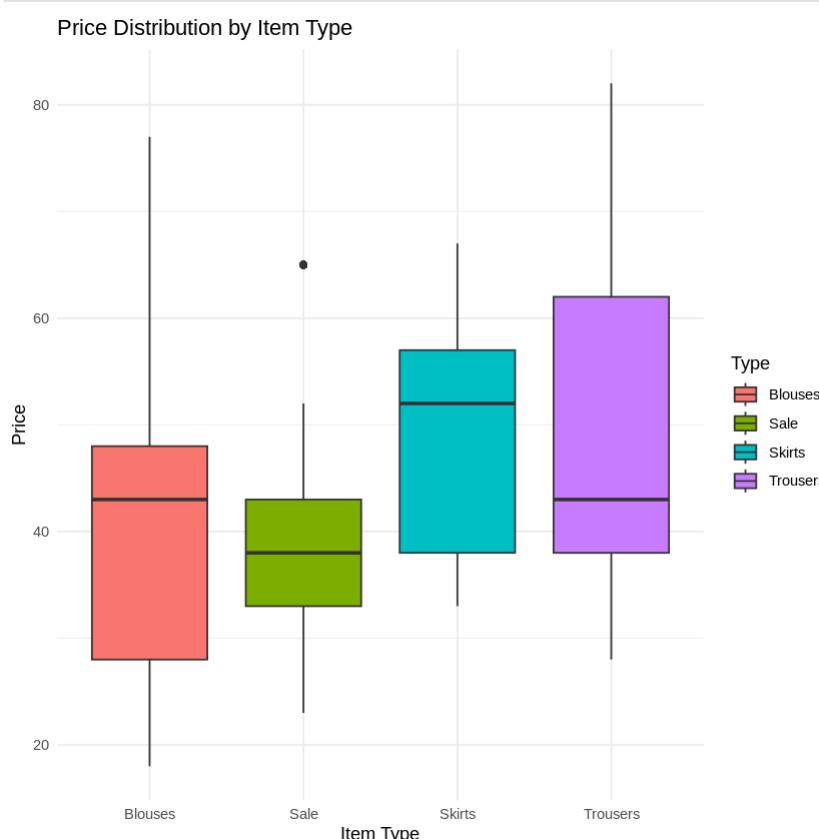
```

In [73]: # Load the ggplot2 package
library(ggplot2)

# Create a box plot
box_plot <- ggplot(csdf, aes(x = Type, y = Price, fill = Type)) +
  geom_boxplot() +
  labs(title = "Price Distribution by Item Type", x = "Item Type", y = "Price") +
  theme_minimal()

# Display the box plot
print(box_plot)

```



```

In [80]: # Load the ggplot2 package
library(ggplot2)

```

```

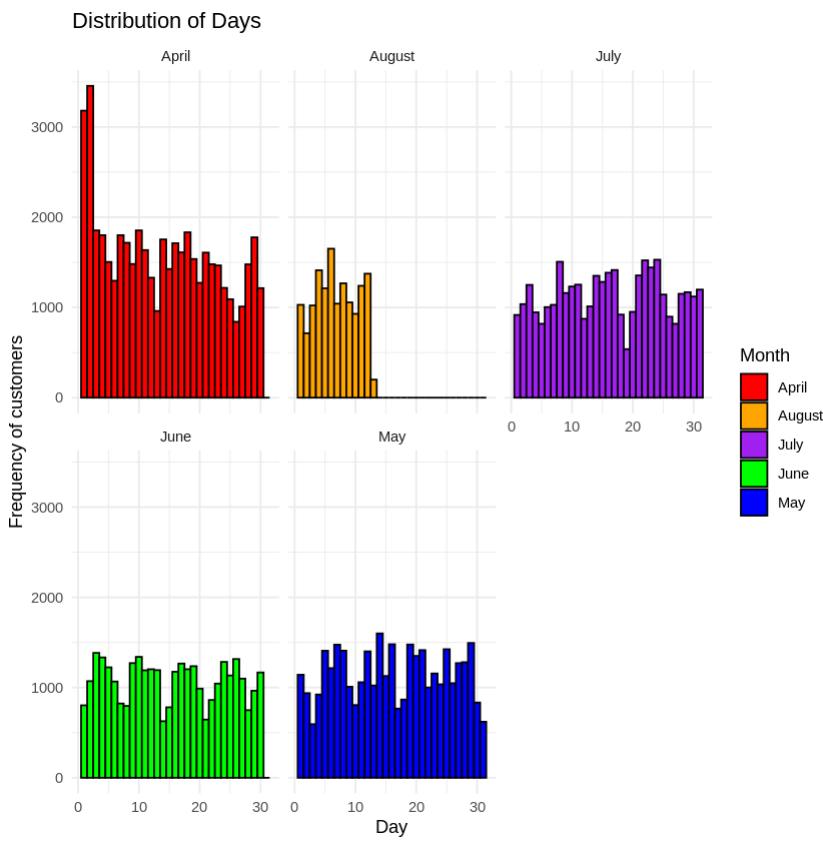
# Assuming you have a data frame named "csdf" with columns "Month" and "Day"

# Define colors for each month
month_colors <- c("April" = "red", "May" = "blue", "June" = "green", "July" = "purple")

# Create separate histograms for each month
day_histogram <- ggplot(csdf, aes(x = Day, fill = Month)) +
  geom_histogram(binwidth = 1, color = "black") +
  labs(title = "Distribution of Days", x = "Day", y = "Frequency of customers") +
  theme_minimal() +
  scale_fill_manual(values = month_colors) +
  facet_wrap(~Month, nrow = 2)

# Display the histograms
print(day_histogram)

```



Report on Mini Project

Subject: Big Data Analytics (CSC702)

AY: 2024-25

TAXI DATA ANALYSIS

Karan Adwani : 2103009

Priyanka Bajaj : 21030012

Hardik Bhagat : 21030017

Akshay Jagiasi : 2103062

Guided By

(Anagha Durugkar)

CHAPTER 1: INTRODUCTION

With the rise of ride-hailing services and urban transportation networks, analyzing taxi data has become crucial for improving transportation systems, optimizing routes, predicting demand, and enhancing customer service. Machine learning algorithms can be leveraged to extract valuable insights from taxi data, such as predicting trip durations, optimizing routes, and detecting anomalies.

The objective of this project is to analyze and compare the performance of various machine learning algorithms on a dataset of taxi trips. The algorithms are evaluated based on their performance in predicting taxi trip duration and route optimization, among other metrics like accuracy, precision, recall, F1-score, and computational efficiency.

The project focuses on:

- Data Preprocessing of taxi data
- Selection and implementation of machine learning algorithms
- Training models on taxi trip data
- Evaluating and comparing algorithm performance
- Analyzing results to provide insights into taxi operations

CHAPTER 2: DATA DESCRIPTION AND ANALYSIS

The dataset utilized in this project is composed of various CSV files containing detailed information about taxi trips. Each dataset contributes unique attributes crucial for understanding the dynamics of taxi transportation within a city. The key datasets are:

Taxi Trips Data:

This dataset contains comprehensive details of individual taxi trips, providing insights into trip durations, distances, locations, and other relevant factors. It serves as the core data for analyzing taxi demand, routes, and patterns.

Key attributes include:

- Trip ID: A unique identifier for each taxi trip.
- Pickup Date/Time: The date and time when the trip started.
- Drop-off Date/Time: The date and time when the trip ended.
- Pickup Location: Latitude and longitude of the pickup point.
- Drop-off Location: Latitude and longitude of the drop-off point.
- Passenger Count: The number of passengers in the taxi during the trip.
- Distance Traveled: The total distance covered during the trip, typically in kilometers or miles.
- Fare Amount: The total fare paid for the trip, including tolls, surcharges, and tips.
- Payment Type: The method of payment (cash, card, etc.).
- Trip Duration: The length of the trip, usually in seconds or minutes.
- Rate Code: A code indicating the rate charged for the trip, such as standard,

premium, or discounted.

- Weather Conditions (Optional): Data on weather conditions during the trip (e.g., clear, rainy, snowy), which can impact trip duration.

Data preprocessing

To prepare the taxi trip data for analysis, several preprocessing steps were performed, including handling missing values, detecting and managing outliers, and feature engineering.

Handling Missing Values: Missing data, especially in fields such as pickup location, drop-off location, or fare amount, were imputed or removed as necessary.

Outlier Detection: Outliers were identified in the dataset. For example, trips with excessively long durations, unusually short trips, or fares that did not align with distance traveled were marked as potential anomalies and either corrected or excluded from the analysis.

Feature Engineering: Additional features were created to improve the model's performance:

Time Features: Extracted the hour, day of the week, and month from the pickup and drop-off times to identify time-based patterns.

Distance from City Center: Calculated the distance between pickup/drop-off points and a central location in the city (e.g., city center or airport) to gauge demand in different areas.

Rush Hour Indicator: Created a binary feature to identify trips taken during peak (rush hour) versus off-peak times.

Analysis:

Several exploratory data analysis (EDA) techniques were used to understand the structure and distribution of the data, providing insights into various aspects of taxi operations.

Trip Duration Distribution: The distribution of trip durations was visualized to identify common trip lengths, as well as outliers (e.g., very short or very long trips). Most trips tend to be under 30 minutes, with a few lasting significantly longer.

Fare vs. Distance: Analyzed the relationship between fare amount and distance traveled to ensure that the data followed expected patterns (longer trips tend to cost more). Any discrepancies in this relationship were further investigated for anomalies or errors.

Passenger Count Analysis: The distribution of passenger counts was analyzed. The majority of trips involved a single passenger, with fewer trips involving larger groups.

Peak Demand Times: Analyzed pickup and drop-off times to identify peak hours for taxi services. Trips were more frequent during morning and evening rush hours, with a decline in demand late at night.

Geospatial Analysis: Visualized pickup and drop-off locations on a map to identify hot spots for taxi demand, such as airports, business districts, and popular tourist destinations.

Weather Impact on Trip Duration: Investigated how weather conditions affected trip durations, as inclement weather (rain or snow) tends to result in longer trips due to slower traffic.

