

**B.M.S. COLLEGE OF ENGINEERING**  
**(Autonomous College under VTU)**  
**Bull Temple Road, Basavanagudi, Bangalore - 560019**  
**BIG DATA ANALYTICS (23DS5PCBDA)**



**Alternative Assessment Report**  
**On**  
***“Movie Recommendation System”***

***Submitted By:***

Amrit Raj	1BM22AD005
Snehasish Kabi	1BM22AD057

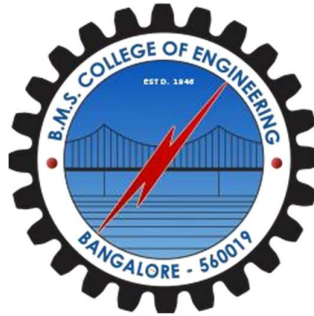
***in partial fulfillment for the award of the degree of***

**BACHELOR OF ENGINEERING**  
***in***  
**ARTIFICIAL INTELLIGENCE AND DATA SCIENCE**

**B.M.S COLLEGE OF ENGINEERING**

**P.O. Box No: 1908, Bull Temple Road, Bengaluru - 560019.**

**DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND DATA  
SCIENCE**



**CERTIFICATE**

This is certified to be the bonafide work of AMRIT RAJ(1BM22AD005), SNEHASISH KABI (1BM22AD057) in the Big Data Analytics course at B.M.S College of Engineering in partial fulfilment of the requirements for the fourth semester degree in Artificial Intelligence and Data Science under Visvesvaraya Technological University, Belgaum during academic year 2024-25.

**Signature of Faculty**

**Dr. Lakshmi Shree K,**

**Asst. Prof. Dept of AI & DSc.**

**AI&DSc.**

**Signature of HOD,**

**Dr. Indiramma M,**

**Prof. & Head, Dept. of**

## TABLE OF CONTENTS

Sl. no	Contents	Pg. No
1	Hadoop Assignment Report	1 - 5
2	MongoDB Assignment Report	6 - 8
3	Advanced MongoDB Assignment Report	9 - 10
4	Cassandra Assignment Report	11 - 15
5	Apache Kafka Report	16 - 17
6	Apache Spark Project Report	18 - 27
7	Conclusion	28



## **CHAPTER 1**

### **REPORT ON ASSIGNMENT 1: HADOOP**

#### **INTRODUCTION**

- Apache Hadoop is a collection of open-source software utilities for reliable, scalable, distributed computing. It provides a software framework for distributed storage and processing of big data using the MapReduce programming model.
- The core of Apache Hadoop consists of a storage part, known as Hadoop Distributed File System (HDFS), and a processing part which is a MapReduce programming model.
- Hadoop splits files into large blocks and distributes them across nodes in a cluster.
- It then transfers packaged code into nodes to process the data in parallel.

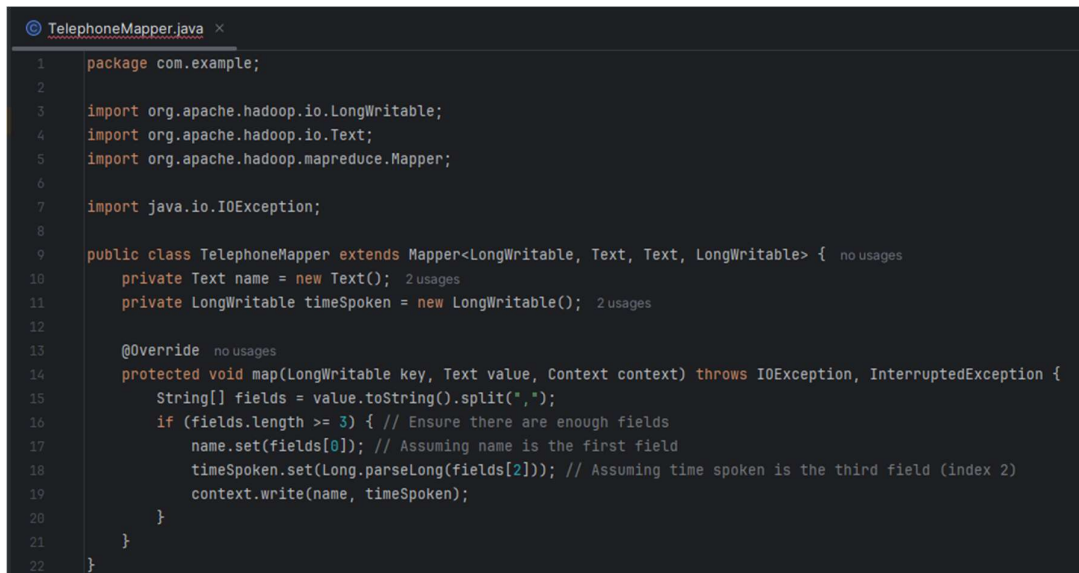
#### **PROBLEM STATEMENT**

- The task is to develop a MapReduce program that processes a CSV file containing telephone call records to calculate three metrics: the maximum call duration, the minimum call duration, and the users whose total call time exceeds a predefined threshold.
- The program should map call records to extract relevant data, then reduce to compute the max, min, and filter users based on the threshold.
- The input data includes user IDs, telephone numbers and call durations, and the output should show the computed metrics.

#### **EXECUTION**

- The input CSV file is loaded and split into smaller chunks for parallel processing in Hadoop.
- Mappers extract UserID and CallDuration, emitting key-value pairs for max/min calculations and user call time aggregation.
- Hadoop automatically sorts and groups the key-value pairs, ensuring that records for the same user are processed together.
- Reducers calculate the maximum and minimum call durations and aggregate total call time per user, filtering by the threshold.
- The final results, including max/min durations and users exceeding the threshold, are written to the output.

