

Movie Ratings Data Analysis Using Apache Sqoop and Apache Hive

**Project submitted to the
SRM UNIVERSITY-AP, Andhra Pradesh
for partial fulfilment of the requirements to award the degree
of
Master of Technology
in
Data Science**

**Submitted by
BANDLA RAVI - AP24122060008
THOTA TARUN KUMAR - AP24122060021
POLAROWTHU JASWANTH -AP24122060023**



**SRM University-AP
Neerukonda, Mangalagiri, Guntur
Andhra Pradesh – 522 240
[DECEMBER, 2024]**

Certificate

Date: 28-Nov-24

This is to certify that the work presented in this Project entitled “Movie Ratings Data Analysis Using Apache Scoop and Apache Hive” has been carried out by Ravi, TarunKumar, Jaswanth under my supervision. The work is genuine, original, and suitable for submission to SRM University for the award of Master of Technology.

Supervisor

DR. Rajiv Senapati

Assistant Professor

Department of Computer Science and Engineering

❖ Acknowledgments

This project was brought to fruition with invaluable support and resources from several key platforms and communities. The Hortonworks Sandbox served as an exceptional environment for big data processing, offering a robust platform for experimenting with tools like HDFS, HIVE and Sqoop. GitHub played a pivotal role by providing a reliable and accessible repository for managing and sharing project data. The open-source community deserves heartfelt gratitude for its unwavering dedication to creating, maintaining, and documenting tools and frameworks, which formed the backbone of this project. Their collective efforts continue to inspire innovation and make advanced technologies accessible to all.

Table of Contents:

Movie Ratings Data Analysis Using Apache Scoop and Apache Hive	1
Certificate	2
Acknowledgement	3
Table of Contents	4
Abstract	5
Introduction	6-7
Dataset Description	8
Implementation	9-38
Conclusion	39
References	40

❖ Abstract

The project showcases a comprehensive end-to-end workflow for big data processing, demonstrating the seamless integration of various tools and platforms to manage and analyze data effectively. The process begins with transferring a dataset from GitHub, a widely used version control and collaboration platform, to the Hadoop Distributed File System (HDFS). This step ensures efficient data storage and scalability, leveraging the distributed architecture of HDFS.

Once stored, the dataset is ingested into MySQL, a robust relational database management system, to facilitate structured storage and relational querying. Apache Sqoop is then utilized as a bridge to move the data seamlessly from MySQL into Apache Hive. Hive, a powerful data warehousing tool, allows the execution of SQL-like queries on the dataset, providing an intuitive interface for big data analytics.

The project's final stage focuses on performing insightful queries in Hive, extracting meaningful patterns and trends from the data. This workflow exemplifies how big data tools collaborate to simplify complex data operations, streamline workflows, and enhance organizational capabilities in deriving actionable insights. It underscores the importance of integrated ecosystems in big data analytics, empowering users to harness the full potential of their data efficiently and effectively for decision-making and innovation.

❖ Introduction

Data-driven decision-making has become a cornerstone of modern organizations, enabling them to derive actionable insights from vast amounts of data. This project focuses on integrating components of the Hadoop ecosystem to demonstrate how big data tools can simplify and enhance the management and analysis of user movie ratings. These components—HDFS, MySQL, Sqoop, and Hive—work together to create an end-to-end data processing pipeline that handles tasks ranging from data storage and transfer to querying and insights generation.

Efficient Data Transfer Across Platforms:

Efficiently transferring data between platforms is a critical aspect of big data workflows. In this project, the dataset is initially sourced from GitHub, a version control platform that allows collaboration and secure storage of data.

1. Transfer to HDFS

The dataset is first moved to the Hadoop Distributed File System (HDFS), which ensures scalability and fault tolerance. HDFS divides the data into blocks and replicates them across different nodes in a cluster, guaranteeing availability even if some nodes fail. This feature is particularly crucial for managing large datasets that traditional file systems might struggle to handle.