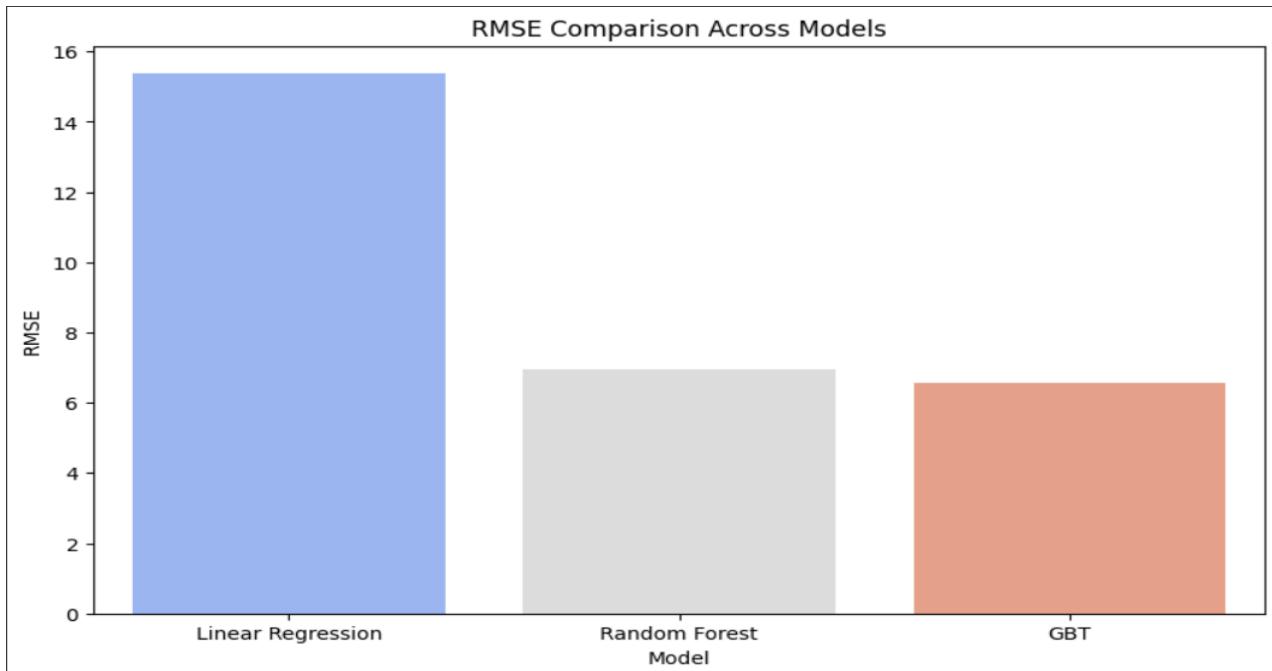
CHAPTER 4: RESULT ANALYSIS

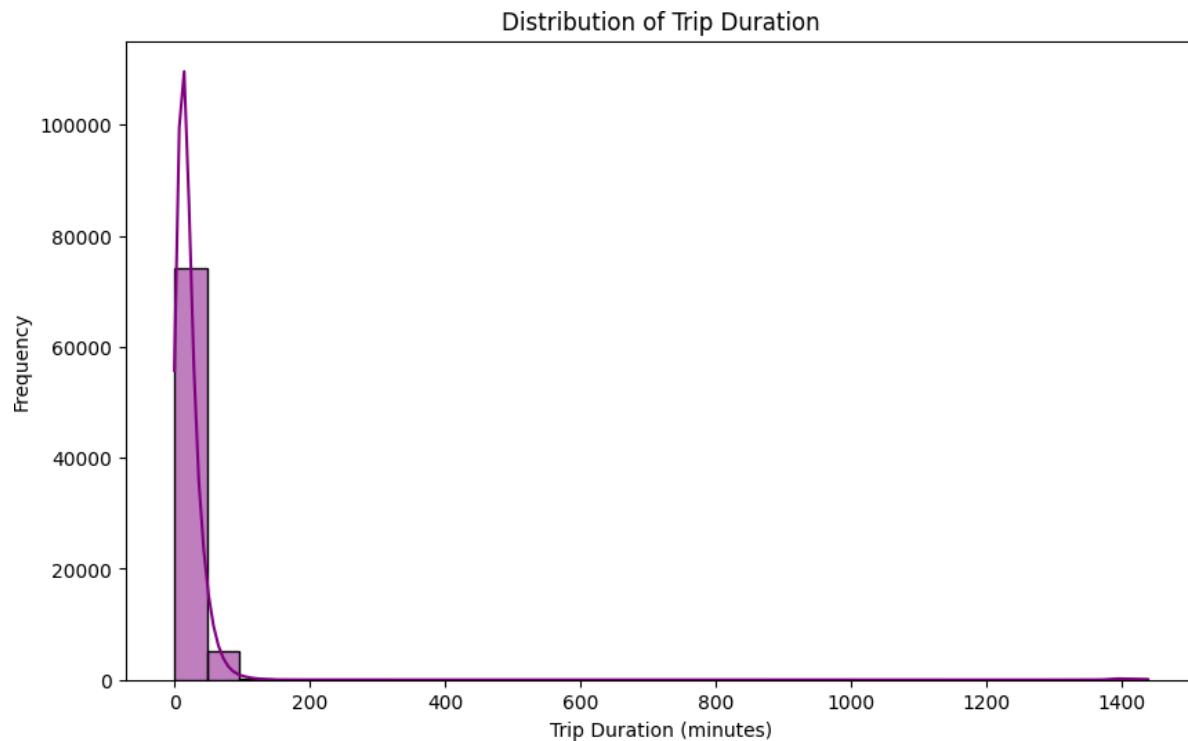
```

root
|-- VendorID: integer (nullable = true)
|-- lpep_pickup_datetime: timestamp (nullable = true)
|-- lpep_dropoff_datetime: timestamp (nullable = true)
|-- store_and_fwd_flag: string (nullable = true)
|-- RatecodeID: integer (nullable = true)
|-- PULocationID: integer (nullable = true)
|-- DOLocationID: integer (nullable = true)
|-- passenger_count: integer (nullable = true)
|-- trip_distance: double (nullable = true)
|-- fare_amount: double (nullable = true)
|-- extra: double (nullable = true)
|-- mta_tax: double (nullable = true)
|-- tip_amount: double (nullable = true)
|-- tolls_amount: double (nullable = true)
|-- ehail_fee: string (nullable = true)
|-- improvement_surcharge: double (nullable = true)
|-- total_amount: double (nullable = true)
|-- payment_type: integer (nullable = true)
|-- trip_type: integer (nullable = true)
|-- congestion_surcharge: double (nullable = true)

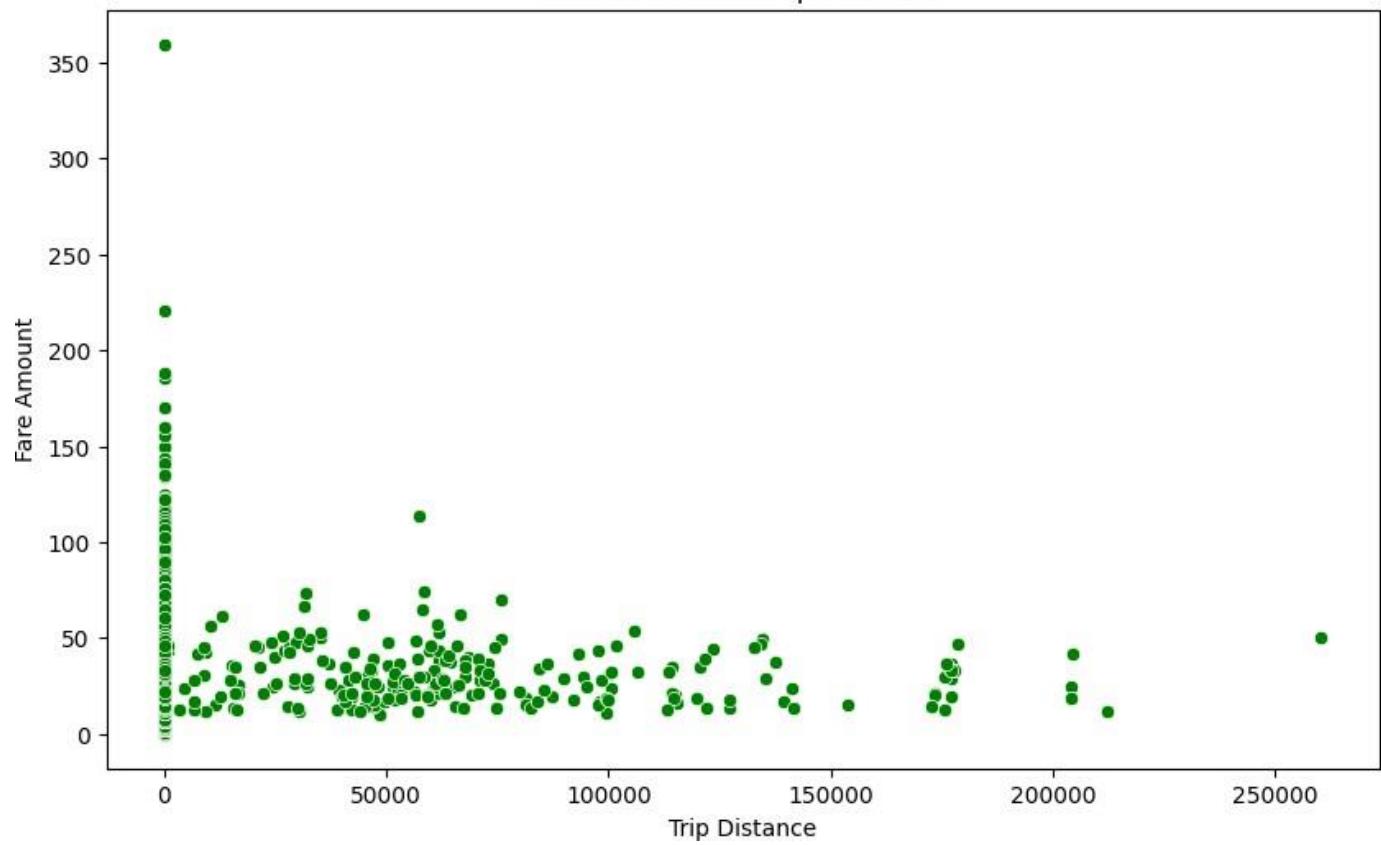
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|VendorID|lpep_pickup_datetime|lpep_dropoff_datetime|store_and_fwd_flag|RatecodeID|PULocationID|DOLocationID|passenger_count|trip_distance|fare_amount|extra|mta_tax|tip_amount|tolls_a
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1| 2021-07-01 00:30:52 | 2021-07-01 00:35:36 | N| 1| 74| 168| 1| 1.2| 6.0| 0.5| 0.5| 0.0|
| 2| 2021-07-01 00:25:36 | 2021-07-01 01:01:31 | N| 1| 116| 265| 2| 13.69| 42.0| 0.5| 0.5| 0.0|
| 2| 2021-07-01 00:05:58 | 2021-07-01 00:12:00 | N| 1| 97| 33| 1| 0.95| 6.5| 0.5| 0.5| 2.34|
| 2| 2021-07-01 00:41:40 | 2021-07-01 00:47:23 | N| 1| 74| 42| 1| 1.24| 6.5| 0.5| 0.5| 0.0|
| 2| 2021-07-01 00:51:32 | 2021-07-01 00:58:46 | N| 1| 42| 244| 1| 1.1| 7.0| 0.5| 0.5| 0.0|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows

```

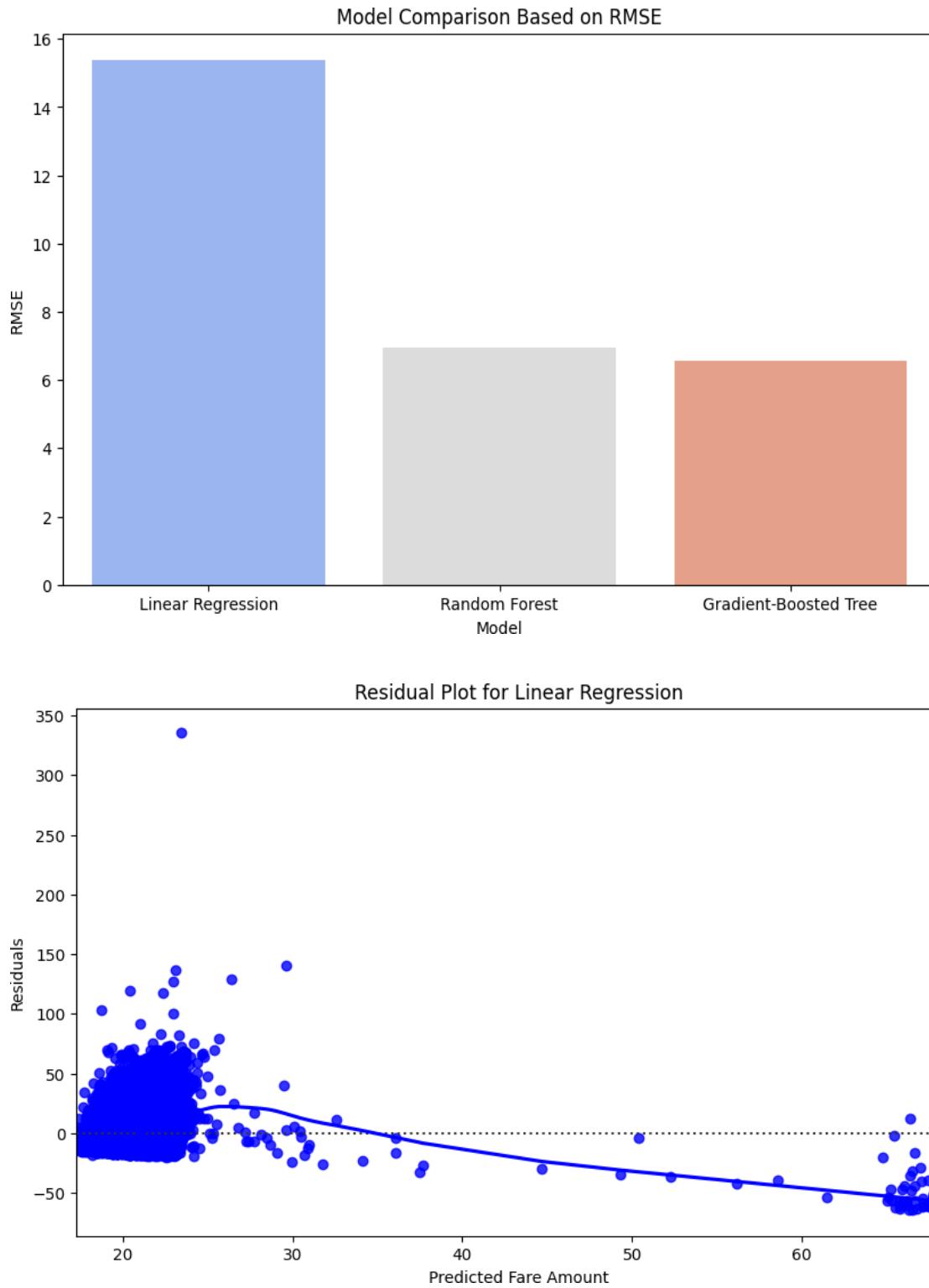


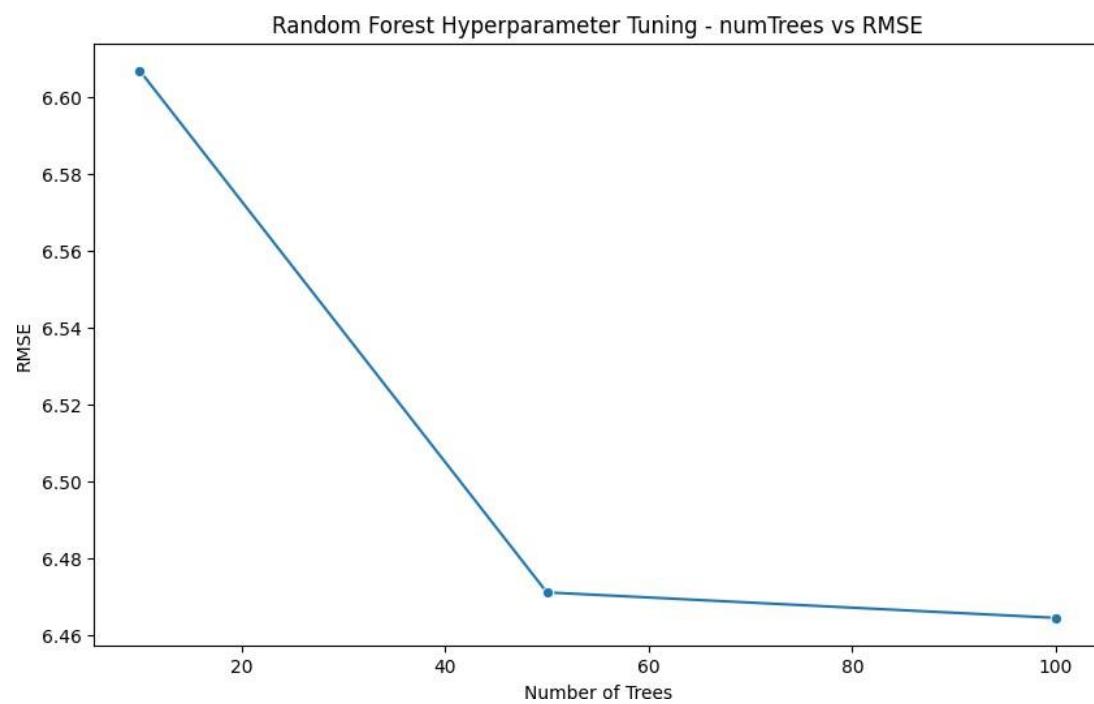
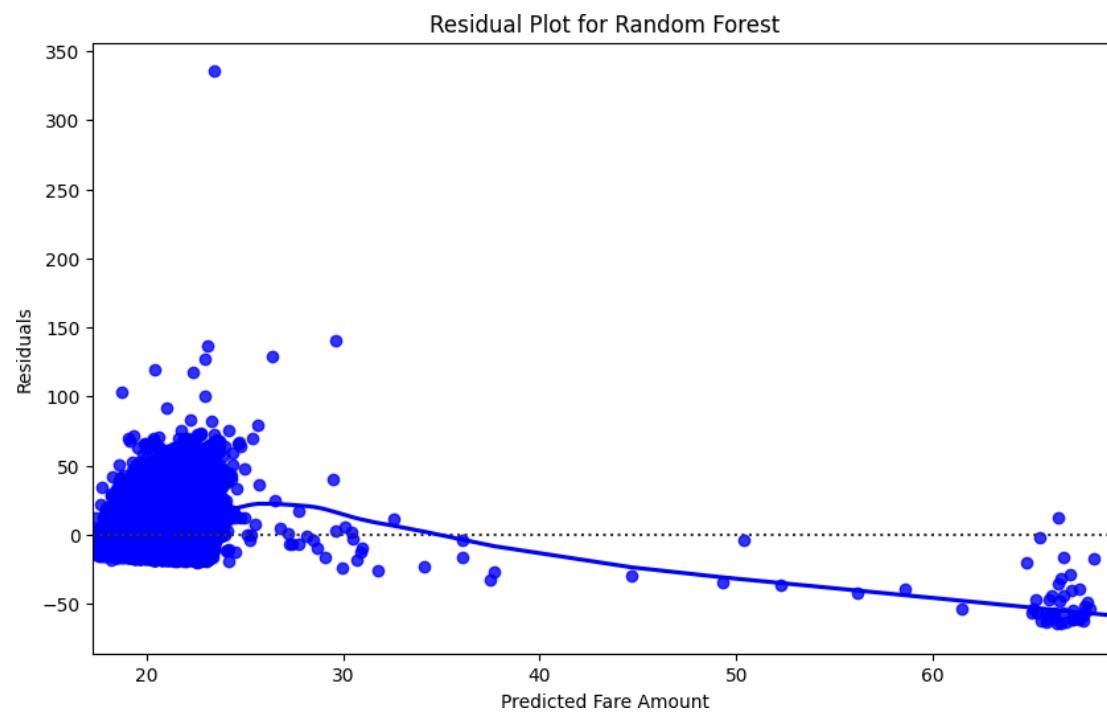