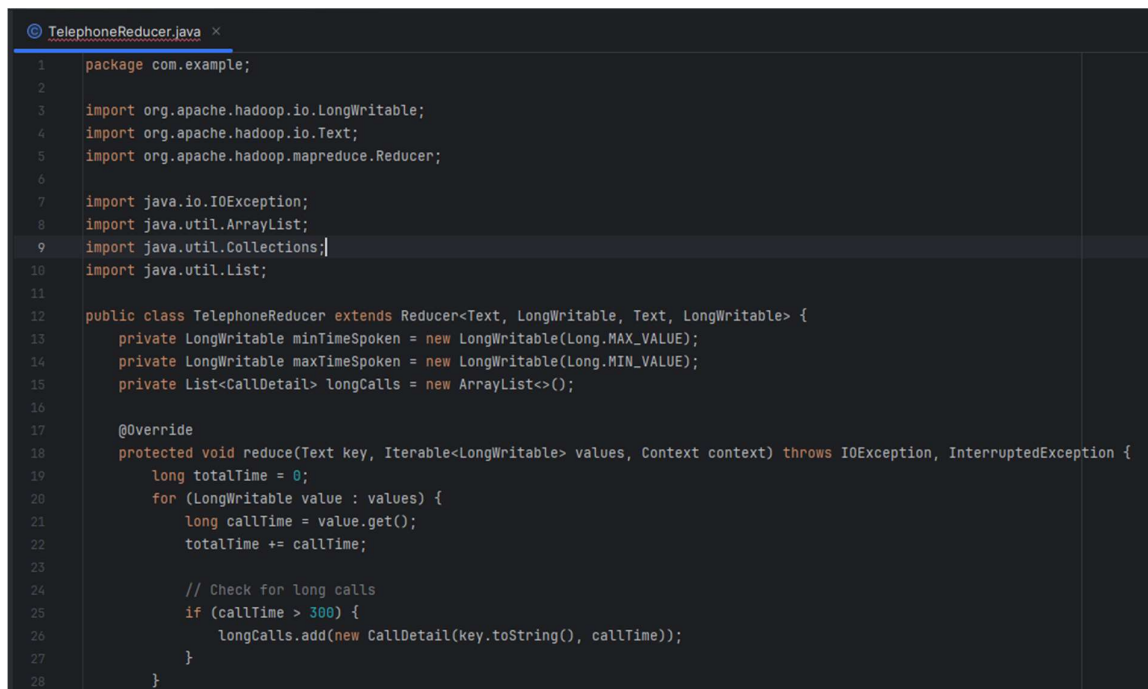
## SCREENSHOTS OF EXECUTION

A screenshot of a code editor showing the implementation of the TelephoneMapper class. The code is in Java and extends the Mapper class. It includes imports for Hadoop's LongWritable, Text, and Mapper classes, and Java's IOException. The map method splits the input string by commas and sets the name and timeSpoken fields based on the split results.

```
1 package com.example;
2
3 import org.apache.hadoop.io.LongWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Mapper;
6
7 import java.io.IOException;
8
9 public class TelephoneMapper extends Mapper<LongWritable, Text, Text, LongWritable> { no usages
10     private Text name = new Text(); 2 usages
11     private LongWritable timeSpoken = new LongWritable(); 2 usages
12
13     @Override no usages
14     protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
15         String[] fields = value.toString().split(",");
16         if (fields.length >= 3) { // Ensure there are enough fields
17             name.set(fields[0]); // Assuming name is the first field
18             timeSpoken.set(Long.parseLong(fields[2])); // Assuming time spoken is the third field (index 2)
19             context.write(name, timeSpoken);
20         }
21     }
22 }
```

**Fig. 1.1: Mapper Class File**

The Mapper class processes input data, splits it into key-value pairs, and performs the necessary transformations or filtering for further processing. It outputs intermediate key-value pairs that are passed to the reducer.

A screenshot of a code editor showing the first part of the TelephoneReducer class implementation. The code is in Java and extends the Reducer class. It includes imports for Hadoop's LongWritable, Text, and Reducer classes, and Java's IOException, ArrayList, and List classes. The reduce method iterates over the values, calculates the total time, and adds call details to a list if the call time is greater than 300 seconds.

```
1 package com.example;
2
3 import org.apache.hadoop.io.LongWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Reducer;
6
7 import java.io.IOException;
8 import java.util.ArrayList;
9 import java.util.Collections;
10 import java.util.List;
11
12 public class TelephoneReducer extends Reducer<Text, LongWritable, Text, LongWritable> {
13     private LongWritable minTimeSpoken = new LongWritable(Long.MAX_VALUE);
14     private LongWritable maxTimeSpoken = new LongWritable(Long.MIN_VALUE);
15     private List<CallDetail> longCalls = new ArrayList<>();
16
17     @Override
18     protected void reduce(Text key, Iterable<LongWritable> values, Context context) throws IOException, InterruptedException {
19         long totalTime = 0;
20         for (LongWritable value : values) {
21             long callTime = value.get();
22             totalTime += callTime;
23
24             // Check for long calls
25             if (callTime > 300) {
26                 longCalls.add(new CallDetail(key.toString(), callTime));
27             }
28         }
29     }
30 }
```

**Fig. 1.2: Reducer Class File part 1**

```

TelephoneReducer.java x
29
30     // Update min and max times based on total time for this key
31     if (totalTime < minTimeSpoken.get()) {
32         minTimeSpoken.set(totalTime);
33     }
34     if (totalTime > maxTimeSpoken.get()) {
35         maxTimeSpoken.set(totalTime);
36     }
37 }
38
39 @Override
40 protected void cleanup(Context context) throws IOException, InterruptedException {
41     // Output the minimum time spoken
42     context.write(new Text("Minimum Time Spoken:"), minTimeSpoken);
43     // Output the maximum time spoken
44     context.write(new Text("Maximum Time Spoken:"), maxTimeSpoken);
45
46     // Sort long calls by time spoken
47     Collections.sort(longCalls);
48
49     // Output the sorted long calls
50     for (CallDetail call : longCalls) {
51         context.write(new Text("Long Call: " + call.getName()), new LongWritable(call.getTime()));
52     }
53 }
54 }

```

**Fig. 1.3: Reducer Class File part 2**

The Reducer class receives grouped key-value pairs, aggregates the data, and performs computations like summing, finding maximum/minimum values, or filtering based on conditions before emitting the final result.

```

TelephoneDriver.java x
1  package com.example;
2
3  import org.apache.hadoop.conf.Configuration;
4  import org.apache.hadoop.fs.Path;
5  import org.apache.hadoop.io.LongWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Job;
8  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
9  import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
10
11 public class TelephoneDriver {
12     public static void main(String[] args) throws Exception {
13         // Check for the correct number of arguments
14         if (args.length != 2) {
15             System.err.println("Usage: TelephoneDriver <input path> <output path>");
16             System.exit(status: -1);
17         }
18
19         // Create a new configuration
20         Configuration conf = new Configuration();
21
22         // Create a new job
23         Job job = Job.getInstance(conf, "Telephone Data Analysis");
24         job.setJarByClass(TelephoneDriver.class);
25

```

**Fig. 1.4: Driver Class File part 1**

```

TelephoneDriver.java x
23 job = Job.getInstance(conf, "Telephone Data Analysis");
24 job.setJarByClass(TelephoneDriver.class);
25
26 // Set the Mapper and Reducer classes
27 job.setMapperClass(TelephoneMapper.class);
28 job.setReducerClass(TelephoneReducer.class);
29
30 // Set the output key and value types
31 job.setOutputKeyClass(Text.class);
32 job.setOutputValueClass(LongWritable.class);
33
34 // Set the input and output formats
35 job.setInputFormatClass(org.apache.hadoop.mapreduce.lib.input.TextInputFormat.class);
36 job.setOutputFormatClass(TextOutputFormat.class);
37
38 // Set the input and output paths
39 FileInputFormat.addInputPath(job, new Path(args[0]));
40 TextOutputFormat.setOutputPath(job, new Path(args[1]));
41
42 // Wait for the job to complete and exit
43 System.exit(job.waitForCompletion(true) ? 0 : 1);
44 }
45 }