2. Ingestion into MySQL

The dataset is then ingested into MySQL, a relational database management system (RDBMS). MySQL is used to store data in structured tables with predefined schemas. This stage facilitates a systematic organization of data, allowing relational operations such as joins, filtering, and aggregation. Relational storage is advantageous for datasets requiring strict schema enforcement and transaction support.

3. Data Movement Using Apache Sqoop

Apache Sqoop bridges the gap between relational databases like MySQL and big data tools like Hive. By transferring structured data from MySQL to Hive, Sqoop allows the dataset to transition seamlessly into the Hadoop ecosystem. This transfer is automated, ensuring minimal manual intervention while maintaining data integrity. Sqoop can handle incremental imports, making it efficient for keeping datasets synchronized between systems.

Integration of Hadoop Ecosystem Components:

Each component of the Hadoop ecosystem plays a specific role in the workflow, and their integration showcases the ecosystem's power:

HDFS → ensures reliable storage of large datasets.

MySQL → provides structured relational storage for initial data organization.

Sqoop → automates data transfer, bridging traditional databases with big data tools.

Hive → facilitates querying and insights generation, enabling non-programmers to interact with big data effectively.

By integrating HDFS, MySQL, Sqoop, and Hive, this project exemplifies the power of the Hadoop ecosystem in managing and analyzing large datasets. It provides a roadmap for leveraging these tools to gain actionable insights, underscoring their importance in the era of data-driven decision-making. This workflow not only highlights the strengths of each component but also demonstrates their synergy in creating robust data pipelines.

❖ Dataset Description

The ratings.data dataset is a rich source of information that captures user behavior and movie preferences. With over 100,000 rows of user ratings, it provides a comprehensive view of interactions between users and movies. The dataset is tab-separated (TSV) and structured as follows:

Dataset Structure

1. **user_id**
 - A unique identifier for each user in the system.
 - Represents the individual providing the movie rating.
 - Use cases:
 - Analyzing user behavior and preferences.
 - Identifying active and passive users based on rating frequency.
2. **movie_id**
 - A unique identifier for each movie rated by users.
 - Represents the entity receiving the rating.
 - Use cases:
 - Measuring the popularity of movies.
 - Grouping ratings by genre or release year (if additional metadata is available).
3. **rating**
 - A numerical score (1–5) indicating the user's opinion of the movie.
 - Represents the quality or enjoyment level perceived by the user.
 - Use cases:
 - Sentiment analysis of user feedback.
 - Identifying top-rated movies or poorly rated ones.
4. **timestamp**
 - The Unix timestamp of when the rating was submitted.
 - Represents the temporal aspect of the data.
 - Use cases:
 - Trend analysis over time.
 - Evaluating the impact of events (e.g., promotions or new releases) on user activity.

❖ Implementation

1. Data Transfer to HDFS:

Objective:

To move the dataset from GitHub to the Hadoop Distributed File System (HDFS) for scalable and fault-tolerant storage.

Steps:

1. Download the Dataset from GitHub

- Use the wget command to fetch the dataset directly from the specified GitHub repository.
- The wget command retrieves files from the web using HTTP, HTTPS, or FTP protocols.

Command Used:

```
wget https://raw.githubusercontent.com/PJ07-09-2003/Big-  
Data/refs/heads/main/ratings.data
```

After executing the above commands, ensure the dataset (ratings.data) is present in the local directory by listing files with the ls command.

Command Used:

```
ls
```

2. Create a Directory in HDFS

- Create a dedicated directory in HDFS to store the dataset. The mkdir command initializes the directory structure.

Command Used:

```
hdfs dfs -mkdir /user/maria_dev/ratings_data
```

The path /user/maria_dev/ratings_data ensures that the dataset is stored under the user-specific directory in HDFS.

3. Transfer the Dataset to HDFS

- Use the put command to move the ratings.data file from the local system to the newly created HDFS directory.