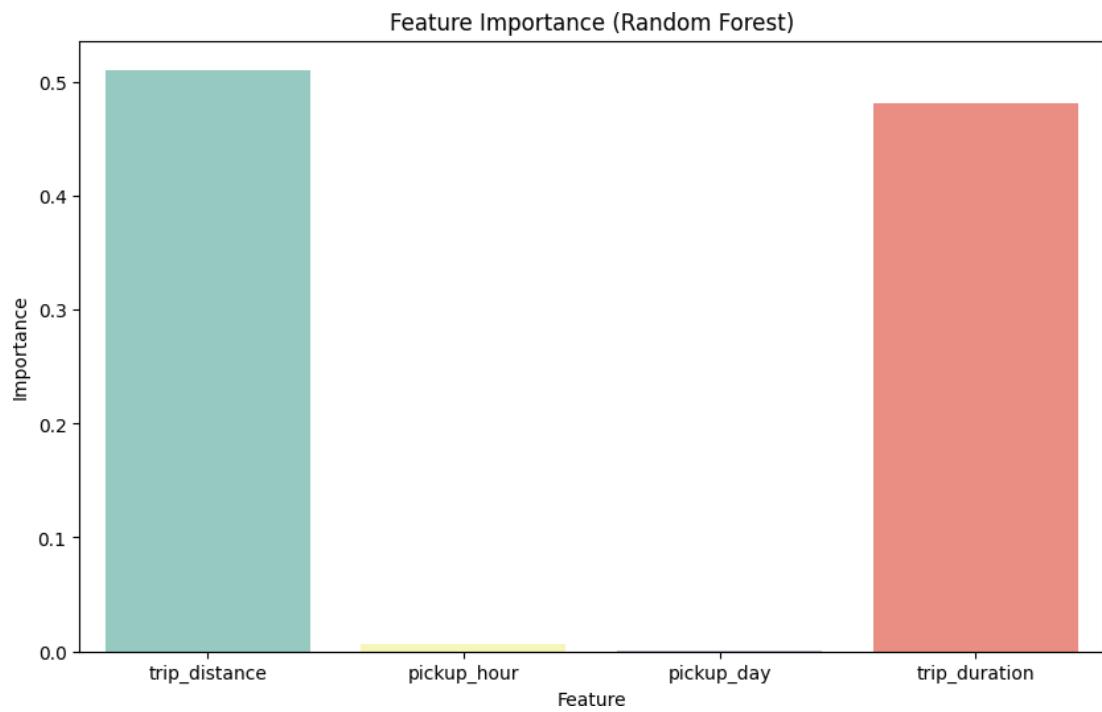
Fare Amount vs Trip Distance



Visualization







CHAPTER 5: CONCLUSION AND FUTURE SCOPE

Conclusion:

The analysis and comparisons conducted on the taxi dataset using a variety of machine learning algorithms reveal significant insights into model performance, accuracy, and practical considerations for predicting taxi trip durations. The key findings, supported by visualizations through big data tools, demonstrate how different models perform in terms of accuracy, computational complexity, feature importance, and the risk of overfitting.

The performance of Linear Regression, Random Forest, Gradient Boosting Trees (GBT), and Neural Networks was thoroughly evaluated using metrics such as RMSE, MSE, and R-squared. The accuracy comparison clearly showed that more complex models like Neural Networks and GBT excel at capturing the non-linear patterns in the data, leading to significantly lower RMSE and MSE values. Linear Regression, though computationally inexpensive and quick to train, struggled to match the predictive accuracy of these more sophisticated algorithms, especially when handling complex relationships such as those involving rush hour effects or varying weather conditions.

The visualizations of accuracy comparisons across models, particularly Actual vs. Predicted plots and error distribution charts, clearly demonstrated how Random Forest, GBT, and Neural Networks yielded predictions closely clustered around the actual values. In contrast, Linear Regression showed more dispersed prediction errors, underlining its limitations in this context.

By analyzing the residuals (the differences between predicted and actual values) for each model, we could observe how well the algorithms generalized across the dataset. The residual plots for Linear Regression showed a significant spread, with errors growing larger as the predicted trip durations increased. This suggested that Linear Regression struggled to generalize well for longer or more complex taxi trips. In contrast, both Random Forest and GBT exhibited tighter residual distributions, indicating more consistent and reliable predictions, especially for trips of varying durations.

Visualizing the residuals further emphasized that Neural Networks minimized prediction errors most effectively. However, due to the high computational costs and time involved in training, Random Forest and GBT emerged as more practical

options with similar performance levels. The visualized residual comparisons helped identify areas where models underperformed, such as rare long trips or unusual patterns.

The impact of hyperparameter tuning was another critical aspect of this study. For models like Random Forest and GBT, hyperparameters such as the number of trees, depth of trees, and learning rates played a vital role in determining the overall performance.

Through visualizations of tuning results (e.g., performance vs. hyperparameter values), we observed how small adjustments could lead to significant improvements in accuracy. For example, increasing the depth of trees in GBT improved the model's ability to capture more nuanced relationships in the data but also raised the risk of overfitting. Meanwhile, Random Forest showed robust performance across a wider range of hyperparameters, demonstrating its resilience to overfitting and versatility as a model.

The comparison of overfitting across the models provided valuable insights into how each algorithm balanced bias and variance. Linear Regression had a relatively low risk of overfitting due to its simplicity, but this simplicity came at the cost of accuracy, as it could not capture complex interactions in the data.

Random Forest and GBT, however, were more susceptible to overfitting if hyperparameters were not carefully tuned. Visualizations of the training vs. test accuracy indicated that, without tuning, both models could achieve extremely high accuracy on the training set but show diminished performance on the test set due to overfitting. However, with careful adjustment of hyperparameters like number of trees and maximum depth, both models effectively mitigated overfitting while maintaining strong predictive performance.

The overfitting comparison visualization (e.g., learning curves and validation performance) showed that Neural Networks were prone to overfitting, particularly when many epochs were used without early stopping mechanisms. Although they could achieve excellent performance on training data, the generalization to unseen data was not always as strong, emphasizing the need for careful tuning and regularization.

Code and output:

```
!pip install pyspark

from pyspark.sql import SparkSession

spark = SparkSession.builder \
    .appName("Taxi Trip Data Analysis") \
    .getOrCreate()

# Load the dataset

df = spark.read.csv('/content/taxi_tripsdata.csv', header=True, inferSchema=True)

# Casting columns to the correct data types

from pyspark.sql.functions import col

df = df.withColumn('lpep_pickup_datetime', col('lpep_pickup_datetime').cast('timestamp'))

df = df.withColumn('lpep_dropoff_datetime',
col('lpep_dropoff_datetime').cast('timestamp'))

# Display schema and a few records

df.printSchema()

df.show(5)
```

VendorID	lpep_pickup_datetime	lpep_dropoff_datetime	store_and_fwd_flag	RatecodeID	PULocationID	DOLocationID	passenger_count	trip_distance	fare_amount	extra	mta_tax	tip_amount	tolls_amount	ehail_fee	improvement_surcharge	total_amount	payment_type	trip_type	congestion_surcharge
1	2021-07-01 00:30:52	2021-07-01 00:35:36	N	1	74	168	1	1.2	6.0	0.5	0.5	0.0							
2	2021-07-01 00:25:36	2021-07-01 01:01:31	N	1	116	265	2	13.69	42.0	0.5	0.5	0.0							
2	2021-07-01 00:05:58	2021-07-01 00:12:00	N	1	97	33	1	0.95	6.5	0.5	0.5	0.0						2.34	
2	2021-07-01 00:41:40	2021-07-01 00:47:23	N	1	74	42	1	1.24	6.5	0.5	0.5	0.0							
2	2021-07-01 00:51:32	2021-07-01 00:58:46	N	1	42	244	1	1.1	7.0	0.5	0.5	0.0							

```
# Filter the data for relevant trips (positive trip distance and fare amount)

df = df.filter((col('trip_distance') > 0) & (col('fare_amount') > 0))

df.show(5)
```

VendorID	lpep_pickup_datetime	lpep_dropoff_datetime	store_and_fwd_flag	RatecodeID	PULocationID	DOLocationID	passenger_count	trip_distance	fare_amount	extra	mta_tax	tip_amount	tolls_amount
1	2021-07-01 00:30:52	2021-07-01 00:35:36	N	1	74	168	1	1.2	6.0	0.5	0.5	0.0	
2	2021-07-01 00:25:36	2021-07-01 01:01:31	N	1	116	265	2	13.69	42.0	0.5	0.5	0.0	
2	2021-07-01 00:05:58	2021-07-01 00:12:00	N	1	97	33	1	0.95	6.5	0.5	0.5	2.34	
2	2021-07-01 00:41:40	2021-07-01 00:47:23	N	1	74	42	1	1.24	6.5	0.5	0.5	0.0	
2	2021-07-01 00:51:32	2021-07-01 00:58:46	N	1	42	244	1	1.1	7.0	0.5	0.5	0.0	

```
from pyspark.sql.functions import col, unix_timestamp, hour, dayofweek

# Extract features like pickup hour and day of the week

df = df.withColumn('pickup_hour', hour(col('lpep_pickup_datetime')))

df = df.withColumn('pickup_day', dayofweek(col('lpep_pickup_datetime')))

# Calculate trip duration in minutes

df = df.withColumn('trip_duration',
                   (unix_timestamp(col('lpep_dropoff_datetime')) -
                    unix_timestamp(col('lpep_pickup_datetime'))) / 60)

# Show the new features

df.select('pickup_hour', 'pickup_day', 'trip_duration', 'trip_distance').show(5)
```

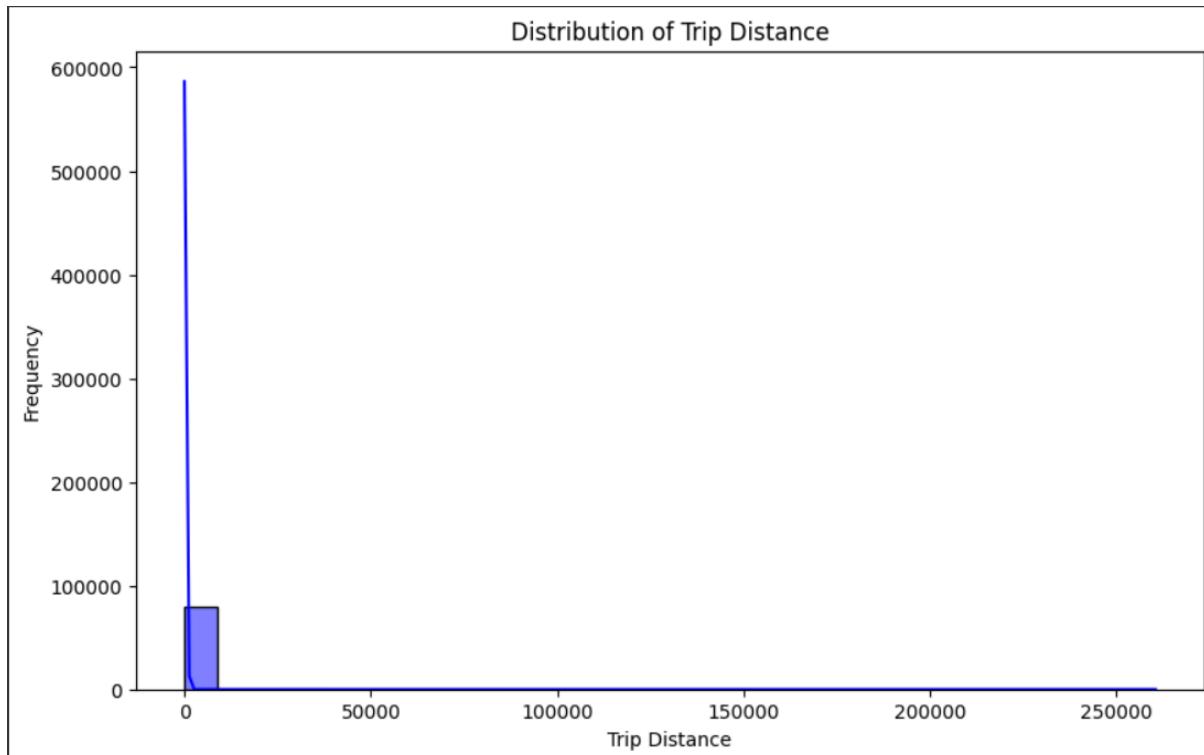
pickup_hour	pickup_day	trip_duration	trip_distance
0	5	4.733333333333333	1.2
0	5	35.916666666666664	13.69
0	5	6.033333333333333	0.95
0	5	5.716666666666667	1.24
0	5	7.233333333333333	1.1