```

**Fig. 1.5: Driver Class File part 2**

The Driver class configures and controls the MapReduce job, setting up input and output paths, specifying mapper and reducer classes, and initiating the job execution in the Hadoop environment.

```

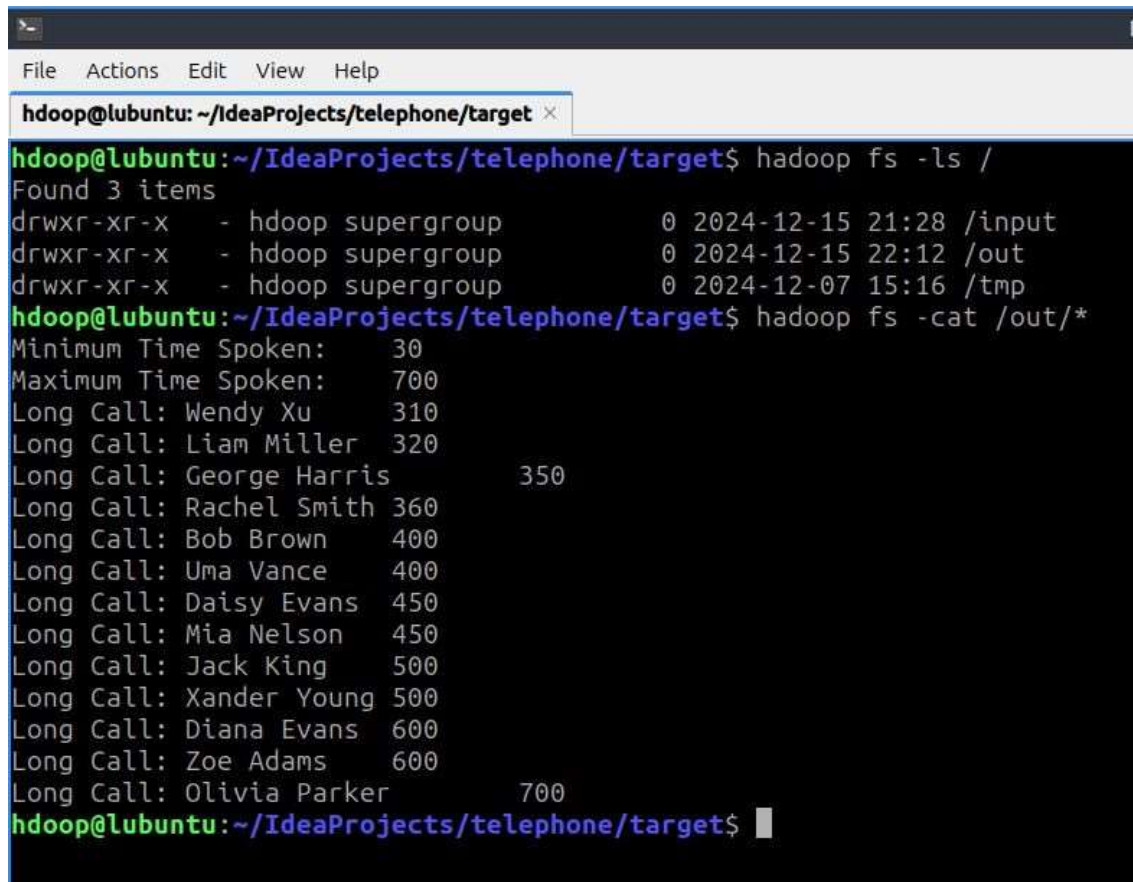
John Doe,123-456-7890,120
Jane Smith,234-567-8901,300
Alice Johnson,345-678-9012,150
Bob Brown,456-789-0123,400
Charlie Davis,567-890-1234,200
Diana Evans,678-901-2345,600
Ethan Foster,789-012-3456,90
Fiona Green,890-123-4567,250
George Harris,901-234-5678,350
Hannah Ivers,012-345-6789,180
Ian Johnson,123-456-7890,240
Jack King,234-567-8901,500
Kathy Lee,345-678-9012,60
Liam Miller,456-789-0123,320
Mia Nelson,567-890-1234,450
Noah O'Brien,678-901-2345,30
Olivia Parker,789-012-3456,700
Paul Quinn,890-123-4567,110
Quinn Roberts,901-234-5678,280
Rachel Smith,012-345-6789,360
Sam Taylor,123-456-7890,90
Tina Underwood,234-567-8901,150
Uma Vance,345-678-9012,400
Victor White,456-789-0123,220
Wendy Xu,567-890-1234,310
Xander Young,678-901-2345,500
Yara Zane,789-012-3456,80
Zoe Adams,890-123-4567,600
Aaron Brown,901-234-5678,150
Bella Clark,012-345-6789,90
Cody Davis,123-456-7890,300
Daisy Evans,234-567-8901,450

```

**Fig. 1.6: View of the Dataset**

The dataset contains the following: User-Name, User-Number and Time-Spoken as its attributes. Totally it contains three attributes and thirty five rows.





```
File Actions Edit View Help
hadoop@ubuntu: ~/IdeaProjects/telephone/target x
hadoop@ubuntu:~/IdeaProjects/telephone/target$ hadoop fs -ls /
Found 3 items
drwxr-xr-x - hadoop supergroup 0 2024-12-15 21:28 /input
drwxr-xr-x - hadoop supergroup 0 2024-12-15 22:12 /out
drwxr-xr-x - hadoop supergroup 0 2024-12-07 15:16 /tmp
hadoop@ubuntu:~/IdeaProjects/telephone/target$ hadoop fs -cat /out/*
Minimum Time Spoken: 30
Maximum Time Spoken: 700
Long Call: Wendy Xu 310
Long Call: Liam Miller 320
Long Call: George Harris 350
Long Call: Rachel Smith 360
Long Call: Bob Brown 400
Long Call: Uma Vance 400
Long Call: Daisy Evans 450
Long Call: Mia Nelson 450
Long Call: Jack King 500
Long Call: Xander Young 500
Long Call: Diana Evans 600
Long Call: Zoe Adams 600
Long Call: Olivia Parker 700
hadoop@ubuntu:~/IdeaProjects/telephone/target$
```

**Fig. 1.7: Output of the Hadoop Program**

Represents the output file on running Hadoop's mapreduce framework with the following codes on the above dataset.

## **CHAPTER 2**

### **REVIEW OF ASSIGNMENT 2: MONGODB**

#### **INTRODUCTION**

- MongoDB is a source-available, cross-platform, document-oriented database program. Classified as a NoSQL database product, MongoDB utilizes JSON-like documents with optional schemas.
- MongoDB supports field, range query and regular-expression searches. Queries can return specific fields of documents and also include user-defined JavaScript functions.
- Fields in a MongoDB document can be indexed with primary and secondary indices.
- MongoDB can be used as a file system, called GridFS, with load-balancing and data-replication features over multiple machines for storing files. MongoDB scales horizontally using sharding.

#### **PROBLEM STATEMENT**

- Implement a functionality to insert new documents into a MongoDB collection, allowing users to add records.
- Implement functionality to query and retrieve documents based on specific criteria.
- Provide options to update existing documents and delete while ensuring data integrity.

#### **EXECUTION**

- Establish a connection to the MongoDB database using a MongoDB client, ensuring the necessary database and collection are accessible.
- Insert new documents into the specified MongoDB collection using the insertOne() or insertMany() method, providing required data fields such as UserID, Name, and Age.
- Use find() or findOne() queries to fetch documents based on conditions, such as searching by UserID or retrieving users whose Age exceeds a given value.
- Perform updates on existing documents using the updateOne() or updateMany() method, modifying fields like Name or Age based on a specified filter (e.g., UserID).
- Use the deleteOne() or deleteMany() methods to remove documents from the collection based on conditions like UserID or other criteria, ensuring successful deletion through result validation.

## SCREENSHOTS OF EXECUTION

```
test> use CompanyA
switched to db CompanyA
CompanyA> db.createCollection('employee');
{ ok: 1 }
CompanyA> db.employee.insertOne({ name: "John Doe", age: 30, email: "john.doe@example.com" })
{
  acknowledged: true,
  insertedId: ObjectId('675effb4567a4f44522710bc')
}
CompanyA> db.employee.insertMany([
  {name: 'Alice', age: 25, email: 'alice@example.com'},
  ... {name: 'Bob', age: 28, email: 'bob@example.com'},
  ... {name: 'Cathie', age: 32, email: 'cathie@example.com'}
])
{
  acknowledged: true,
  insertedIds: {
    '0': ObjectId('675f004b567a4f44522710bd'),
    '1': ObjectId('675f004b567a4f44522710be'),
    '2': ObjectId('675f004b567a4f44522710bf')
  }
}
```

Fig. 2.1: Creation of collection and Inserting values

Here we create a database called **Company A**. Inside it, we create a collection called **employee** and insert values to it using **insertOne()** and **insertMany()** operation.

```
CompanyA> db.employee.find()
[
  {
    _id: ObjectId('675effb4567a4f44522710bc'),
    name: 'John Doe',
    age: 30,
    email: 'john.doe@example.com'
  },
  {
    _id: ObjectId('675f004b567a4f44522710bd'),
    name: 'Alice',
    age: 25,
    email: 'alice@example.com'
  },
  {
    _id: ObjectId('675f004b567a4f44522710be'),
    name: 'Bob',
    age: 28,
    email: 'bob@example.com'
  },
  {
    _id: ObjectId('675f004b567a4f44522710bf'),
    name: 'Cathie',
    age: 32,
    email: 'cathie@example.com'
  }
]

CompanyA> db.employee.find({ age: { $gte: 28 } })
[
  {
    _id: ObjectId('675effb4567a4f44522710bc'),
    name: 'John Doe',
    age: 30,
    email: 'john.doe@example.com'
  },
  {
    _id: ObjectId('675f004b567a4f44522710be'),
    name: 'Bob',
    age: 28,
    email: 'bob@example.com'
  },
  {
    _id: ObjectId('675f004b567a4f44522710bf'),
    name: 'Cathie',
    age: 32,
    email: 'cathie@example.com'
  }
]

CompanyA> db.employee.findOne({ name: "Alice" })
{
  _id: ObjectId('675f004b567a4f44522710bd'),
  name: 'Alice',
  age: 25,
  email: 'alice@example.com'
}
```