Command Used:

```
hdfs dfs -put ratings.data /user/maria_dev/ratings_data
```

This command ensures that the dataset is now accessible in the distributed storage environment.

4. Verify the File in HDFS

- List the contents of the HDFS directory to confirm the successful transfer of the dataset.

Command Used:

```
hdfs dfs -ls /user/maria_dev/ratings_data
```

2. MySQL Integration

Objective:

To temporarily move the dataset to a local directory and ingest it into a MySQL database for relational storage and validation.

Steps:

1. Retrieve Dataset from HDFS to Local System

- Use the get command to move the dataset back to the local directory for MySQL ingestion.

Command Used:

```
hdfs dfs -get /user/maria_dev/ratings_data/ratings.data /tmp/ratings.data
```

The /tmp/ratings.data path ensures the file is stored in a temporary local directory.

2. Open MySQL Shell

To open the MySQL shell in the Hadoop sandbox, you can use the following command:

Command Used:

```
mysql -u root -p
```

Once executed, you will be prompted to enter the MySQL root user's password. In the Hadoop sandbox, the default password is usually:

Password is: Hadoop

After entering the password, you'll enter the MySQL shell where you can execute SQL commands.

3. Create and Use a Database in MySQL

- Launch MySQL and create a database named bigdata_db to store the ratings data.

Command Used:

```
CREATE DATABASE bigdata_db;  
USE bigdata_db;
```

4. Create a Table and Load Data

- Create a table named ratings to store the data. The table schema matches the dataset structure:
 - user_id (INT)
 - movie_id (INT)
 - rating (FLOAT)
 - timestamp (BIGINT)

- Load the data into the table using the LOAD DATA LOCAL INFILE command.

Command Used:

```
LOAD DATA LOCAL INFILE '/tmp/ratings.data'  
INTO TABLE ratings  
FIELDS TERMINATED BY '\t'  
LINES TERMINATED BY '\n'  
(user_id, movie_id, rating, timestamp);
```

5. Validate Data Ingestion

- Check the total number of records ingested into the table using the COUNT(*) query.

Command Used:

```
SELECT COUNT(*) FROM ratings;
```

This ensures that the dataset has been successfully loaded into MySQL.

3. Data Import with Sqoop

Objective:

To move data from MySQL to Hive for distributed storage and querying using Apache Sqoop.

Steps:

1. Set Up Sqoop Import Command

- Use the sqoop import command to transfer data from the MySQL database to Hive.

Command Used:

```
sqoop import \
--connect "jdbc:mysql://sandbox.hortonworks.com:3306/bigdata_db" \
--username root --password hadoop \
--table ratings --hive-import \
--hive-table ratings --delete-target-dir --m 1
```

Explanation:

- --connect: Specifies the JDBC URL to connect to the MySQL database.
- --username and --password: Provide authentication credentials for MySQL.
- --table: Specifies the table to be imported (ratings).
- --hive-import: Indicates that data should be imported into Hive.
- --hive-table: Names the Hive table where data will be stored.
- --delete-target-dir: Ensures the target directory is cleared before importing.
- --m 1: Specifies a single mapper job for simplicity.

2. Verify Data in Hive Using Ambari Dashboard

1. **Log into Ambari Dashboard:** Open your web browser and go to the Ambari dashboard URL.
2. **Access Hive View:** In the Ambari dashboard, find the **Hive** service in the left-hand navigation panel. Click on it to open the Hive service view.
3. **Go to the Files View:** In the Hive view, navigate to the **Files** tab or **Files View** section. This view will give you a file browser interface where you can see all the files and directories in HDFS.

1.1 Listing Databases in MySQL Using Sqoop

To identify the available databases in a MySQL server, the sqoop list-databases command is used. This helps users understand the data structure and available tables before performing operations.

Command Used:

```
sqoop list-databases --connect jdbc:mysql://localhost/bigdata_db --username root -P
```

Purpose: Lists all databases in the MySQL server.