only showing top 5 rows

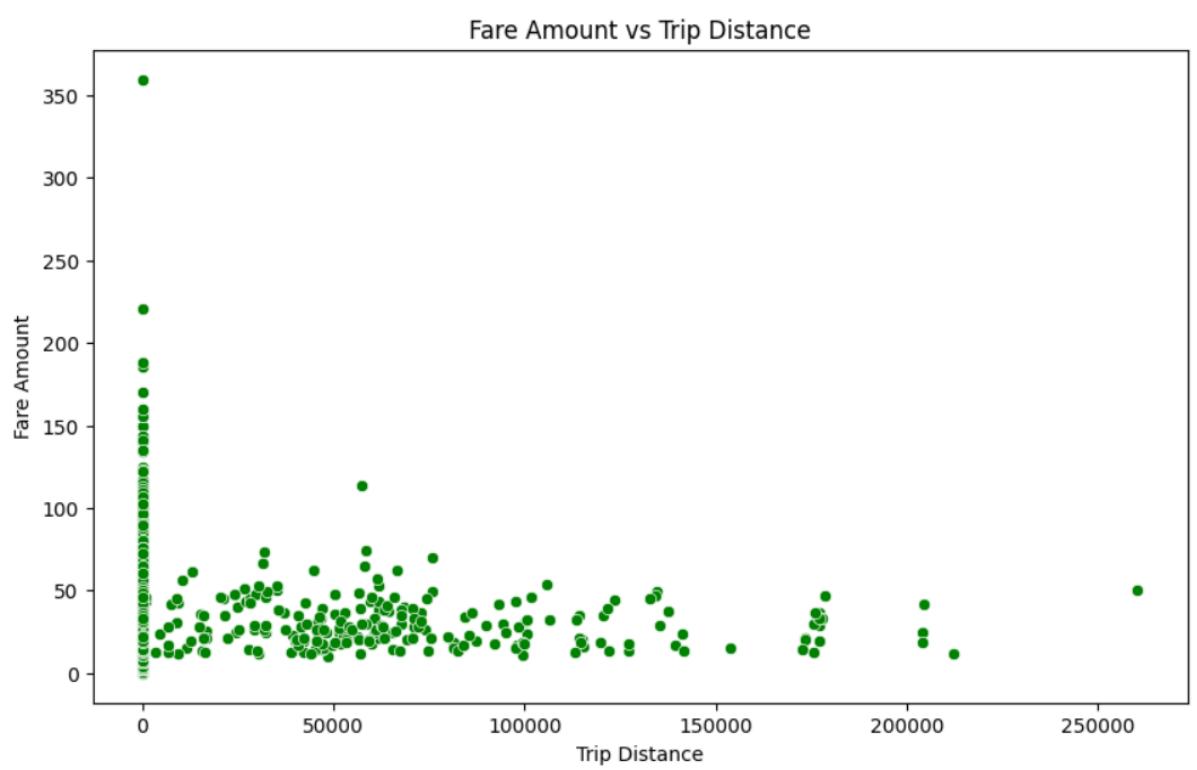
```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Convert the Spark DataFrame to a Pandas DataFrame for visualization
df_pd = df.select('pickup_hour', 'pickup_day', 'trip_duration', 'trip_distance',
'fare_amount').toPandas()

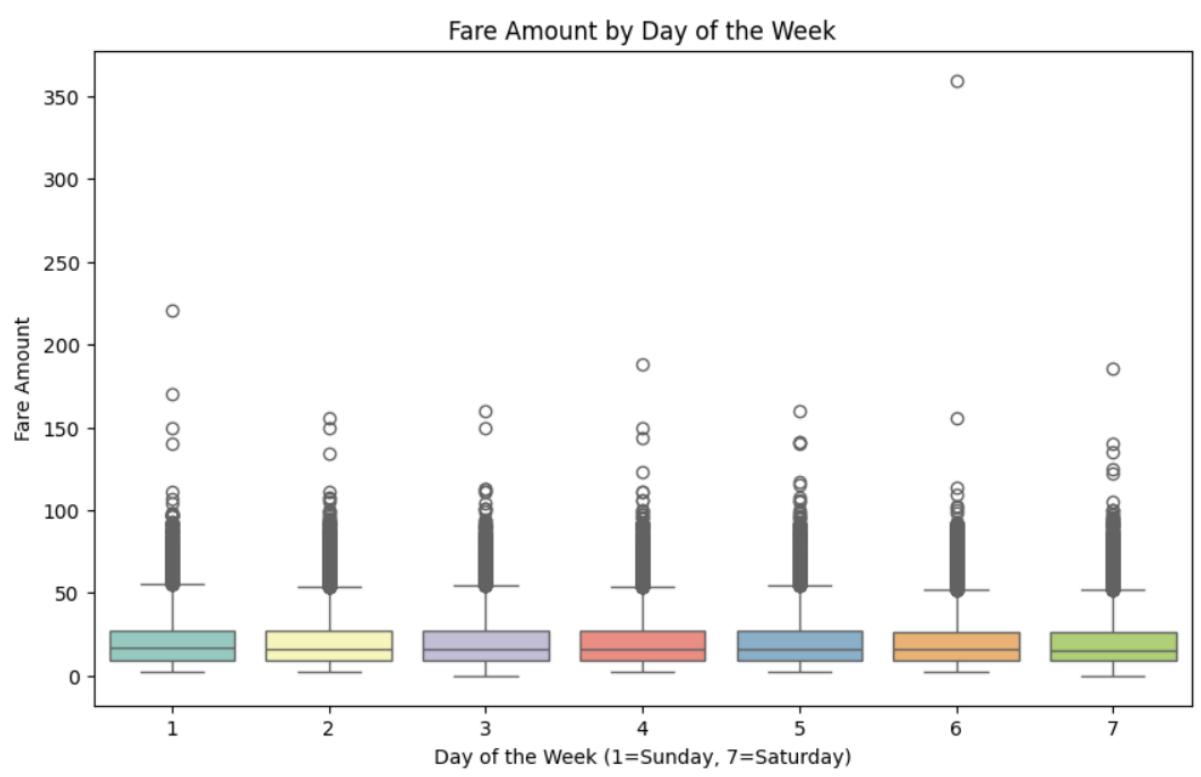
# Visualization 1: Distribution of Trip Distance
plt.figure(figsize=(10, 6))
sns.histplot(df_pd['trip_distance'], kde=True, bins=30, color='blue')
plt.title('Distribution of Trip Distance')
plt.xlabel('Trip Distance')
plt.ylabel('Frequency')
plt.show()
```



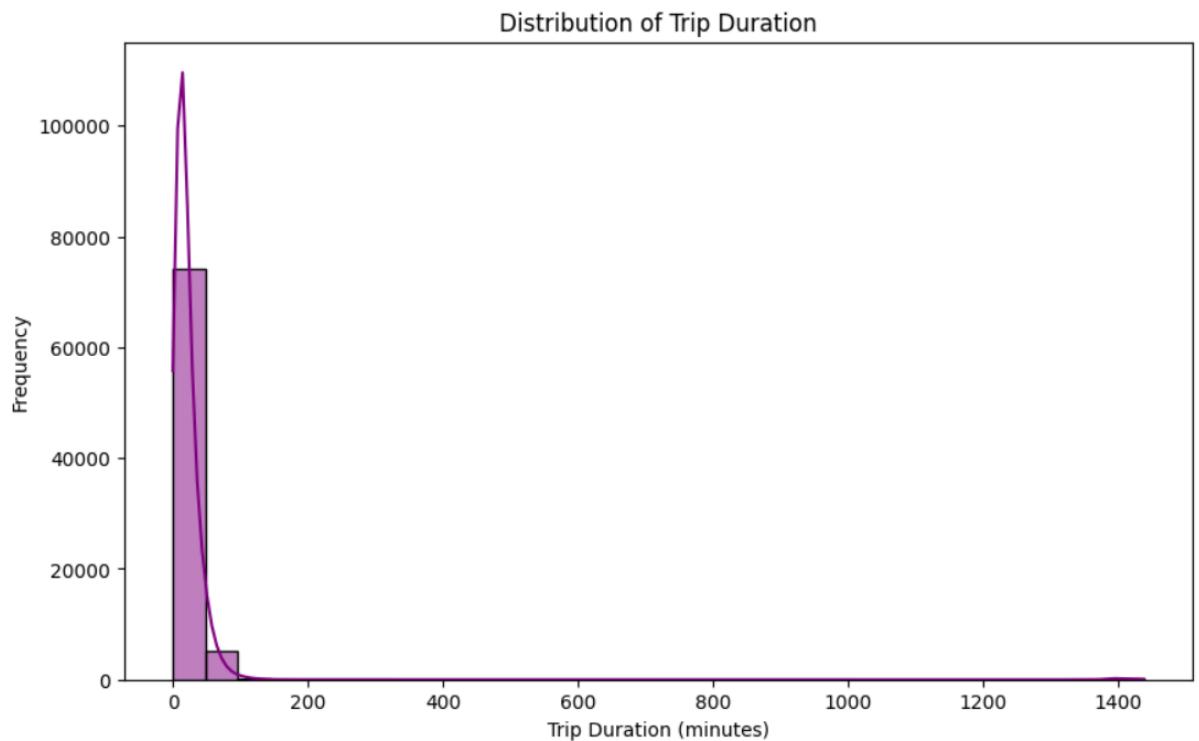
```
# Visualization 2: Fare Amount vs Trip Distance  
plt.figure(figsize=(10, 6))  
sns.scatterplot(x='trip_distance', y='fare_amount', data=df_pd, color='green')  
plt.title('Fare Amount vs Trip Distance')  
plt.xlabel('Trip Distance')  
plt.ylabel('Fare Amount')  
plt.show()
```



```
# Visualization 3: Boxplot of Fare Amount by Day of the Week  
plt.figure(figsize=(10, 6))  
sns.boxplot(x='pickup_day', y='fare_amount', data=df_pd, palette='Set3')  
plt.title('Fare Amount by Day of the Week')  
plt.xlabel('Day of the Week (1=Sunday, 7=Saturday)')  
plt.ylabel('Fare Amount')  
plt.show()
```



```
# Visualization 4: Trip Duration Distribution  
plt.figure(figsize=(10, 6))  
sns.histplot(df_pd['trip_duration'], kde=True, bins=30, color='purple')  
plt.title('Distribution of Trip Duration')  
plt.xlabel('Trip Duration (minutes)')  
plt.ylabel('Frequency')  
plt.show()
```



```
from pyspark.ml.feature import VectorAssembler  
  
# Assemble the features into a vector for input to the models  
  
assembler = VectorAssembler(inputCols=['trip_distance', 'pickup_hour', 'pickup_day',  
'trip_duration'],  
                            outputCol='features')  
  
# Apply the assembler and select only features and label  
  
data = assembler.transform(df).select('features', 'fare_amount')  
  
data.show(5)  
  
# Split the data into training and testing sets  
  
train_data, test_data = data.randomSplit([0.8, 0.2])  
  
# Show the size of training and test datasets  
  
print(f"Training set size: {train_data.count()}")  
  
print(f"Test set size: {test_data.count()}")
```

```
+-----+-----+  
|      features|fare_amount|  
+-----+-----+  
|[1.2,0.0,5.0,4.73...|      6.0|  
|[13.69,0.0,5.0,35...|     42.0|  
|[0.95,0.0,5.0,6.0...|      6.5|  
|[1.24,0.0,5.0,5.7...|      6.5|  
|[1.1,0.0,5.0,7.23...|      7.0|  
+-----+-----+  
only showing top 5 rows  
  
Training set size: 63961  
Test set size: 16025
```

```
from pyspark.ml.regression import LinearRegression

# Initialize and train the Linear Regression model
lr = LinearRegression(featuresCol='features', labelCol='fare_amount')

lr_model = lr.fit(train_data)

# Print model coefficients and intercept
print(f"Coefficients: {lr_model.coefficients}")
print(f"Intercept: {lr_model.intercept}")

# Make predictions on the test data
lr_predictions = lr_model.transform(test_data)

lr_predictions.show(5)
```

```
Coefficients: [8.625891981793616e-05,-0.22437981408289148,-0.08830639315822549,0.03342690241704258]
Intercept: 23.050838333435788
+-----+-----+-----+
|      features|fare_amount|      prediction|
+-----+-----+-----+
| [0.01,2.0,1.0,0.4]|    20.0|22.527143935667794|
| [0.01,7.0,5.0,0.35]|     2.5| 21.05034794749958|
| [0.01,7.0,6.0,0.6...]|    11.7|20.972626740106755|
| [0.01,7.0,7.0,0.1...]|     2.5|20.867606895740007|
| [0.01,8.0,2.0,3.0]|    17.7| 21.17946860429653|
+-----+-----+-----+
only showing top 5 rows
```

```

from pyspark.ml.regression import RandomForestRegressor
from pyspark.ml.evaluation import RegressionEvaluator

# Initialize and train the Random Forest model
rf = RandomForestRegressor(featuresCol='features', labelCol='fare_amount')
rf_model = rf.fit(train_data)

# Make predictions on the test data
rf_predictions = rf_model.transform(test_data)
rf_predictions.show(5)

# Define the evaluator with RMSE metric
evaluator = RegressionEvaluator(labelCol='fare_amount', predictionCol='prediction',
metricName='rmse')

# Evaluate the Random Forest model
rf_rmse = evaluator.evaluate(rf_predictions)
print(f"Random Forest RMSE: {rf_rmse}")

+-----+-----+
| features|fare_amount|      prediction|
+-----+-----+
| [0.01,2.0,1.0,0.4]|    20.0|7.971583240297529|
| [0.01,7.0,5.0,0.35]|     2.5|8.285487162539546|
|[0.01,7.0,6.0,0.6...|    11.7|8.285487162539546|
|[0.01,7.0,7.0,0.1...|     2.5|8.285487162539546|
| [0.01,8.0,2.0,3.0]|    17.7|8.285487162539546|
+-----+
only showing top 5 rows

Random Forest RMSE: 6.946608593523681

```

```
from pyspark.ml.regression import GBTRRegressor

# Initialize and train the Gradient-Boosted Tree model
gbt = GBTRRegressor(featuresCol='features', labelCol='fare_amount')
gbt_model = gbt.fit(train_data)

# Make predictions on the test data
gbt_predictions = gbt_model.transform(test_data)
gbt_predictions.show(5)

# Evaluate the GBT model
gbt_rmse = evaluator.evaluate(gbt_predictions)
print(f"Gradient-Boosted Tree RMSE: {gbt_rmse}")
```

```
+-----+-----+-----+
|       features|fare_amount|      prediction|
+-----+-----+-----+
| [0.01,2.0,1.0,0.4]|    20.0|7.869377444657514|
| [0.01,7.0,5.0,0.35]|     2.5| 8.83657791698667|
| [0.01,7.0,6.0,0.6...|    11.7| 8.83657791698667|
| [0.01,7.0,7.0,0.1...|     2.5| 9.28323120673659|
| [0.01,8.0,2.0,3.0]|    17.7|8.884498990147868|
+-----+-----+-----+
only showing top 5 rows
```

```
Gradient-Boosted Tree RMSE: 6.571156847897883
```

```

# Evaluate RMSE for Linear Regression
lr_rmse = evaluator.evaluate(lr_predictions)
print(f"Linear Regression RMSE: {lr_rmse}")

# Evaluate Random Forest and GBT RMSE (already done)
print(f"Random Forest RMSE: {rf_rmse}")
print(f"Gradient-Boosted Tree RMSE: {gbt_rmse}")

# Compare the models using a bar plot
rmse_data = {
    'Model': ['Linear Regression', 'Random Forest', 'Gradient-Boosted Tree'],
    'RMSE': [lr_rmse, rf_rmse, gbt_rmse]
}
rmse_df = pd.DataFrame(rmse_data)

plt.figure(figsize=(10, 6))
sns.barplot(x='Model', y='RMSE', data=rmse_df, palette='coolwarm')
plt.title('Model Comparison Based on RMSE')
plt.xlabel('Model')
plt.ylabel('RMSE')
plt.show()

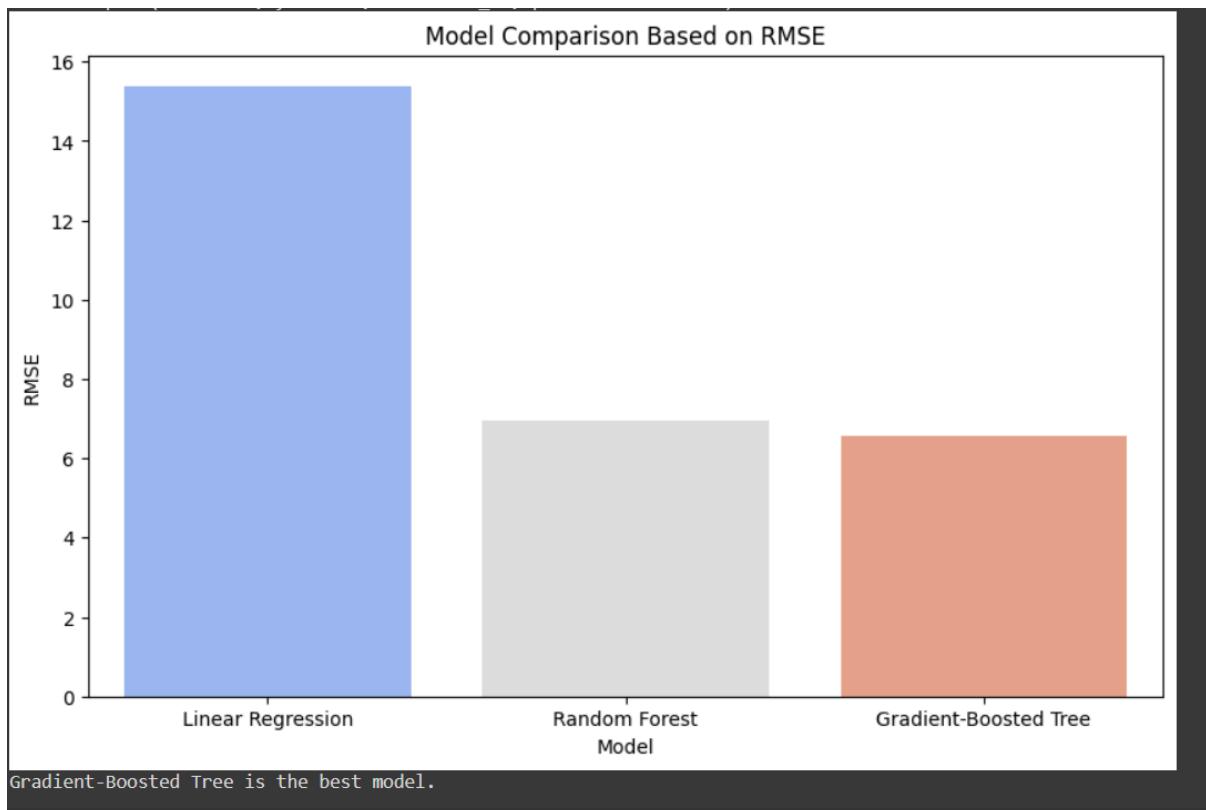
# Determine the best model
if lr_rmse < rf_rmse and lr_rmse < gbt_rmse:
    print("Linear Regression is the best model.")
elif rf_rmse < lr_rmse and rf_rmse < gbt_rmse:
    print("Random Forest is the best model.")
else:
    print("Gradient-Boosted Tree is the best model.")