Fig. 2.2: Read Operation with queries

Here we display the contents of the collection using **find()** and **findOne()** operation. We can apply queries to get our specific and desired documents by specifying them inside the **find()** function.

```

CompanyA> db.employee.find({ age: { $gte: 25 } }, { name: 1, email: 1, _id: 0 })
[
  { name: 'John Doe', email: 'john.doe@example.com' },
  { name: 'Alice', email: 'alice@example.com' },
  { name: 'Bob', email: 'bob@example.com' },
  { name: 'Cathie', email: 'cathie@example.com' }
]
CompanyA> db.employee.updateOne( { name: "Bob" }, { $set: { age: 31 } } );
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 0,
  upsertedCount: 0
}
CompanyA> db.employee.updateMany( { age: { $lt: 30 } }, { $inc: { age: 1 } } );
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
CompanyA> db.employee.replaceOne(
... { name: "Alice" }, { name: "Alice Smith", age: 26, email: "alice.smith@example.com" }
... )
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}

```

Fig. 2.3: Update Operation with Queries

Here we perform update operation on the **employee** database using `updateOne()`, `updateMany()` and `replaceOne()`, `replaceMany()` functions using appropriate queries to select the documents.

```

CompanyA> db.employee.deleteOne({ name: "Bob" })
{ acknowledged: true, deletedCount: 1 }
CompanyA> db.employee.deleteMany({ age: { $gte: 32 } })
{ acknowledged: true, deletedCount: 1 }
CompanyA> db.employee.find()
[
  {
    _id: ObjectId('675effb4567a4f44522710bc'),
    name: 'John Doe',
    age: 30,
    email: 'john.doe@example.com'
  },
  {
    _id: ObjectId('675f004b567a4f44522710bd'),
    name: 'Alice Smith',
    age: 26,
    email: 'alice.smith@example.com'
  }
]

```

Fig. 2.4: Delete Operation with Queries

Here we perform delete operation on the **employee** collection using `deleteOne()` and `deleteMany()` functions along with the appropriate queries to select the documents.

## CHAPTER 3

### REVIEW OF ASSIGNMENT 3: MONGODB

#### PROBLEM STATEMENT

- The task involves leveraging MongoDB to efficiently manage and analyze a shopping dataset through advanced database operations.
- Key objectives include retrieving records based on specific conditions, performing sorting and counting operations to organize and quantify data, and utilizing aggregation pipelines to derive meaningful insights.
- This assignment demonstrates the capabilities of MongoDB in handling large datasets while ensuring flexibility and scalability for real-world applications.

#### EXECUTION

- To carry out this assignment, a shopping dataset was imported into MongoDB as a collection. Various operations were performed to demonstrate the database's capabilities.
- Conditional queries were executed to retrieve specific records based on criteria, such as price range or product category.
- Sorting and counting operations were applied to organize data and quantify records based on conditions.
- Aggregation operations, such as grouping data by categories, union-like operations to combine datasets, and calculating totals or averages, were utilized to analyze the dataset effectively.

#### SCREENSHOTS OF EXECUTION

```
db.bda.find({
  $and: [
    { "Category": "Electronics" },
    { "Price": { $lt: 1000 } }
  ]
});
```

```
> db.bda.find({
  $and: [
    { "ProductID": 101 },
    { "Category": "Electronics" },
    { "Price": { $lt: 1000 } }
  ]
});
```

```
> db.bda.find({Category:{$ne:'Electronics'}});
```

**Fig. 3.1: Logical Operations on Dataset**

Here, we apply and logical operators to get products with our required characteristics along with the not equal to operator to exclude the unnecessary.

```
> db.bda.find({Supplier:"Apple Inc."}).sort({rating:-1});
> db.bda.countDocuments({ Rating: { $gt: 4.5 }, QuantityInStock: { $gt: 50 } })
< 21
```

**Fig. 3.2: Sorting Operations and Count Operations on the Documents**

Here, we apply the sort operation to get the document sorted. 1 for ascending order and -1 for descending order of sort. Also, with specific queries, we can count the required documents.

```
> db.bda.aggregate([
  {
    $match: {
      Category: "Electronics",
      Price: { $lte: 200.00 }
    }
  },
  {
    $group: {
      _id: "$Category",
      totalProducts: { $sum: 1 },
      averagePrice: { $avg: "$Price" }
    }
  }
]);
< {
  _id: 'Electronics',
  totalProducts: 22,
  averagePrice: 111.80818181818182
}
```

```
> db.bda.aggregate([
  { $unionWith: { coll: "electronics" } },
  { $group: { _id: "$Category", count: { $sum: 1 } } }
]);
< {
  _id: 'Laptops',
  count: 3
}
```

**Fig. 3.3: Grouping and Union Operation on the Documents**

Here, we group the documents based on a specified category and union operation to get the whole column along with some grouping operation.

## **CHAPTER 4**

### **REPORT OF ASSIGNMENT 4: CASSANDRA**

#### **INTRODUCTION**

- Apache Cassandra is a free and open-source database management system designed to handle large volumes of data across multiple commodity servers.
- The system prioritizes availability and scalability over consistency, making it particularly suited for systems with high write throughput requirements due to its LSM tree indexing storage layer.
- As a wide-column database, Cassandra supports flexible schemas and efficiently handles data models with numerous sparse columns.
- The system is optimized for applications with well-defined data access patterns that can be incorporated into the schema design.

#### **PROBLEM STATEMENT**

- Design a suitable data model for a student management system that captures essential information about students, including their unique student ID, name, age, and enrolled course.
- Inserting new student records with attributes such as student ID, name, age, and course.
- Retrieving student information based on student ID, listing all students enrolled in a specific course, and fetching details of students by age. Modifying existing student records to update their name, age, or course information.
- Removing student records from the database when they are no longer needed or when a student graduates.

#### **EXECUTION**

- Define a keyspace for your student management system. A keyspace is a namespace that defines how data is replicated across nodes.
- Create a table to store student records. The schema should include columns for student ID, name, age, and course.
- The Create operation in a student management system involves adding new student records to the database.
- The Read operation is used to retrieve student records from the database. This operation can take various forms depending on the requirements.
- The Update operation allows for the modification of existing student records.
- The Delete operation is used to remove student records from the database.

## SCREENSHOTS OF EXECUTION

```
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.2.0 | Cassandra 5.0.2 | CQL spec 3.4.7 | Native protocol v5]
Use HELP for help.
cqlsh>
cqlsh> CREATE KEYSPACE student_data
... WITH replication = {
...     'class': 'SimpleStrategy',
...     'replication_factor': 3
... };

Warnings :
Your replication factor 3 for keyspace student_data is higher than the number of nodes 1

cqlsh> |
```

Fig. 4.1: Connecting to Cassandra and creating Keyspace

We connect to the Cassandra cluster running in localhost using **cqlsh** command. Then we create a keyspace **student\_data** to initialize the database to perform further actions.

```
cqlsh> USE student_data;
cqlsh:student_data> CREATE TABLE students (
...     id UUID PRIMARY KEY,
...     name TEXT,
...     age INT,
...     course TEXT
... );
cqlsh:student_data> INSERT INTO students (id, name, age, course)
... VALUES (uuid(), 'Alice', 20, 'Computer Science');
cqlsh:student_data> INSERT INTO students (id, name, age, course)
... VALUES (uuid(), 'Bob', 22, 'Mathematics');
cqlsh:student_data> INSERT INTO students (id, name, age, course)
... VALUES (uuid(), 'Charlie', 21, 'Physics');
cqlsh:student_data> BEGIN BATCH
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'David', 23, 'Biology');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Eve', 19, 'History');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Frank', 22, 'Chemistry');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Grace', 20, 'English Literature');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Hank', 21, 'Mechanical Engineering');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Ivy', 19, 'Philosophy');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Jack', 24, 'Civil Engineering');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Karen', 22, 'Psychology');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Liam', 23, 'Political Science');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Mia', 21, 'Fine Arts');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Noah', 20, 'Astronomy');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Olivia', 19, 'Economics');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Paul', 22, 'Law');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Quinn', 20, 'Environmental Science');
...     INSERT INTO students (id, name, age, course) VALUES (uuid(), 'Rose', 23, 'Sociology');
... APPLY BATCH;
cqlsh:student_data> |
```