1.2 Executing Queries Using Sqoop Eval

The sqoop eval command allows you to execute SQL queries on a database directly from Sqoop. This is useful for validating data before importing or exporting.

Command Used:

```
sqoop eval --connect jdbc:mysql://localhost/bigdata_db --username root -P --query  
"select * from ratings;"
```

Purpose: Executes the query SELECT * FROM ratings; to display all records from the ratings table in MySQL.

1.3 Importing Full Data from MySQL to HDFS

The sqoop import command transfers an entire table from MySQL to HDFS.

Command Used:

```
sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password  
hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver --m 1
```

Key Options:

- --connect: Specifies the JDBC URL to connect to the database.
- --username and --password: Authentication details for MySQL.
- --table: Specifies the table to import (ratings).
- --warehouse-dir: Defines the HDFS directory to store the data.
- --m 1: Uses a single mapper for the task.

Significance:

- Ensures the ratings table is imported into HDFS for processing.
- Creates an HDFS directory (/user/ratings) containing the imported data.

1.3.1 Validating Data in HDFS

a) Listing the Imported Data

To check the files created during import:

Command Used:

```
hdfs dfs -ls /user/ratings
```

b) Viewing Sample Data

To inspect the first few records in the file:

Command Used:

```
hdfs dfs -cat /user/ratings/part-m-00000 | head
```

c) Checking HDFS File Count

To validate the file count and sizes:

Command Used:

```
hdfs dfs -count /user/ratings
```

1.4 Importing Incremental Data (Append Mode)

Incremental imports allow adding only new or updated records to HDFS. This is critical for regularly updated datasets.

Command Used:

```
sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver --m 1 --incremental append --check-column rating --last-value 5
```

- **Key Options:**

- --incremental append: Adds new records based on the check-column.
- --check-column rating: Uses the rating column to identify new records.
- --last-value 5: Specifies the last imported value for rating.

Significance: Incremental imports are efficient for synchronizing only new data into HDFS without re-importing the entire table.

1.5 Exporting Data from HDFS to MySQL

To export processed or transformed data from HDFS back to MySQL, the sqoop export command is used.

Command Used:

```
sqoop export --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --export-dir /user/ratings/part-m-00000 --driver com.mysql.jdbc.Driver --m 1
```

- **Key Options:**

- --export-dir: Specifies the HDFS directory containing the data to export.
- --table: The destination table in MySQL.

Significance:

- Enables sharing processed data with downstream systems or applications.
- Demonstrates a complete data pipeline from MySQL to HDFS and back.