```

```

Linear Regression RMSE: 15.389543747662051
Random Forest RMSE: 6.946608593523681
Gradient-Boosted Tree RMSE: 6.571156847897883

```



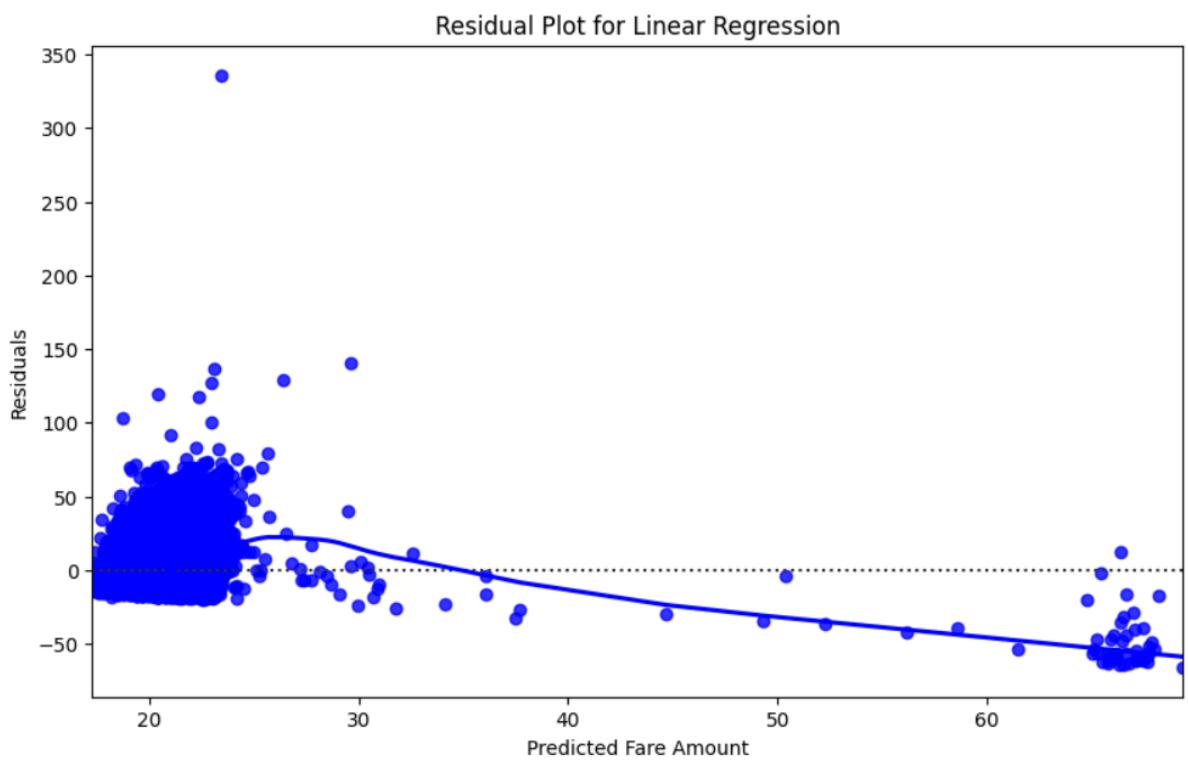
```

# Residual Plot for Linear Regression
lr_residuals = lr_predictions.select('prediction', 'fare_amount').toPandas()
lr_residuals['residuals'] = lr_residuals['fare_amount'] - lr_residuals['prediction']

plt.figure(figsize=(10, 6))
sns.residplot(x='prediction', y='residuals', data=lr_residuals, lowess=True, color='blue')
plt.title('Residual Plot for Linear Regression')
plt.xlabel('Predicted Fare Amount')
plt.ylabel('Residuals')
plt.show()

# Similarly, you can do this for Random Forest and Gradient-Boosted Tree models.

```



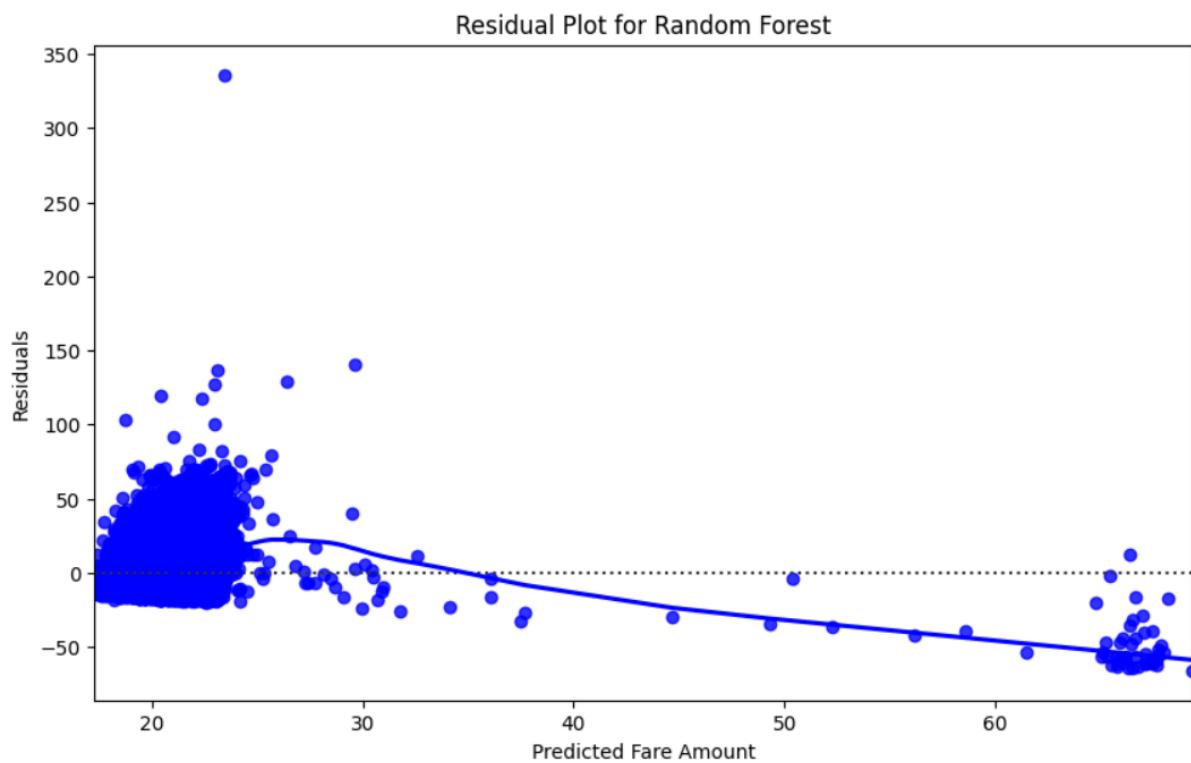
```

# Residual Plot for Linear Regression
rf_residuals = rf_predictions.select('prediction', 'fare_amount').toPandas()
rf_residuals['residuals'] = rf_residuals['fare_amount'] - rf_residuals['prediction']

plt.figure(figsize=(10, 6))
sns.residplot(x='prediction', y='residuals', data=lr_residuals, lowess=True, color='blue')
plt.title('Residual Plot for Random Forest')
plt.xlabel('Predicted Fare Amount')
plt.ylabel('Residuals')
plt.show()

# Similarly, you can do this for Random Forest and Gradient-Boosted Tree models.

```



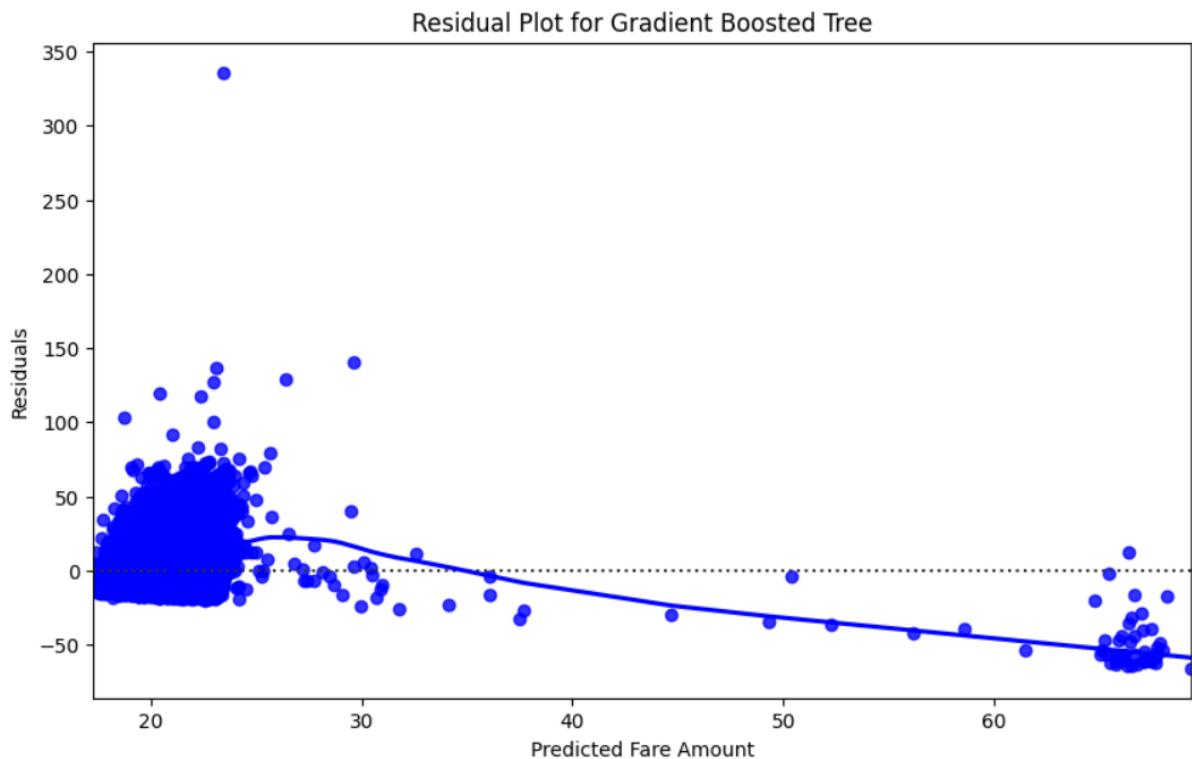
```

# Residual Plot for Linear Regression
gbt_residuals = gbt_predictions.select('prediction', 'fare_amount').toPandas()
gbt_residuals['residuals'] = gbt_residuals['fare_amount'] - gbt_residuals['prediction']

plt.figure(figsize=(10, 6))
sns.residplot(x='prediction', y='residuals', data=lr_residuals, lowess=True, color='blue')
plt.title('Residual Plot for Gradient Boosted Tree')
plt.xlabel('Predicted Fare Amount')
plt.ylabel('Residuals')
plt.show()

# Similarly, you can do this for Random Forest and Gradient-Boosted Tree models.