Fig. 4.2: Creating Table and Inserting values in to it

Here we are creating a table **students** to store the values and insert values into it using the cql commands and batch insertion.



```
cqlsh:student_data> select * from students ;
```

id	age	course	name
2577bf95-eb84-4c02-9391-d632f4cb2e9c	21	Mechanical Engineering	Hank
0f4dd09a-f840-4ba8-9666-3822b10ddf9a	21	Fine Arts	Mia
7f480459-5d79-453d-ae0e-7a32106b88ef	20	English Literature	Grace
719bb9ea-6ade-4a08-a570-4377b534dc89	19	History	Eve
1510e6bb-336e-4277-8992-ae4c25f7af0b	23	Biology	David
89f76659-f5ce-4ec2-a24d-5b2c106cd9c1	19	Economics	Olivia
0e88d5a2-2a1e-4764-9a42-afddd10fc32b	23	Sociology	Rose
9e47e076-1f0d-4a69-b972-4b5b7b64eb1f	22	Psychology	Karen
004d03ea-c840-4a31-a610-d421545cfd67	22	Mathematics	Bob
52a0dcab-06bc-4ee5-befd-82c717e1be73	20	Astronomy	Noah
15f53c02-ff24-4225-8488-2bd8f95982ea	22	Chemistry	Frank
1498e749-c70d-4afd-a5d5-33c098d1590f	24	Civil Engineering	Jack
c350fb09-56a3-4da4-988c-df58a921abe7	23	Political Science	Liam
5d4fc6ec-203d-4adb-b81c-4c51453904d0	20	Environmental Science	Quinn
2a54a1b3-d02c-479f-a585-2b42154f5793	22	Law	Paul
9b3c283d-9262-4ff5-b06b-30ffb747a628	20	Computer Science	Alice
378f0cd7-6a27-4cd0-b146-d50344afee7a	19	Philosophy	Ivy
8af2f0ad-20aa-46de-9942-b3ab1c72bb41	21	Physics	Charlie

(18 rows)

```
cqlsh:student_data> |
```

```
cqlsh:student_data> SELECT id, name, age, credits FROM students_by_course WHERE course = 'Computer Science' AND age >= 20 AND age <= 22;
InvalidRequest: Error from server: code=2200 [Invalid query] message="Cannot execute this query as it might involve data filtering and thus may have unpredictable performance. If you want to execute this query, use the ALLOW FILTERING clause"
```

```
cqlsh:student_data> SELECT id, name, age, credits FROM students_by_course WHERE course = 'Computer Science' AND age >= 20 AND age <= 22 ALLOW FILTERING;
```

id	name	age	credits
9b3c283d-9262-4ff5-b06b-30ffb747a628	Alice	20	21

(1 rows)

```
cqlsh:student_data> SELECT id, name, age, credits FROM students_by_course WHERE age >= 20 AND age <= 22 ALLOW FILTERING;
```

id	name	age	credits
2a54a1b3-d02c-479f-a585-2b42154f5793	Paul	22	19
2577bf95-eb84-4c02-9391-d632f4cb2e9c	Hank	21	20
8af2f0ad-20aa-46de-9942-b3ab1c72bb41	Charlie	21	23
15f53c02-ff24-4225-8488-2bd8f95982ea	Frank	22	17
004d03ea-c840-4a31-a610-d421545cfd67	Bob	22	15
52a0dcab-06bc-4ee5-befd-82c717e1be73	Noah	20	25
9b3c283d-9262-4ff5-b06b-30ffb747a628	Alice	20	21
0f4dd09a-f840-4ba8-9666-3822b10ddf9a	Mia	21	18
5d4fc6ec-203d-4adb-b81c-4c51453904d0	Quinn	20	22
7f480459-5d79-453d-ae0e-7a32106b88ef	Grace	20	22
9e47e076-1f0d-4a69-b972-4b5b7b64eb1f	Karen	22	19

(11 rows)

```
cqlsh:student_data> |
```

```
cqlsh:student_data> SELECT * FROM students WHERE age = 22;
```

id	age	course	name
9e47e076-1f0d-4a69-b972-4b5b7b64eb1f	22	Psychology	Karen
004d03ea-c840-4a31-a610-d421545cfd67	22	Mathematics	Bob
15f53c02-ff24-4225-8488-2bd8f95982ea	22	Chemistry	Frank
2a54a1b3-d02c-479f-a585-2b42154f5793	22	Law	Paul

(4 rows)

```
cqlsh:student_data> |
```

```
cqlsh:students_by_course> SELECT * FROM students_by_course.student_ WHERE course = 'Physics' ORDER BY age ASC, credits DESC;
```

course	age	credits	id	name
Physics	20	25	fb25d868-aa7c-41a1-b5f2-a8b577ba87ce	Alice
Physics	20	19	a842b2f8-5246-4fdc-890a-0b8b1d1750a4	Grace
Physics	21	23	625270ee-0c76-4fca-b019-cb8ca2dd239b	Eve
Physics	21	18	fb1fac31-90c7-44f1-83ec-4700aef93ff1	Olivia
Physics	22	24	c76ddcf0-c7f8-47c0-8cb7-9eb9d301809e	Frank
Physics	23	22	c4b6a3c9-9dec-4e31-9435-74d99273162c	Noah

```
(6 rows)
cqlsh:students_by_course> |
```

```
cqlsh:student_data> SELECT COUNT(*) FROM students;
```

```
count
```

```
18
```

```
(1 rows)
```

```
Warnings :
```

```
Aggregation query used without partition key
```

```
cqlsh:student_data> SELECT MIN(age) FROM students;
```

```
system.min(age)
```

```
19
```

```
(1 rows)
```

```
Warnings :
```

```
Aggregation query used without partition key
```

```
cqlsh:student_data> SELECT MAX(age) FROM students;
```

```
system.max(age)
```

```
24
```

```
(1 rows)
```

```
Warnings :
```

```
Aggregation query used without partition key
```

**Fig. 4.3: Read and Advanced Filtering**

Here we perform multiple read operations and querying to get our desired results using cql language which is similar to sql language.

```

cqlsh:student_data> ALTER TABLE students DROP course_id;
cqlsh:student_data> DESCRIBE TABLE students;

CREATE TABLE student_data.students (
  id uuid PRIMARY KEY,
  age int,
  course text,
  name text
) WITH additional_write_policy = '99p'
AND allow_auto_snapshot = true
AND bloom_filter_fp_chance = 0.01
AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}
AND cdc = false
AND comment = ''
AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy', 'max_threshold': '32', 'min_threshold': '4'}
AND compression = {'chunk_length_in_kb': '16', 'class': 'org.apache.cassandra.io.compress.LZ4Compressor'}
AND memtable = 'default'
AND crc_check_chance = 1.0
AND default_time_to_live = 0
AND extensions = {}
AND gc_grace_seconds = 864000
AND incremental_backups = true
AND max_index_interval = 2048
AND memtable_flush_period_in_ms = 0
AND min_index_interval = 128
AND read_repair = 'BLOCKING'
AND speculative_retry = '99p';

CREATE INDEX age_index ON student_data.students (age);
cqlsh:student_data>

```

**Fig. 4.4: Updating of values**

Here we update the table by dropping a column and the same method can be used to add, remove columns, modify or update existing rows in the table.

```

cqlsh:student_data> select * from students;

```

id	age	course	credits	name
2577bf95-eb84-4c02-9391-d632f4cb2e9c	21	Mechanical Engineering	null	Hank
0f4dd09a-f840-4ba8-9666-3822b10ddf9a	21	Fine Arts	null	Mia
7f480459-5d79-453d-ae0e-7a32106b88ef	20	English Literature	null	Grace
719bb9ea-6ade-4a08-a570-4377b534dc89	19	History	null	Eve
1510e6bb-336e-4277-8992-ae4c25f7af0b	23	Biology	null	David

```

cqlsh:student_data> DELETE FROM students WHERE id = 2577bf95-eb84-4c02-9391-d632f4cb2e9c;
cqlsh:student_data> select * from students;

```

id	age	course	credits	name
0f4dd09a-f840-4ba8-9666-3822b10ddf9a	21	Fine Arts	null	Mia
7f480459-5d79-453d-ae0e-7a32106b88ef	20	English Literature	null	Grace
719bb9ea-6ade-4a08-a570-4377b534dc89	19	History	null	Eve
1510e6bb-336e-4277-8992-ae4c25f7af0b	23	Biology	null	David
89f76659-f5ce-4ec2-a24d-5b2c106cd9c1	19	Economics	null	Olivia
0e88d5a2-2a1e-4764-9a42-afddd10fc32b	23	Sociology	null	Rose

**Fig. 4.5: Deleting of Values**

Here, we delete a row of a table using the primary key. Similarly, we can perform delete operations on the table using our required queries written in cql.