Overall Outputs of Execution:

```
raw.githubusercontent.com/P... | Hortonworks Sandbox with HDI | Ambari - Sandbox | root@sandbox:~ Shell In A Box | 127.0.0.1:4200

sandbox login: root
root@sandbox,hortonworks.com's password:
Last login: Tue Nov 26 18:42:52 2024 from 172.17.0.2
[root@sandbox ~]# git clone <https://raw.githubusercontent.com/PJ07-09-2003/Big-Data/refs/heads/main/ratings.data>
-bash: syntax error near unexpected token `newline'
[root@sandbox ~]# ls
anaconda-ks.cfg build.out hdp install.log install.log.syslog ratings.data sandbox.info start_ambari.sh start_hbase.sh
[root@sandbox ~]# vi ratings.data
[root@sandbox ~]# hdfs dfs -mkdir /user/root/ratings_data
[root@sandbox ~]# /user/root/ratings_data: No such file or directory
[root@sandbox ~]# hdfs dfs -ls /user
Found 13 items
drwxr-xr-x  - admin    hdfs      0 2016-10-25 08:11 /user/admin
drwxrwx---  - ambari-qa hdfs      0 2016-10-25 07:47 /user/ambari-qa
drwxr-xr-x  - amy_ds   hdfs      0 2016-10-25 08:02 /user/amy_ds
drwxr-xr-x  - hbase    hdfs      0 2016-10-25 07:48 /user/hbase
drwxr-xr-x  - hcat     hdfs      0 2016-10-25 07:51 /user/hcat
drwxr-xr-x  - hive     hdfs      0 2016-10-25 08:10 /user/hive
drwxr-xr-x  - holger_gov hdfs      0 2016-10-25 08:03 /user/holger_gov
drwxrwxr-x  - ivy      hdfs      0 2016-10-25 07:49 /user/ivy
drwxr-xr-x  - maria_dev hdfs      0 2016-10-25 07:58 /user/maria_dev
drwxrwxr-x  - oozie    hdfs      0 2016-10-25 07:52 /user/oozie
drwxr-xr-x  - raj_ops   hdfs      0 2016-10-25 08:04 /user/raj_ops
drwxrwxr-x  - spark    hdfs      0 2016-10-25 07:48 /user/spark
drwxr-xr-x  - zeppelin hdfs      0 2016-10-25 07:50 /user/zeppelin
[root@sandbox ~]# hdfs dfs -mkdir /user/maria_dev/ratings_data
[root@sandbox ~]# hdfs dfs -put ratings.data /user/maria_dev/ratings_data
[root@sandbox ~]# mysql -u root -p
Enter password:
Welcome to the MySQL monitor. Commands end with ; or \g.
Your MySQL connection id is 40
Server version: 5.6.34 MySQL Community Server (GPL)

Copyright (c) 2000, 2016, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

04:25 PM 27-11-2024
```

```
raw.githubusercontent.com/P... | Hortonworks Sandbox with HDI | Ambari - Sandbox | root@sandbox:~ Shell In A Box | 127.0.0.1:4200

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> CREATE DATABASE bigdata_db;
Query OK, 1 row affected (0.01 sec)

mysql> USE bigdata_db;
Database changed
mysql> CREATE TABLE ratings ( user_id INT NOT NULL, movie_id INT NOT NULL, rating INT NOT NULL, timestamp BIGINT NOT NULL );
Query OK, 0 rows affected (0.33 sec)

mysql> exit
Bye
[root@sandbox ~]# hdfs dfs -get /user/maria_dev/ratings_data/ratings.data /tmp/ratings.data
[root@sandbox ~]# mysql -u root -p
Enter password:
Welcome to the MySQL monitor. Commands end with ; or \g.
Your MySQL connection id is 42
Server version: 5.6.34 MySQL Community Server (GPL)

Copyright (c) 2000, 2016, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> USE bigdata_db; LOAD DATA LOCAL INFILE '/tmp/ratings.data' INTO TABLE ratings FIELDS TERMINATED BY '\t' LINES TERMINATED BY '\n' (user_id, movie_id, rating, timestamp);
-- Verify the data SELECT COUNT(*) FROM ratings;
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
Query OK, 100003 rows affected (2.07 sec)
Records: 100003 Deleted: 0 Skipped: 0 Warnings: 0

mysql> SELECT COUNT(*) FROM ratings;
+-----+
| COUNT(*) |
+-----+
| 100003 |
+-----+
04:25 PM 27-11-2024
```

```
raw.githubusercontent.com/PIC | Hortonworks Sandbox with HDI | Ambari - Sandbox | root@sandbox-- Shell In A Box | +  
127.0.0.1:4200  
Copyright (c) 2000, 2016, Oracle and/or its affiliates. Other names may be trademarks of their respective owners.  
Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.  
mysql> hostname  
->  
-> Ctrl-C -- exit!  
Aborted  
[root@sandbox ~]# hostname  