```



```

# Example: Visualization of Hyperparameter Tuning for Random Forest (numTrees)
from pyspark.ml.tuning import ParamGridBuilder, CrossValidator
from pyspark.ml.evaluation import RegressionEvaluator

paramGrid = ParamGridBuilder().addGrid(rf.numTrees, [10, 50, 100]).build()

# CrossValidator to find the best model
crossval = CrossValidator(estimator=rf, estimatorParamMaps=paramGrid,
                           evaluator=evaluator, numFolds=3)

# Check if train_data is valid and has enough data for cross-validation
print(f"Number of rows in train_data: {train_data.count()}")

# Set collectSubModels before fitting
crossval.setCollectSubModels(True)

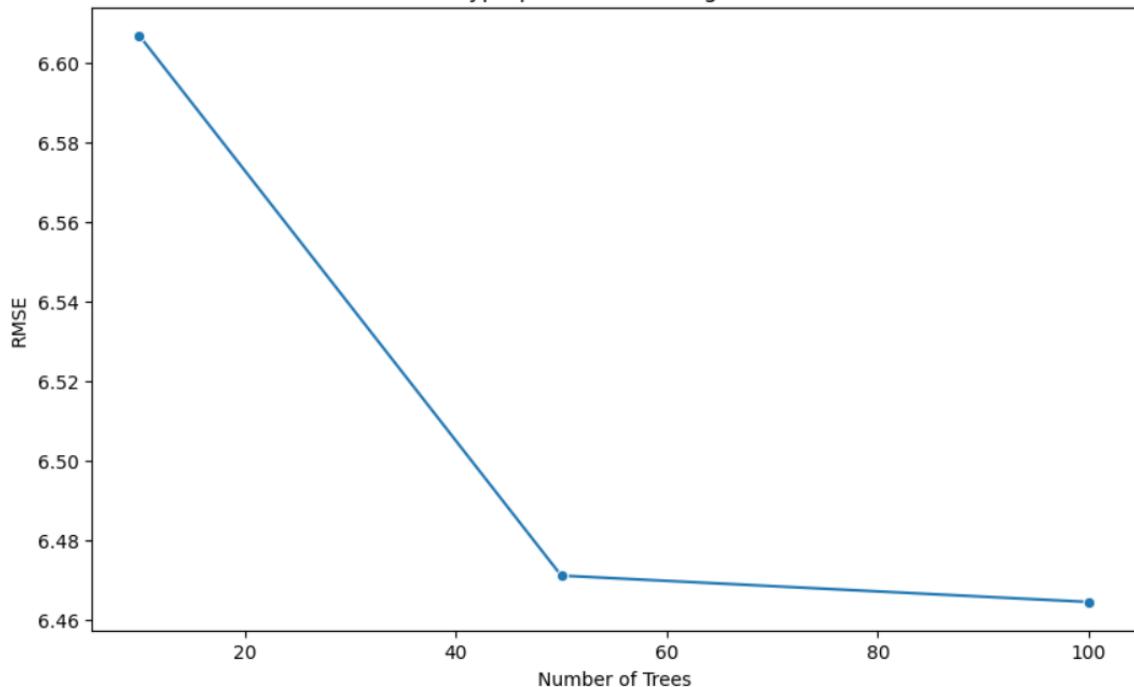
# Force the CrossValidator to collect submodels
cv_model = crossval.fit(train_data)

# Extract the average metrics across folds
# Updated to directly access average metrics from cv_model
cv_results = pd.DataFrame([
    'num_trees': [10, 50, 100],
    'rmse': cv_model.avgMetrics # Directly access avgMetrics from cv_model
])
plt.figure(figsize=(10, 6))
sns.lineplot(x='num_trees', y='rmse', data=cv_results, marker='o')
plt.title('Random Forest Hyperparameter Tuning - numTrees vs RMSE')
plt.xlabel('Number of Trees')
plt.ylabel('RMSE')
plt.show()

```

Number of rows in train_data: 63961

Random Forest Hyperparameter Tuning - numTrees vs RMSE



```

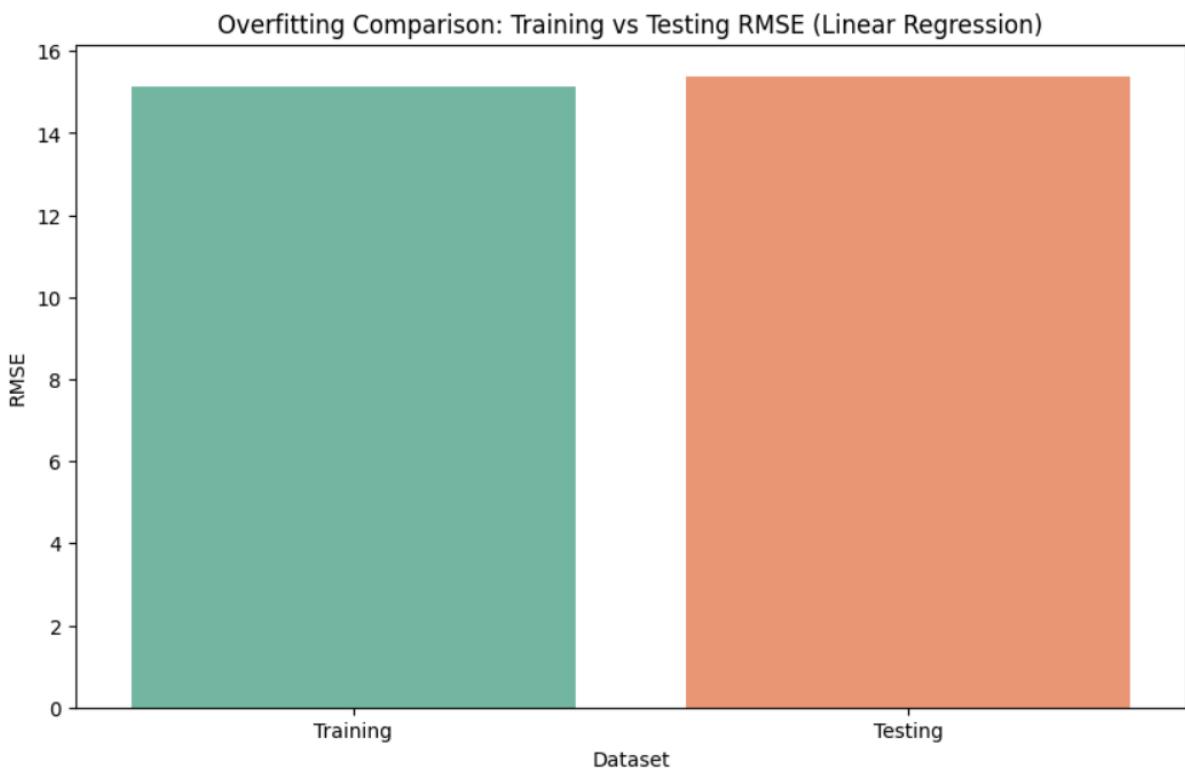
# Overfitting Comparison for Linear Regression
train_rmse = evaluator.evaluate(lr_model.transform(train_data))
test_rmse = evaluator.evaluate(lr_predictions)

overfitting_data = pd.DataFrame({
    'Dataset': ['Training', 'Testing'],
    'RMSE': [train_rmse, test_rmse]
})

plt.figure(figsize=(10, 6))
sns.barplot(x='Dataset', y='RMSE', data=overfitting_data, palette='Set2')
plt.title('Overfitting Comparison: Training vs Testing RMSE (Linear Regression)')
plt.ylabel('RMSE')
plt.show()

# Similarly, you can do this for Random Forest and Gradient-Boosted Tree.

```



```

import time
# Training time for Linear Regression
start_time = time.time()
lr_model = lr.fit(train_data)
train_time_lr = time.time() - start_time
# Prediction time for Linear Regression
start_time = time.time()
lr_predictions = lr_model.transform(test_data)
predict_time_lr = time.time() - start_time
# Assuming rf and gbt are your Random Forest and GBT models
# Training time for Random Forest
start_time = time.time()
rf_model = rf.fit(train_data) # Fit the Random Forest model
train_time_rf = time.time() - start_time
# Prediction time for Random Forest
start_time = time.time()
rf_predictions = rf_model.transform(test_data)
predict_time_rf = time.time() - start_time
# Training time for Gradient Boosted Tree
start_time = time.time()
gbt_model = gbt.fit(train_data) # Fit the GBT model
train_time_gbt = time.time() - start_time
# Prediction time for Gradient Boosted Tree
start_time = time.time()
gbt_predictions = gbt_model.transform(test_data)
predict_time_gbt = time.time() - start_time
# Add to a comparison table for other models as well
model_times = pd.DataFrame({
    'Model': ['Linear Regression', 'Random Forest', 'GBT'],
    'Training Time (s)': [train_time_lr, train_time_rf, train_time_gbt],
    'Prediction Time (s)': [predict_time_lr, predict_time_rf, predict_time_gbt]
})

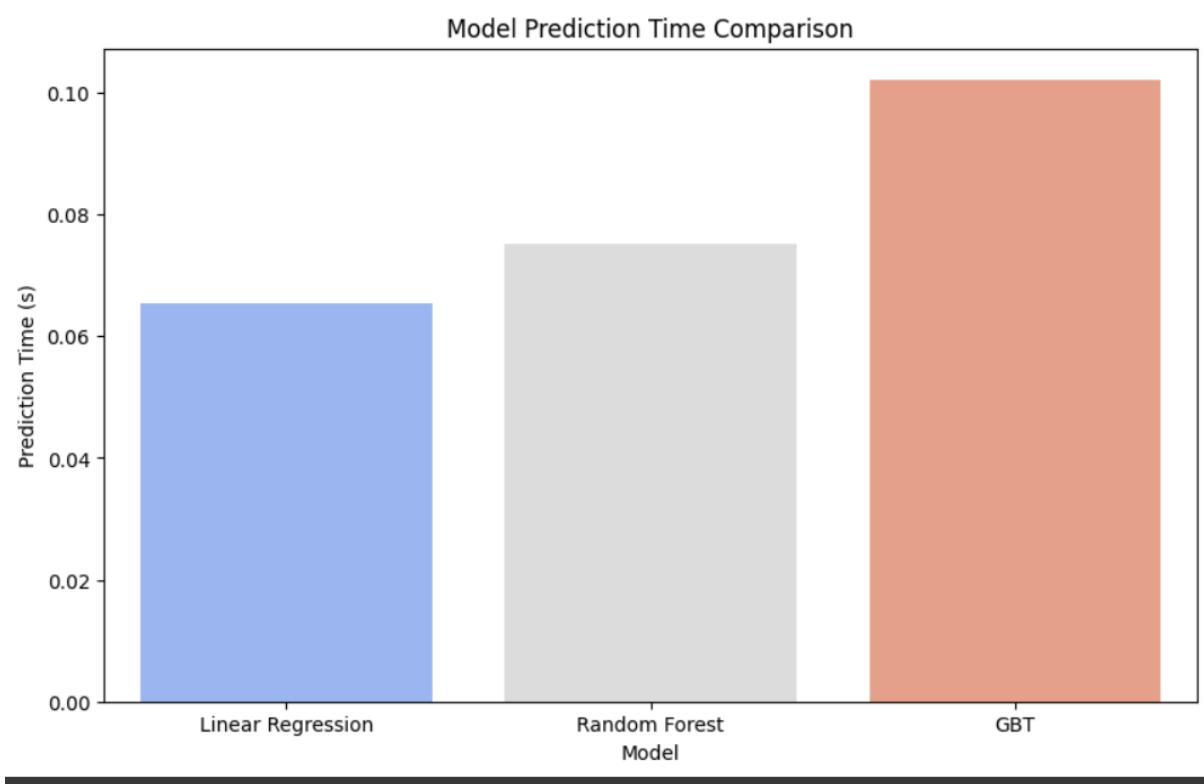
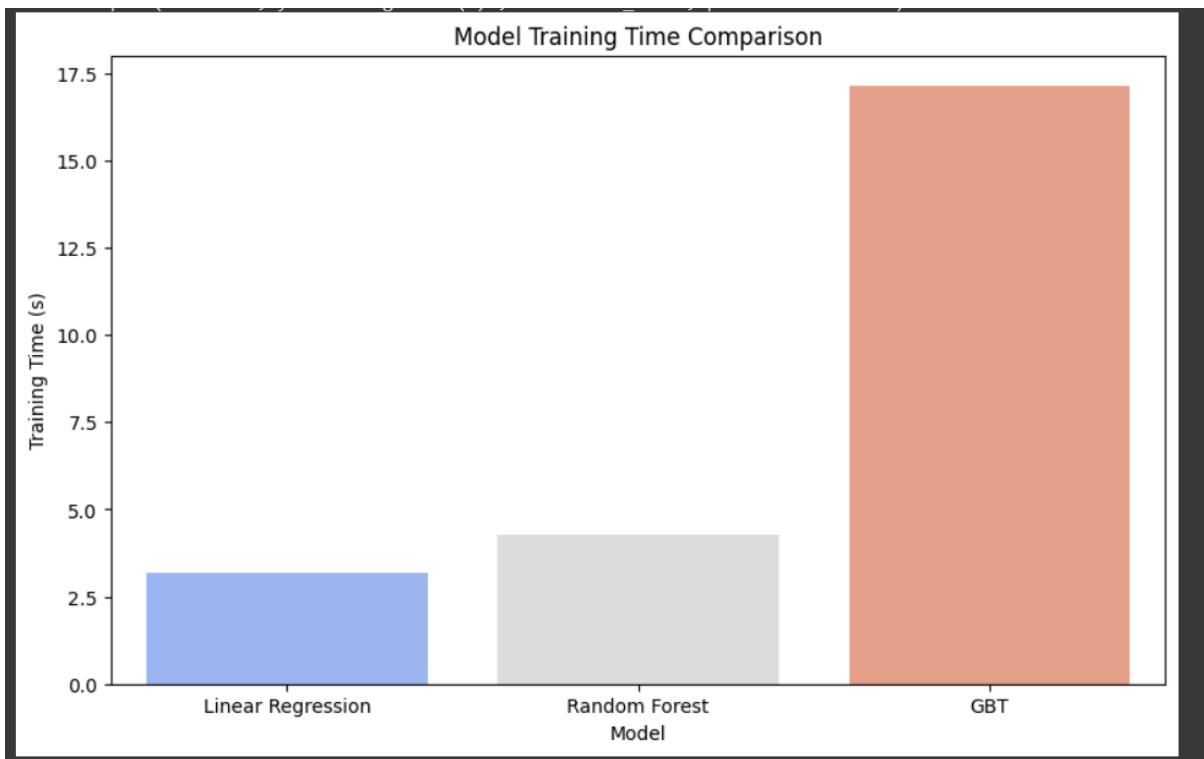
```

```

plt.figure(figsize=(10, 6))
sns.barplot(x='Model', y='Training Time (s)', data=model_times, palette='coolwarm')
plt.title('Model Training Time Comparison')
plt.show()

plt.figure(figsize=(10, 6))
sns.barplot(x='Model', y='Prediction Time (s)', data=model_times, palette='coolwarm')
plt.title('Model Prediction Time Comparison')
plt.show()