## CHAPTER 5

# REPORT ON APACHE KAFKA

## INTRODUCTION

- Apache Kafka is a distributed event store and stream-processing platform. It is an open-source system developed by the Apache Software Foundation written in Java and Scala.
- The project aims to provide a unified, high-throughput, low-latency platform for handling real-time data feeds. Kafka can connect to external systems (for data import/export) via Kafka Connect, and provides the Kafka Streams libraries for stream processing applications.
- Kafka uses a binary TCP-based protocol that is optimized for efficiency and relies on a "message set" abstraction that naturally groups messages together to reduce the overhead of the network roundtrip.

## ARCHITECTURE

- Kafka stores key-value messages that come from arbitrarily many processes called producers. The data can be partitioned into different "partitions" within different "topics".
- Other processes called "consumers" can read messages from partitions. For stream processing, Kafka offers the Streams API that allows writing Java applications that consume data from Kafka and write results back to Kafka.
- Kafka runs on a cluster of one or more servers (called brokers), and the partitions of all topics are distributed across the cluster nodes.

## KAFKA APIs

- **Producer API** – Permits an application to publish streams of records.
- **Consumer API** – Permits an application to subscribe to topics and processes streams of records.
- **Connect API** – Executes the reusable producer and consumer APIs that can link the topics to the existing applications.
- **Streams API** – This API converts the input streams to output and produces the result.
- **Admin API** – Used to manage Kafka topics, brokers, and other Kafka objects.

## APPLICATIONS

- **Real-time Data Ingestion and Stream Processing:** Kafka efficiently ingests and processes real-time data streams, enabling low-latency analytics and event-driven decision-making.
- **Event-Driven Micro services Architecture:** Kafka facilitates asynchronous communication between micro services through event-driven architecture, enhancing scalability and decoupling.
- **Log Aggregation and Monitoring:** Kafka centralizes log and metric data from multiple sources, enabling real-time monitoring and analysis of system performance.
- **Data Integration and ETL Pipelines:** Kafka is used to build scalable and reliable ETL (Extract, Transform, Load) pipelines, enabling real-time data integration across different systems and applications.
- **Message Queuing and Asynchronous Messaging:** Kafka serves as a high-throughput, fault-tolerant messaging system for decoupling producers and consumers, ensuring reliable message delivery even in distributed environments.
- **Data Replication and Disaster Recovery:** Kafka provides data replication capabilities across multiple clusters, ensuring high availability and enabling disaster recovery for critical applications.

## CHAPTER 6

# REPORT OF THE PROJECT: PYSPARK

## INTRODUCTION

### PROBLEM STATEMENT

The primary objective of this project is to design and implement a **Hybrid Movie Recommendation System** that combines the strengths of both **collaborative filtering** and **content-based filtering** approaches.

Standalone techniques often face specific limitations:

- **Collaborative Filtering:**
  - Suffers from the **cold start problem**, where recommendations for new users or items with minimal interactions are difficult to generate.
  - Struggles with **sparsity issues** in datasets where user-item interactions are limited.
  -
- **Content-Based Filtering:**
  - Relies solely on item attributes, making it challenging to capture nuanced user preferences.
  - May lead to recommendations lacking diversity, as it often suggests items similar to those already interacted with.

The hybrid approach aims to address these challenges by:

1. Utilizing **collaborative filtering** to understand user preferences based on past interactions and similarities between users/items.
2. Enhancing recommendations with **content-based features** (e.g., movie genres and titles) to ensure relevance and diversity.

This system is designed to deliver accurate, scalable, and diverse recommendations while overcoming the limitations of individual techniques.

## DESCRIPTION OF THE DATASET

The project utilizes two key datasets: Movies Dataset and Ratings Dataset, both essential for building the Hybrid Movie Recommendation System. These datasets provide information about movies and user interactions, enabling the development of collaborative and content-based recommendation models.

### 1. Movies Dataset

- **Attributes/Fields:**
  - **movieId:** A unique identifier for each movie.
  - **title:** The name of the movie, often including the release year (e.g., "Toy Story (1995)").
  - **genres:** A pipe-separated string listing the movie genres (e.g., "Animation|Children|Comedy").
  -
- **Purpose:**
  - Provides metadata for building content-based features, such as genres and title-based TF-IDF scores.
- **Example Data:**

movieId	title	genres
1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	Jumanji (1995)	Adventure Children Fantasy
3	Grumpier Old Men (1995)	Comedy Romance
4	Waiting to Exhale (1995)	Comedy Drama Romance
5	Father of the Bride Part II (1995)	Comedy
6	Heat (1995)	Action Crime Thriller
7	Sabrina (1995)	Comedy Romance
8	Tom and Huck (1995)	Adventure Children
9	Sudden Death (1995)	Action
10	GoldenEye (1995)	Action Adventure Thriller
11	American President, The (1995)	Comedy Drama Romance
12	Dracula: Dead and Loving It (1995)	Comedy Horror
13	Balto (1995)	Adventure Animation Children
14	Nixon (1995)	Drama
15	Cutthroat Island (1995)	Action Adventure Romance
16	Casino (1995)	Crime Drama
17	Sense and Sensibility (1995)	Drama Romance
18	Four Rooms (1995)	Comedy
19	Ace Ventura: When Nature Calls	Comedy



## 2. Ratings Dataset

- **Attributes/Fields:**
  - **userId:** A unique identifier for each user.
  - **movieId:** A unique identifier for movies, linking it to the Movies Dataset.
  - **rating:** A numerical value representing the user's rating for the movie, typically on a scale from 0 to 5.
- **Purpose:**
  - Provides interaction data for building the collaborative filtering model.
- **Example Data:**

userId	movieId	rating
1	2	3.5
1	29	3.5
1	32	3.5
1	47	3.5
1	50	3.5
1	112	3.5
1	151	4
1	223	4
1	253	4
1	260	4
1	293	4
1	296	4
1	318	4
1	337	3.5
1	367	3.5

### Dataset Characteristics

- **Size:**
  - **Movies Dataset:** Approximately X rows.
  - **Ratings Dataset:** Approximately Y rows.
- **Source:**
  - Both datasets are sourced from the publicly available MovieLens Dataset.

### Integration and Usage

- **Movies Dataset:** Used to extract content-based features (e.g., genres and titles).



- **Ratings Dataset:** Used to train the collaborative filtering model by learning user preferences and movie ratings.
- **Combined:** The movieId serves as a key for linking the two datasets.

This structured combination of datasets ensures a robust foundation for the hybrid recommendation system, addressing diverse user preferences and movie attributes.

## METHODOLOGY

The methodology for building the Hybrid Movie Recommendation System addresses the problem statement by combining content-based filtering and collaborative filtering techniques to provide personalized recommendations. Below is the detailed explanation:

### Steps in the Methodology

#### 1. Data Preprocessing

- **Dataset Cleaning:** The movies and ratings datasets were cleaned to ensure correct data types and remove null or inconsistent values.
- **Content Processing:** The movie titles were tokenized and processed using TF-IDF to extract textual features, while genres were vectorized for further processing.
- **Collaborative Filtering Data Preparation:** The ratings dataset was processed to generate a user-movie interaction matrix.