sandbox.hortonworks.com  
[root@sandbox ~]# mysql -u root -p -h sandbox.hortonworks.com  
Enter password:  
Welcome to the MySQL monitor. Commands end with ; or \g.  
Your MySQL connection id is 49  
Server version: 5.6.34 MySQL Community Server (GPL)  
  
Copyright (c) 2000, 2016, Oracle and/or its affiliates. All rights reserved.  
  
Oracle is a registered trademark of Oracle Corporation and/or its affiliates. Other names may be trademarks of their respective owners.  
  
Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.  
mysql> exit  
Bye  
[root@sandbox ~]# sqoop import \ --connect "jdbc:mysql://sandbox.hortonworks.com:3306/bigdata_db" \ --username root \ --password hadoop \ --table ratings \ --hive-import \ --hive-table ratings \ --delete-target-dir \ --m 1  
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.  
Please set $ACCUMLIO_HOME to the root of your Accumulo installation.  
24/11/27 10:19:30 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Error parsing arguments for import:  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --connect  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: jdbc:mysql://sandbox.hortonworks.com:3306/bigdata_db  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --username  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: root  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --password  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: hadoop  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --table  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --hive-import  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --hive-table  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --delete-target-dir  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: --m  
24/11/27 10:19:30 ERROR tool.BaseSqoopTool: Unrecognized argument: 1  
Breaking news  
Eknath Shinde ...  
Search  04:26 PM  
27-11-2024
```

```
raw.githubusercontent.com/PIC | Hortonworks Sandbox with HDI | Ambari - Sandbox | root@sandbox-- Shell In A Box | +  
127.0.0.1:4200  
24/11/27 10:23:26 INFO mapreduce.Job: map 100% reduce 0%  
24/11/27 10:23:28 INFO mapreduce.Job: Job job_1732701689398_0001 completed successfully  
24/11/27 10:23:29 INFO mapreduce.Job: Counters: 30  
File System Counters  
FILE: Number of bytes read=0  
FILE: Number of bytes written=162763  
FILE: Number of read operations=0  
FILE: Number of large read operations=0  
FILE: Number of write operations=0  
HDFS: Number of bytes read=87  
HDFS: Number of bytes written=1979226  
HDFS: Number of read operations=4  
HDFS: Number of large read operations=0  
HDFS: Number of write operations=2  
Job Counters  
Launched map tasks=1  
Other local map tasks=1  
Total time spent by all maps in occupied slots (ms)=11649  
Total time spent by all reduces in occupied slots (ms)=0  
Total time spent by all map tasks (ms)=11649  
Total vcore-milliseconds taken by all map tasks=11649  
Total megabyte-milliseconds taken by all map tasks=2912250  
Map-Reduce Framework  
Map input records=100003  
Map output records=100003  
Input split bytes=87  
Spilled Records=0  
Failed Shuffles=0  
Merged Map outputs=0  
GC time elapsed (ms)=495  
CPU time spent (ms)=6440  
Physical memory (bytes) snapshot=145846272  
Virtual memory (bytes) snapshot=1941893120  
Total committed heap usage (bytes)=45613056  
File Input Format Counters  
Bytes Read=0  
File Output Format Counters  
Bytes Written=1979226  
24/11/27 10:23:29 INFO mapreduce.ImportJobBase: Transferred 1.8875 MB in 100.0079 seconds (19.3269 KB/sec)  
24/11/27 10:23:29 INFO mapreduce.ImportJobBase: Retrieved 100003 records.  
24/11/27 10:23:29 INFO mapreduce.TwoPointJobBase: Publishing Hive/Hcat import job data to listeners  
29C  
Partly sunny  
Search  05:33 PM  
27-11-2024
```