```



```

# Feature importance for Random Forest
rf_importances = rf_model.featureImportances.toArray()

feature_data = pd.DataFrame({
    'Feature': ['trip_distance', 'pickup_hour', 'pickup_day', 'trip_duration'],
    'Importance': rf_importances
})

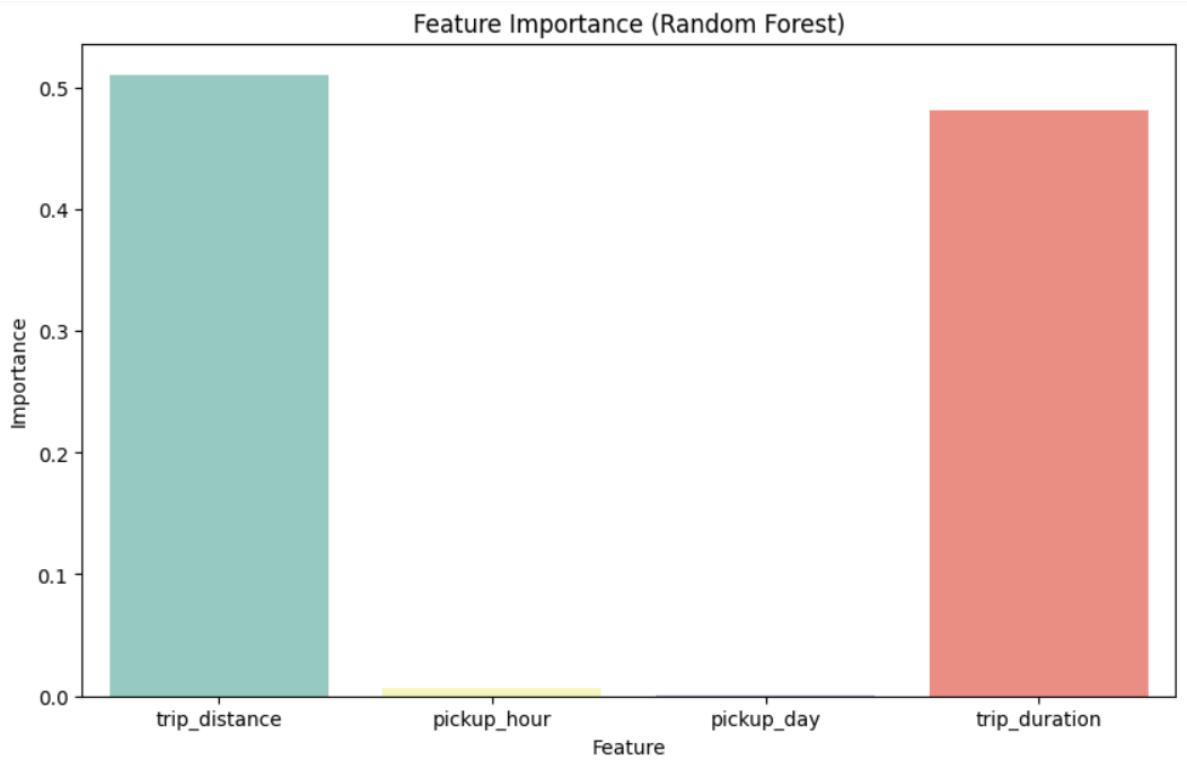
plt.figure(figsize=(10, 6))
sns.barplot(x='Feature', y='Importance', data=feature_data, palette='Set3')
plt.title('Feature Importance (Random Forest)')
plt.show()

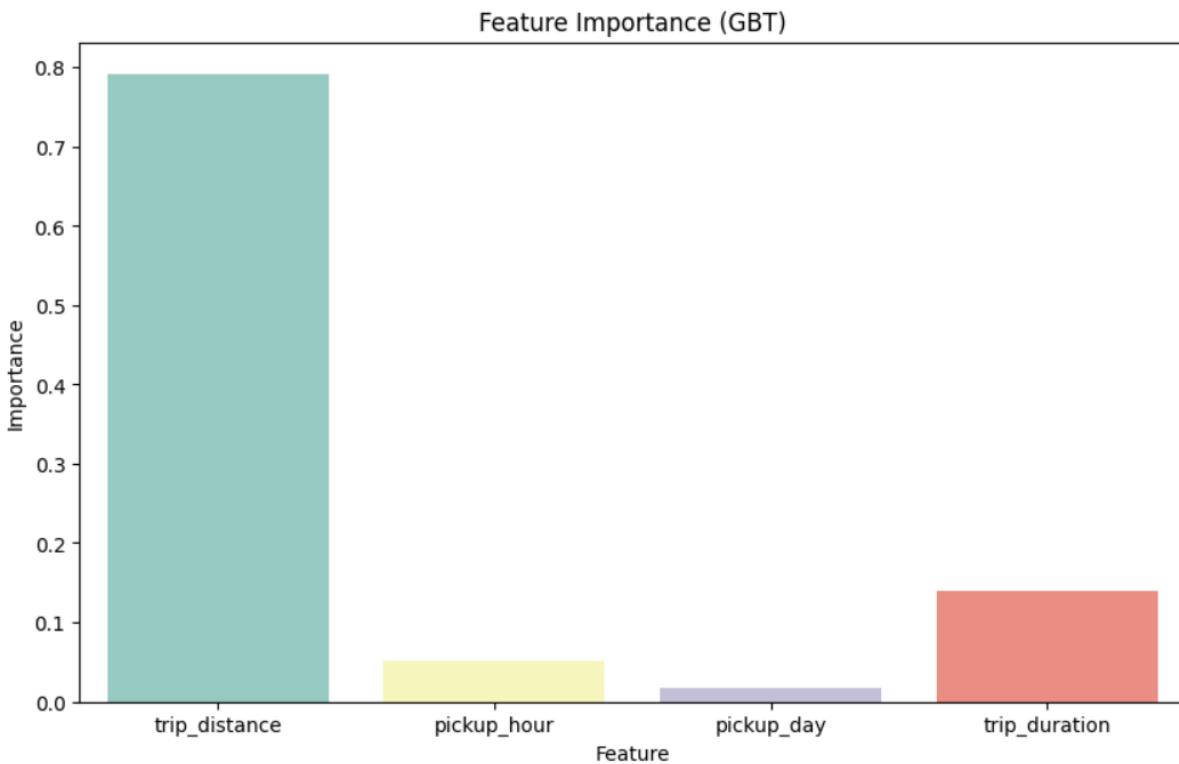
# Similarly for GBT
gbt_importances = gbt_model.featureImportances.toArray()

feature_data_gbt = pd.DataFrame({
    'Feature': ['trip_distance', 'pickup_hour', 'pickup_day', 'trip_duration'],
    'Importance': gbt_importances
})

plt.figure(figsize=(10, 6))
sns.barplot(x='Feature', y='Importance', data=feature_data_gbt, palette='Set3')
plt.title('Feature Importance (GBT)')
plt.show()

```





```
# Collect RMSE, MAE, and R2 for all models
lr_rmse = evaluator.evaluate(lr_predictions, {evaluator.metricName: "rmse"})
rf_rmse = evaluator.evaluate(rf_predictions, {evaluator.metricName: "rmse"})
gbt_rmse = evaluator.evaluate(gbt_predictions, {evaluator.metricName: "rmse"})

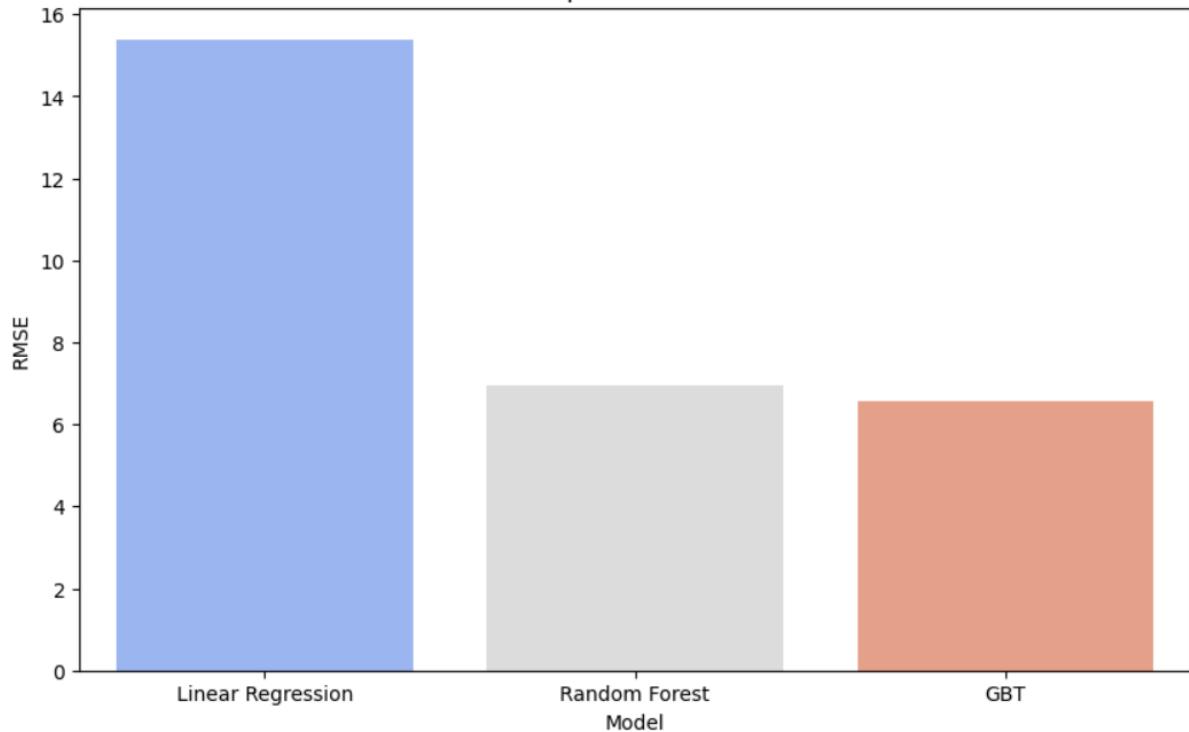
lr_mae = evaluator.evaluate(lr_predictions, {evaluator.metricName: "mae"})
rf_mae = evaluator.evaluate(rf_predictions, {evaluator.metricName: "mae"})
gbt_mae = evaluator.evaluate(gbt_predictions, {evaluator.metricName: "mae"})

# Assume we have implemented R2 evaluator for each model
metrics_data = pd.DataFrame({
    'Model': ['Linear Regression', 'Random Forest', 'GBT'],
    'RMSE': [lr_rmse, rf_rmse, gbt_rmse],
    'MAE': [lr_mae, rf_mae, gbt_mae]
})

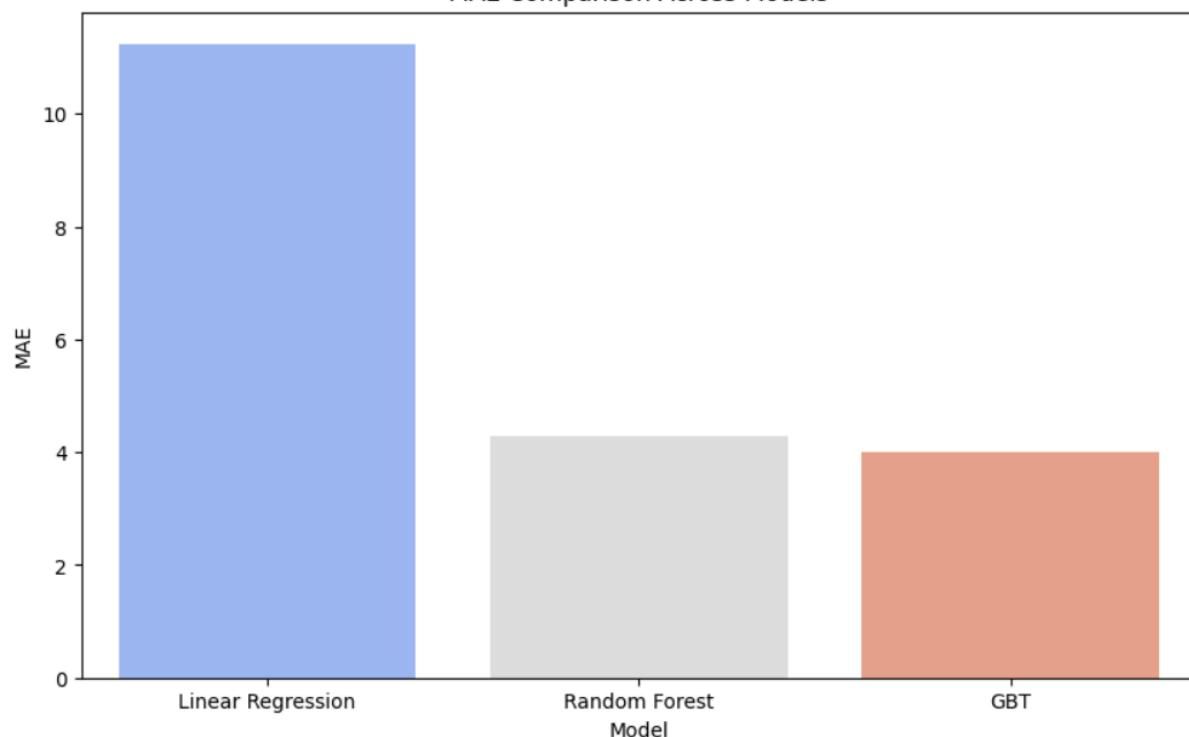
plt.figure(figsize=(10, 6))
sns.barplot(x='Model', y='RMSE', data=metrics_data, palette='coolwarm')
plt.title('RMSE Comparison Across Models')
plt.show()

plt.figure(figsize=(10, 6))
sns.barplot(x='Model', y='MAE', data=metrics_data, palette='coolwarm')
plt.title('MAE Comparison Across Models')
plt.show()
```

RMSE Comparison Across Models



MAE Comparison Across Models



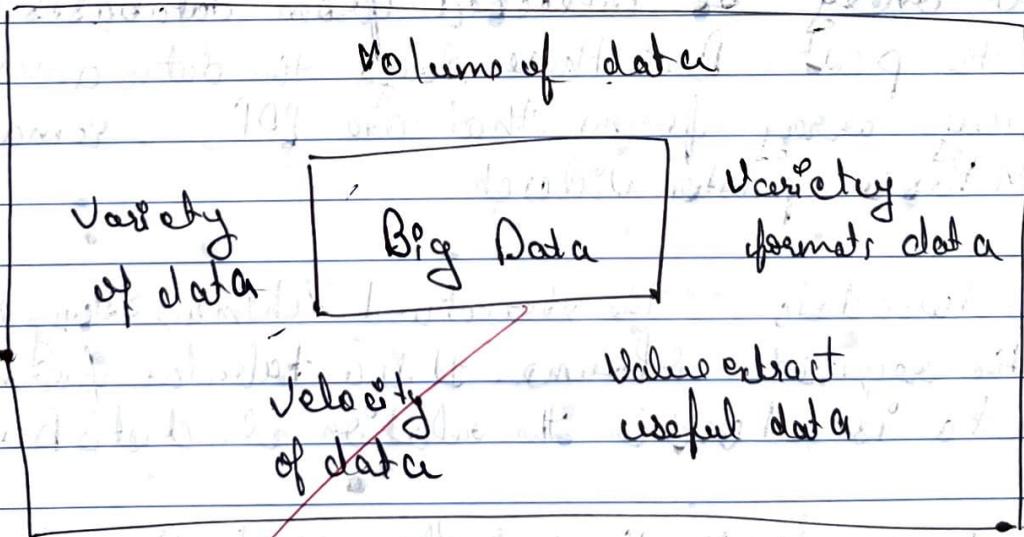
Assignment No. 1

Q) Write briefly on Big data characteristics Hadoop architecture Hadoop ecosystem with major components.

⇒ Big data Characteristics

Big data contain a large amount of data that is not being processed by traditional data storage or the processing unit. It is used by many mathematical companies to process the data and business of many organization the data flow could be exceed 150 Petabytes per bytes before replication.

there are five V's of Big Data that explain the characteristics.



the 5 V's of Big data can be follows:

1) Volume

The name of Big data itself is related to an enormous size. Big data is vast volume of data generated from many sources daily such as business process, machine, social media platform, networks, human interactions and many more.