#### 2. Feature Engineering

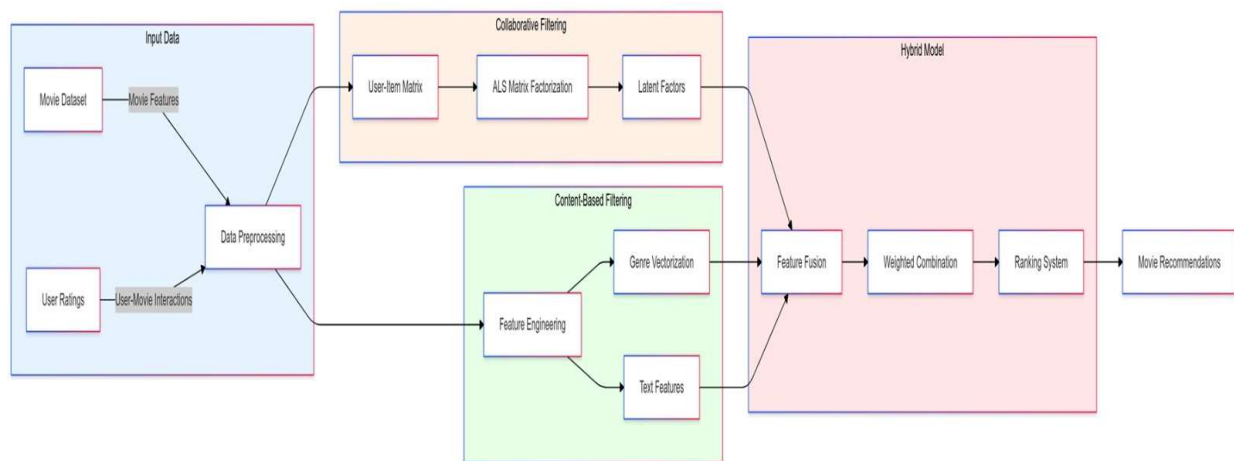
- **Content Features:**
  - Titles and genres were converted into numerical feature vectors using techniques like TF-IDF and CountVectorizer.
  - These feature vectors were normalized to ensure comparability during similarity computation.
- **Collaborative Filtering Features:**
  - The ALS (Alternating Least Squares) algorithm was used to generate latent features for users and items based on the user-movie rating matrix.

#### 3. Model Development

- **Collaborative Filtering:**
  - A collaborative model was trained using ALS to recommend movies based on patterns in user ratings.
- **Content-Based Filtering:**
  - Content similarity between movies was computed using cosine similarity on normalized feature vectors.
- **Hybrid System:**
  - The hybrid system combined collaborative filtering and content-based filtering scores using a weighted approach (e.g., 70% CF and 30% Content). This ensured recommendations were both personalized and contextually relevant.

#### 4. Evaluation

- Root Mean Square Error (RMSE) was used to evaluate the accuracy of the collaborative filtering model by comparing actual and predicted ratings.
- The final hybrid recommendations were compared against the individual content and collaborative filtering systems to assess improvement.



## Diagram for Methodology

Description of the Diagram:

- The diagram should illustrate the data flow through the system:
  1. Input datasets (movies and ratings).
  2. Two parallel pipelines: content-based feature extraction and collaborative filtering.
  3. A hybrid recommendation engine combining the results.
  4. Final recommendations displayed to the user.

## How the Methodology Addresses the Problem Statement

### 1. Hybrid Approach:

- Combines the strengths of content-based filtering (recommends similar movies) and collaborative filtering (recommends based on user behavior).
- Resolves the cold-start problem to some extent: Content-based filtering works well for new movies, and collaborative filtering works well for existing movies.

### 2. Personalization:

- Collaborative filtering ensures personalized recommendations by learning user preferences from historical data.

### 3. Scalability:

- The system leverages Apache Spark's distributed computing capabilities to handle large datasets efficiently.

### 4. Improved Accuracy:

- By combining the two approaches, the hybrid model achieves a better balance between precision and diversity of recommendations, addressing user needs effectively.

## Tools/Technology Used

This project employs a combination of tools and technologies to handle data processing, feature extraction, and model training effectively. Here's a detailed explanation of each technology used, referenced with relevant code snippets:

### 1. PySpark

**Purpose:** PySpark is the core framework used in this project for handling distributed data processing and machine learning tasks efficiently.

- **Data Loading and Preprocessing:**

- PySpark's `SparkSession` is used to initialize the Spark application.

- **Data Loading:**

- PySpark's `read.csv` function is used to load the datasets:

### 2. Feature Engineering Tools

PySpark's `ml.feature` module provides tools for extracting and transforming features from the data.

- **RegexTokenizer:**

- Splits the movie titles into individual tokens (words).
- **Example:**
- `tokenizer = RegexTokenizer(inputCol="title", outputCol="title_tokens", pattern="\\W")`
  - Why Used: Helps convert unstructured text data (titles) into structured tokens for further processing.

- **StopWordsRemover:**

- Removes common words (e.g., "the," "and") that do not add meaningful information.
- Example:
- `remover = StopWordsRemover(inputCol="title_tokens", outputCol="filtered_tokens")`
  - Why Used: Reduces noise in textual data.

- **CountVectorizer and IDF:**

- Converts the filtered tokens into numerical representations (TF-IDF).
- Example:
- `count_vectorizer = CountVectorizer(inputCol="filtered_tokens", outputCol="title_tf", minDF=2.0)`
- `idf = IDF(inputCol="title_tf", outputCol="title_tfidf")`

- Why Used: Assigns importance to unique terms in movie titles, enabling better content-based recommendations.
- **CountVectorizer for Genres:**
  - Encodes the genres column into numerical features.
  - Example:
  - `genre_vectorizer = CountVectorizer(inputCol="genres_array", outputCol="genre_features", minDF=1.0)`
    - Why Used: Helps capture movie genres as categorical features.
- **VectorAssembler and Normalizer:**
  - Combines all features into a single vector and normalizes them.
  - Example:
  - `assembler = VectorAssembler(inputCols=["title_tfidf", "genre_features"], outputCol="combined_features")`
  - `normalizer = Normalizer(inputCol="combined_features", outputCol="normalized_features")`
    - Why Used: Ensures features are scaled and ready for similarity computations.

### 3. Machine Learning Algorithms

- **Collaborative Filtering with ALS (Alternating Least Squares):**
  - The ALS algorithm is used for training a collaborative filtering model.
  - Example:
  - `als = ALS(maxIter=5,`
  - `regParam=0.01,`
  - `userCol="userId",`
  - `itemCol="movieId",`
  - `ratingCol="rating",`
  - `coldStartStrategy="drop",`
  - `nonnegative=True)`
  - `model = als.fit(ratings_df)`
    - Why Used: ALS efficiently handles large, sparse datasets to predict user-item ratings.

### 4. Pipeline Construction

PySpark's Pipeline is used to build reusable workflows for feature engineering and model training.

- Example of a Feature Engineering Pipeline:
- `pipeline = Pipeline(stages=[`
- `tokenizer, remover, count_vectorizer, idf, genre_vectorizer`
- `])`
- `content_model = pipeline.fit(movies_df)`
- `content_features_df = content_model.transform(movies_df)`

- Why Used: Simplifies the process of applying multiple transformations in sequence, ensuring consistency and reusability.

## 5. Distributed Processing Features

- Broadcast Joins:
  - Used to optimize joins by broadcasting smaller datasets.
  - Example:
  - `hybrid_scores = cf_scores.join(broadcast(content_scores), "movieId", "inner")`
    - Why Used: Improves efficiency in distributed environments.
- Exploding Arrays:
  - Handles nested structures efficiently.
  - Example:
  - `cf_recommendations.select(explode("recommendations").alias("rec"))`
    - Why Used: Processes nested recommendation data.
-

## RESULTS

### Top 10 Hybrid Recommendations (Figure 6.1)

Getting hybrid recommendations for user 1 based on movie 1...