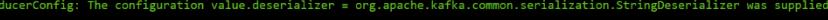
```
Try --help for usage instructions.
[root@sandbox ~]# sqoop import --connect "jdbc:mysql://sandbox.hortonworks.com:3306/bigdata_db" --username root --password hadoop --table ratings --hive-import --hive-table ratings --delete-target-dir -- -m 1
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please see $ACCUMULO_HOME to the root of your Accumulo installation.
24/11/27 18:21:35 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-2.5.0.0-1245
24/11/27 18:21:35 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
24/11/27 18:21:35 INFO tool.BaseSqoopTool: Using Hive-specific delimiters for output. You can override
24/11/27 18:21:35 INFO tool.BaseSqoopTool: delimiters with --fields-terminated-by, etc.
24/11/27 18:21:36 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
24/11/27 18:21:36 INFO tool.CodeGenTool: Beginning code generation
24/11/27 18:21:38 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM 'ratings' AS t LIMIT 1
24/11/27 18:21:39 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t LIMIT 1
24/11/27 18:21:39 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/hdp/2.5.0.0-1245/hadoop-mapreduce
Note: /tmp/sqoop-root/compile/1c0665b09c67a9eba76bd0536f990f2b/ratings.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
24/11/27 18:21:46 INFO tool.ImportTool: Writing jar file: /tmp/sqoop-root/compile/1c0665b09c67a9eba76bd0536f990f2b/ratings.jar
24/11/27 18:21:49 INFO tool.ImportTool: Destination directory rating is not present, hence not deleting.
24/11/27 18:21:49 WARN manager.MySQLManager: It looks like you are importing from mysql.
24/11/27 18:21:49 WARN manager.MySQLManager: This transfer can be faster! Use the --direct
24/11/27 18:21:49 WARN manager.MySQLManager: option to exercise the MySQL-specific fast path.
24/11/27 18:21:49 INFO manager.MySQLManager: Converting from NETMF behavior to convertToNull (mysql)
24/11/27 18:21:49 INFO manager.MySQLManager: Beginning import of ratings
24/11/27 18:21:51 INFO impl.TimelineClientImpl: Timeline service address: http://sandbox.hortonworks.com:8188/ws/v1/timeline/
24/11/27 18:21:51 INFO client.RMProxy: Connecting to ResourceManager at sandbox.hortonworks.com/172.17.0.2:8080
24/11/27 18:21:51 INFO client.AHSProxy: Connecting to Application History server at sandbox.hortonworks.com/172.17.0.2:10200
24/11/27 18:22:02 INFO db.DBInputFormat: Using read committed transaction isolation
24/11/27 18:22:02 INFO mapreduce.JobSubmitter: number of splits: 1
24/11/27 18:22:03 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1732701689398_0001
24/11/27 18:22:06 INFO impl.YarnClientImpl: Submitted application application_1732701689398_0001
24/11/27 18:22:06 INFO mapreduce.Job: The url to track the job: http://sandbox.hortonworks.com:8088/proxy/application_1732701689398_0001/
24/11/27 18:22:06 INFO mapreduce.Job: Running job: job_1732701689398_0001
24/11/27 18:23:11 INFO mapreduce.Job: Job job_1732701689398_0001 running in uber mode : false
24/11/27 18:23:11 INFO mapreduce.Job: map % 0% reduce 0%
24/11/27 18:23:26 INFO mapreduce.Job: map 100% reduce 0%
24/11/27 18:23:28 INFO mapreduce.Job: Job job_1732701689398_0001 completed successfully
24/11/27 18:23:29 INFO mapreduce.Job: Counters: 30
File System Counters
  FILE: Number of bytes read=0
  FILE: Number of bytes written=162763