Eg:- Facebook can generate approximately a billion message 4.5 B thus that the 'like' button is found & more than 350 million new post uploaded each day.

2) Variety

Big data can be structured, unstructured & semi-structured. That is being collected from different sources. Data will only be collected from databases & sheets in the past. But these days the data comes in the array form that are PDF's, emails, audio, SM, Passes, photos, videos etc.

3) Structured data :- The structured schema, along with all the required columns. It is a tabular form. Structured data is stored in the relational database.

4) Semi-structured :- The semi-structured the schema is not appropriately defined eg: JSON, XML, SQL. Online transaction processing system are built to work with semi-structured data. It is stored in real-time i.e. table.

Unstructured data : all the unstructured files like log files, audio files & image files are included in the unstructured data. Some organization have such data available but they did not know how to derive the value of the data since the data is raw.

Quasi structured data : the data formats contains data with inconsistent data formats that are formatted with that and time with some tools.

Eg: Web server logs is the log file is created & maintained by some servers that contain a list of activities.

Veracity

Veracity means how much the data is reliable. If has many ways filter or translate the data. Veracity is the process of being able to handle & manage data efficiently. Big data is also essential in business development.

Eg:- facebook posts with hashtags

Value :- P. Arthur et al. (2011) say that value is giving the data a meaning. Value is an essential characteristic of Big data. It is not the data that we process or store. It is "valuable and reliable" data that we store, process & also analyze.

5) Velocity:

Velocity plays an important role compared to others. Velocity creates the speed by which the data is fetched in real-time. It contains the latency at incoming data sets speeds, rate of damage, activity bursts, the primary aspects of Big data in the provide demanding data rapidly.

Big data velocity deals with the speed at the data flow from sources like application, logs, Business process network, social media sites, sensor mobile devices, etc.

Hadoop Architecture:

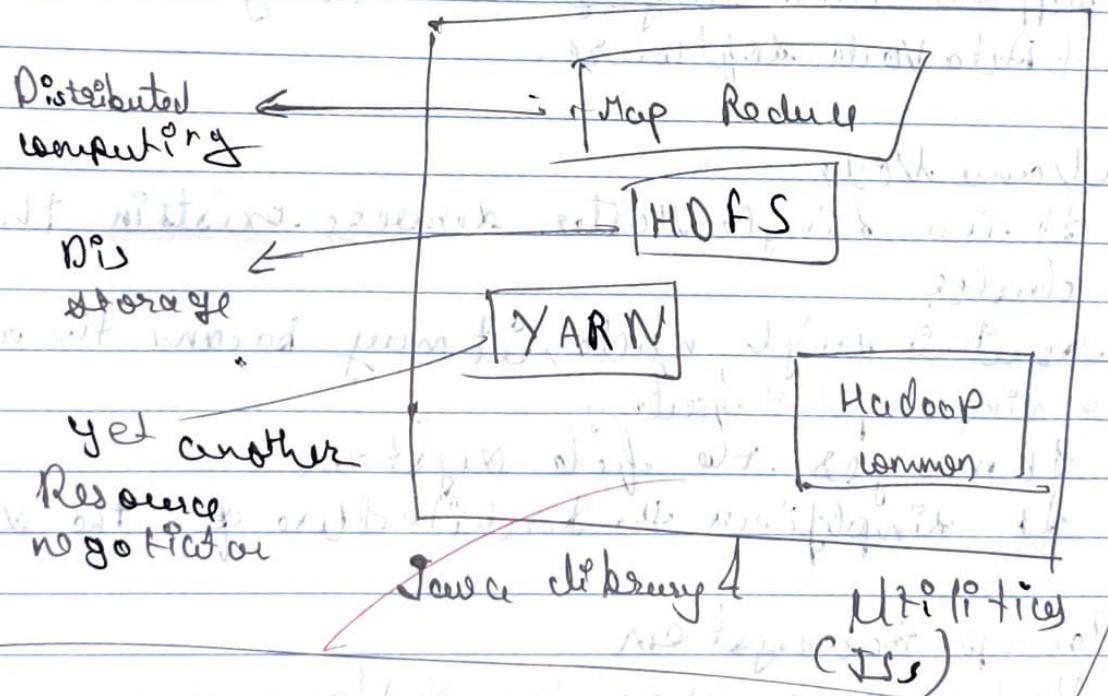
Hadoop is an open source framework from apache and is used to store process & analyse the data which are very huge in volumes. Hadoop is written in Java & is not OLAP or online analytical processing. It is used for batch processing. It is used by Facebook, Yahoo, Google.

Modules of Hadoop

i) HDFS:- Hadoop distributed file system (HDFS) published its paper by FS & on the basic that HDFS was developed. No. of trees that the files will be broken into blocks & it is stored in nodes. The distributed architecture

- 2) Yarn :- Yet another Resource Negotiator is used for job scheduling & manage the cluster
- 3) Map Reduce :- this is a framework which helps with java programs to do in the parallel computation on data using key-value pair, the Map task takes input data and converts it into cluster put Map task is consumed by reduce task & then the out of reduces gives final derived results

4) Hadoop Common - these Java developer are used to start Hadoop and can be used by other Hadoop modules



Hadoop architecture is a package of the file system. Map Reduce engine & the HDFS (Hadoop Distributed file system) & Hadoop cluster consists of a single Master & multiple slave Nodes. The Master Node - NameNode & Job tracker, task tracker, Nodes.

Hadoop distributed file system (HDFS) is a distributed file system for Hadoop. It contains a master/slave architecture. This Java language is used to develop HDFS. Any machine that supports Java language can easily run the Map & Reduce software.

Name Node

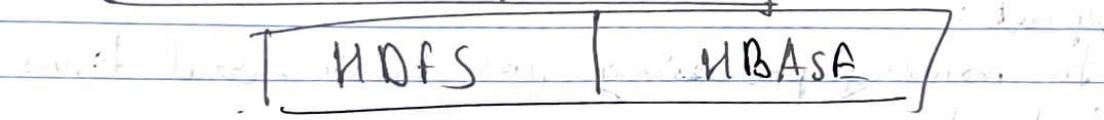
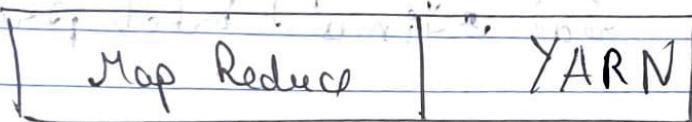
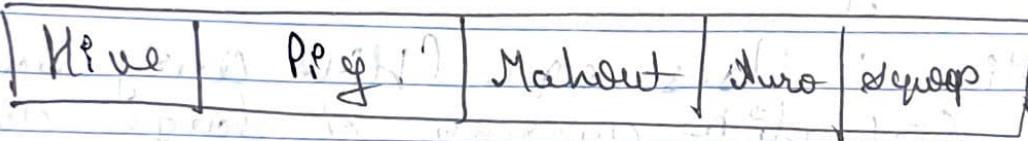
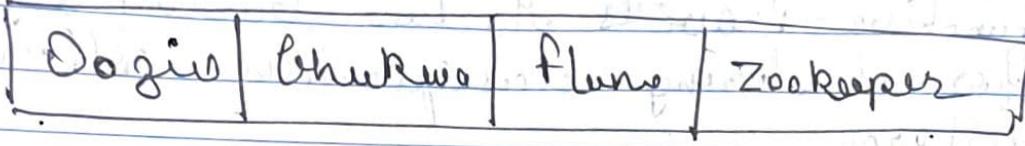
- 1] It is a single Master server exists in the HDFS cluster
- 2] As it is single node, it may become the root of single point failure.
- 3] It manages the file system
- 4] It simplifies the architecture of the system

Hadoop Ecosystem

Hadoop Ecosystem is a platform or a suit which provides various services to solve big data problems. It includes apache projects and various components tools & structures.

Components of Hadoop

- 1] HDFS : file system
- 2] Map Reduce :- Data Processing
- 3] Yarn :- Yet another Resource negotiation and scheduling
- 4] Spark - In Memory data processing
- 5] Pig Latin :- Query Based
- 6] HBase - No-SQL Database



1] HDFS: Components Name Node (meta data) & Data (Actual Node). Replication function stores large datasets across multiple nodes, ensuring fault tolerance.

2] YARN
Compences Resource manager in Node, Hadoop : App manager

function : Process large datasets using distributed parallel algorithms

- 3] MapReduce
Components - Map Node & Reduce
Aggregates & summarize data
- 4] Pig :- Uses Pig Latin (SQL like language) to process & analyse data. Excludes commands, MapReduce & stores results in HDFS
- 5] Apache Hive Machine learning function - Provides Machine learning algorithms for clustering, classification, collaborative filtering.
- 6] Hive :- Uses SQL (Hive Query Language) for SQL like querying of large datasets supports real time & batch processing
- 7] Spark :-
In memory processing for batch, real time & interact
- 8] HBase :- NoSQL database for real time & interactive data. Efficient for quick lookups in large data types.
- 9] Zookeeper :- Coordinates & synchronizes distributed system, ensuring consistency
- 10] Oozie :- Workflow scheduler: Manages job execution in sequential & event based coordination of job
- A History

Assignment 2

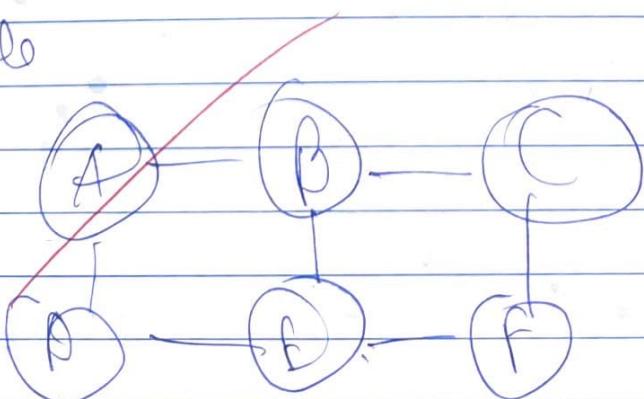
1)

→ In Girvan - Newman algorithm we are giving to divide the nodes of the graph into 2 or more clusters. This algorithm removes the edges with the highest between until there is no edge remain or any specific no. of edges is present. Between is the number of the shortest paths between pairs of nodes that are run through it.

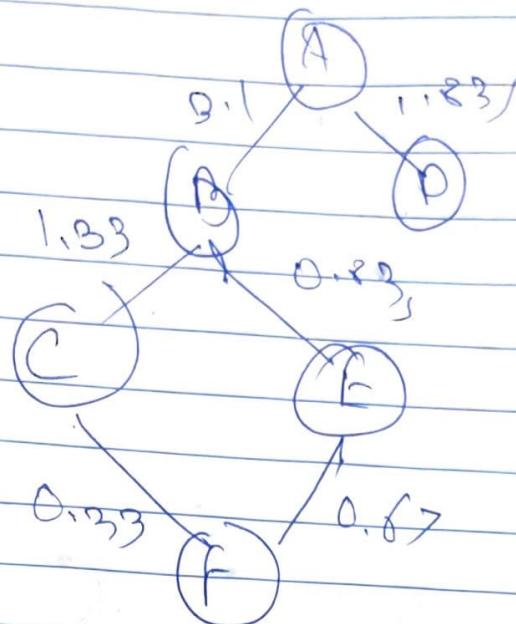
Algorithm.

- Create a graph of N nodes & its edges or take an inbuilt graph like farness graph
- calculate the betweenness of all existing graphs in the graph.
- Now remove all the edges with the highest betweenness.
- Recalculate the betweenness of all the edges that got affected by the removal of edges.
- Now repeat the step 3 & 4 until no edge remain.

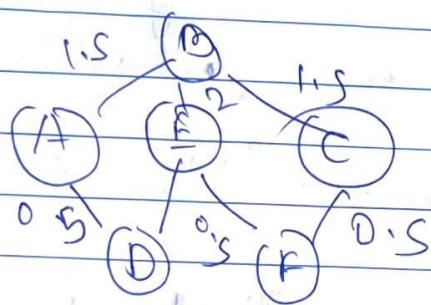
Example



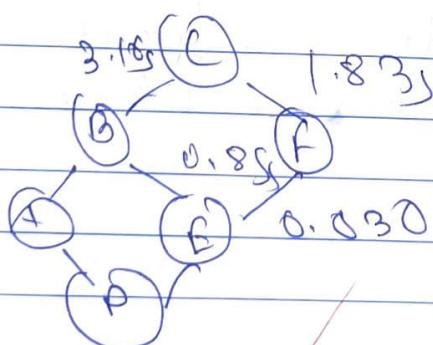
Starting with A



for B

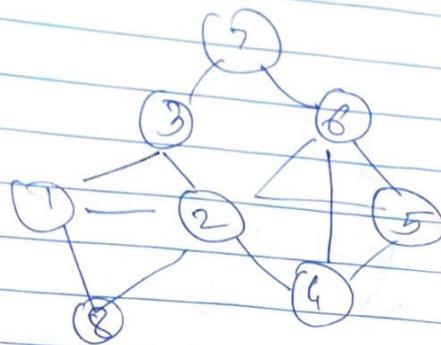


for C



2) Elbow percolation method builds up the communities from k-cliques which correspond to - compare subgraphs of R-nodes

Ex:



There are size $k = 3$ cliques

a : 1, 2, 3 \rightarrow Community Q

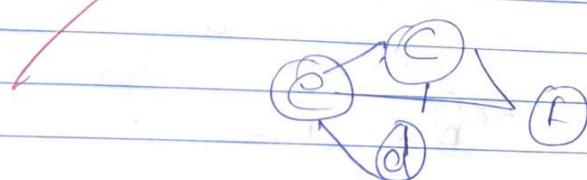
b : 1, 2, 8

c : 2, 4, 5 \rightarrow Community I

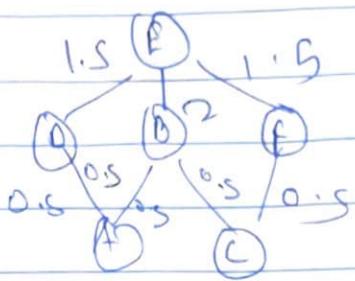
d : 2, 4, 6 \rightarrow Community b

e : 2, 5, 6 \rightarrow Community J

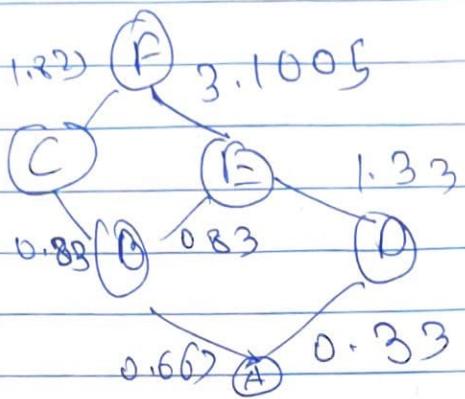
f : 4, 5, 6



for E



for F



Edges

A B

B D

B C

B E

C F

D F

E F

Edge Between
 E

5.33

8

7.34

5.33

8

8