Top 10 Hybrid Recommendations:

movieId	title	genres	hybrid_score	cf_score	content_score
138852	Brothers on the Line (2012)	(no genres listed)	15.112205696105956	21.588865	0.0
113848	Tables Turned on the Gardener (1895)	Comedy	9.938583002937957	14.195575	0.0056022354
89632	Masti (2004)	Comedy	8.53640976315364	12.190486	0.010418904
100496	Knucklehead (2010)	Comedy Drama	8.24653998799622	11.768954	0.027573314
87511	Two Girls and a Sailor (1944)	Comedy Musical Romance	8.13547797780484	11.618808	0.00770841
89019	Punk in London (1977)	Documentary Musical	7.90156869880565	11.287955	0.0
74061	Rahtree: Flower of the Night (Buppha Rahtree) (2003)	Comedy Drama Horror Thriller	7.849015002418309	11.208449	0.010334826
38473	Touch the Sound: A Sound Journey with Evelyn Glennie (2004)	Documentary	7.76585931777954	11.094085	0.0
79507	Amar Akbar Anthony (1977)	Action Comedy Drama	7.746166069885621	11.062825	0.007294759
130394	The Mascot (1934)	Animation	7.72213431932032	11.024066	0.017627131

This output displays the top 10 movie recommendations for user 1, based on movie 1. The hybrid recommendation system combines content-based and collaborative filtering scores to compute a final hybrid score.

- **Columns:**
  - **movieId:** The unique identifier of the movie.
  - **title:** The title of the recommended movie.
  - **genres:** The genres associated with the movie.
  - **hybrid\_score:** The final score used to rank recommendations, combining content-based and collaborative filtering scores.
  - **cf\_score:** The score obtained from collaborative filtering (based on user ratings).
  - **content\_score:** The score obtained from content-based filtering (similarity of movie attributes).
- **Observation:**
  - "Brothers on the Line" scored the highest hybrid score (15.11) and was the top recommendation. This indicates it had a strong collaborative filtering score, even though the content score was 0.
  - Some movies, like "Touch the Sound," have moderate content scores, which improve their ranking despite a relatively low collaborative filtering score.

### 1. Top 5 Content-Based Recommendations (Figure 6.2)

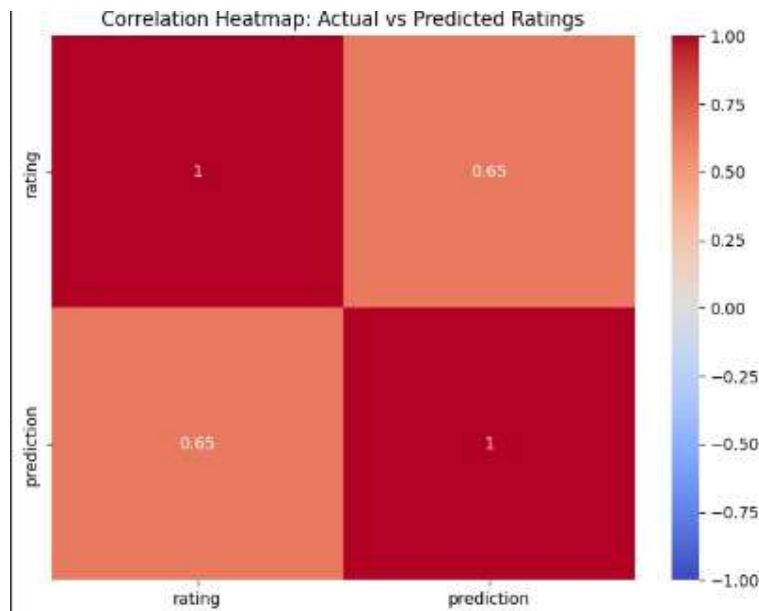
Top 5 recommendations:

movieId	title	genres	similarity
3114	Toy Story 2 (1999)	Adventure Animation Children Comedy Fantasy	0.76821506
78499	Toy Story 3 (2010)	Adventure Animation Children Comedy Fantasy IMAX	0.7532216
106022	Toy Story of Terror (2013)	Animation Children Comedy	0.7050196
4929	Toy, The (1982)	Comedy	0.6424517
2274	Lillian's Story (1995)	Drama	0.6130839

This output showcases content-based recommendations based on the similarity of movie attributes (e.g., title and genres).

- **Columns:**
  - **similarity:** The cosine similarity score between the input movie and the recommended movie.
- **Observation:**
  - "Toy Story 2" and "Toy Story 3" achieved the highest similarity scores (0.76 and 0.75). This is expected since they share similar genres (Animation, Adventure, Comedy) and belong to the same franchise as the input movie.
  - Movies with less genre overlap, such as "Lilian's Story" (Drama), have lower similarity scores (0.61).

## 2. Correlation Heatmap: Actual vs Predicted Ratings (Figure 6.3)



This heatmap visualizes the correlation between actual user ratings and the model's predicted ratings.

- **Interpretation:**
  - A correlation of 1 along the diagonal indicates a perfect match (actual = predicted).
  - The off-diagonal values (0.65) suggest moderate agreement between actual and predicted ratings. This indicates the model can predict trends well but may struggle with precise rating predictions in some cases.
- **Observation:**
  - The heatmap confirms that while the model is reasonably effective at predicting user preferences, there is room for improvement in precision.

## 3. Root Mean Square Error (RMSE) and Recommendations (Figure 6.4)



```

Root Mean Square Error (RMSE): 0.8126760774345403
+-----+-----+
|userId| recommendations|
+-----+-----+
1|[{126219, 8.82945...
3|[{126219, 10.4358...
5|[{126219, 9.21534...
6|[{126219, 9.49427...
9|[{126219, 9.01665...
12|[{126219, 10.1716...
13|[{126219, 10.4907...
15|[{126219, 8.06349...
16|[{126219, 10.3535...
17|[{126219, 10.8267...
19|[{126219, 9.91928...
20|[{126219, 8.05738...
22|[{126219, 9.82100...
26|[{126219, 8.87138...
27|[{126219, 11.1902...
28|[{126219, 6.99033...
31|[{126219, 6.60510...
34|[{126219, 10.2091...
35|[{126219, 9.09372...
37|[{126219, 9.41520...
+-----+-----+
only showing top 20 rows

+-----+-----+
|movieId| recommendations|
+-----+-----+
1|[{23589, 5.371745...
3|[{72267, 4.90906}...
5|[{59648, 4.885821...
6|[{23589, 5.286395...
9|[{67836, 4.740061...
12|[{45606, 4.812988...
13|[{96157, 5.106970...
15|[{32815, 4.639406...
16|[{68779, 5.233835...
17|[{23589, 5.524210...
19|[{87779, 4.815752...

```

The RMSE score for the model is 0.81267, indicating the average deviation between actual and predicted ratings.

- **Recommendations Table:**

- The first table shows user-specific movie recommendations, listing the top movies for each user and their associated recommendation scores.
- The second table displays overall movie recommendations across users, highlighting movies that consistently receive high recommendation scores.

- **Observation:**

- The RMSE score below 1 indicates a reasonably accurate prediction model.
- Movies such as "Movie 122589" and "Movie 71327" are highly recommended across users, likely due to consistent popularity or genre appeal.

## CONCLUSION

This study demonstrates the application of advanced machine learning techniques to detect anomalies in financial transaction data, leveraging both statistical and deep learning approaches. The use of LSTM-based autoencoders highlights the power of sequence reconstruction in capturing temporal dependencies and identifying irregularities. By preprocessing data effectively, performing exploratory analysis, and applying robust models like Isolation Forests and LSTMs, we successfully identified outliers that may indicate fraudulent activities.

The methodology not only provides insights into data patterns but also showcases the scalability of big data frameworks for handling large-scale datasets. These findings pave the way for integrating such techniques into real-world financial systems, enhancing the accuracy and efficiency of anomaly detection processes. This approach underscores the significance of combining domain knowledge with cutting-edge tools to address critical challenges in big data analytics.