```

```
raw.githubusercontent.com/P... Hortonworks Sandbox with HDI | Ambari - Sandbox | root@sandbox:~ Shell In A Box +  
29°C Partly sunny  
27.11.2024 05:33 PM  
24/11/27 16:23:29 INFO mapreduce.Job: Counters: 36  
File System Counters  
  FILE: Number of bytes read=0  
  FILE: Number of bytes written=162763  
  FILE: Number of read operations=0  
  FILE: Number of large read operations=0  
  FILE: Number of write operations=0  
  HDFS: Number of bytes read=87  
  HDFS: Number of bytes written=1979226  
  HDFS: Number of read operations=4  
  HDFS: Number of large read operations=0  
  HDFS: Number of write operations=2  
Job Counters  
  Launched map tasks=1  
    Other local map tasks=1  
    Total time spent by all maps in occupied slots (ms)=11649  
    Total time spent by all reduces in occupied slots (ms)=0  
    Total time spent by all map tasks (ms)=11649  
    Total vcore-milliseconds taken by all map tasks=11649  
    Total megabyte-milliseconds taken by all map tasks=2912250  
Map-Reduce Framework  
  Map input records=100003  
  Map output records=100003  
  Input split bytes=87  
  Spilled Records=0  
  Failed Shuffles=0  
  Merged Map outputs=0  
  GC time elapsed (ms)=495  
  CPU time spent (ms)=6440  
  Physical memory (bytes) snapshot=145846272  
  Virtual memory (bytes) snapshot=1941893120  
  Total committed heap usage (bytes)=45613056  
File Input Format Counters  
  Bytes Read=0  
File Output Format Counters  
  Bytes Written=1979226  
24/11/27 16:23:29 INFO mapreduce.ImportJobBase: Transferred 1.8875 MB in 100.0079 seconds (19.3269 KB/sec)  
24/11/27 16:23:29 INFO mapreduce.ImportJobBase: Retrieved 180003 records.  
24/11/27 16:23:29 INFO mapreduce.ImportJobBase: Publishing Hive/Hcat import job data to Listeners  
24/11/27 16:23:31 INFO atlas.ApplicationProperties: Looking for atlas-application.properties in classpath  
24/11/27 16:23:31 INFO atlas.ApplicationProperties: Loading atlas-application.properties from file:/etc/sqoop/2.5.0.0-1245/0/atlas-application.properties
```

```
Bytes Written=1979226
24/11/27 10:23:29 INFO mapreduce.ImportJobBase: Transferred 1.8875 MB in 100.0079 seconds (19.3269 KB/sec)
24/11/27 10:23:29 INFO mapreduce.ImportJobBase: Retrieved 100003 records.
24/11/27 10:23:29 INFO mapreduce.ImportJobBase: Publishing Hive/Hcat import job data to Listeners
24/11/27 10:23:31 INFO atlas.ApplicationProperties: Looking for atlas-application.properties in classpath
24/11/27 10:23:31 INFO atlas.ApplicationProperties: Loading atlas-application.properties from file:/etc/sqoop/2.5.0-e-1245/0/atlas-application.properties
24/11/27 10:23:32 ERROR security.InMemoryJAASConfiguration: Unable to add JAAS configuration for client [KafkaClient] as it is missing param [atlas.jaas.KafkaClient.log.inModuleFileName]. Skipping JAAS config for [KafkaClient]
24/11/27 10:23:34 INFO producer.ProducerConfig: ProducerConfig values:
metric.reporters = []
metadata.max.age.ms = 300000
reconnect.backoff.ms = 50
sasl.kerberos.ticket.renew.window.factor = 0.8
bootstrap.servers = [sandbox.hortonworks.com:6667]
ssl.keystore.type = JKS
sasl.mechanism = GSSAPI
max.block.ms = 60000
interceptor.classes = null
ssl.truststore.password = null
client.id =
ssl.endpoint.identification.algorithm = null
request.timeout.ms = 30000
acks = 1
receive.buffer.bytes = 32768
ssl.truststore.type = JKS
retries = 0
ssl.truststore.location = null
ssl.keystore.password = null
send.buffer.bytes = 131072
compression.type = none
metadata.fetch.timeout.ms = 60000
retry.backoff.ms = 100
sasl.kerberos.kinit.cmd = /usr/bin/kinit
buffer.memory = 33554432
timeout.ms = 30000
key.serializer = class org.apache.kafka.common.serialization.StringSerializer
sasl.kerberos.service.name = null
sasl.kerberos.ticket.renew.jitter = 0.05
ssl.trustmanager.algorithm = PKIX
29°C
Partly sunny
05:33 PM
27-11-2024
```

```
ssl.trustmanager.algorithm = PKIX
block.on.buffer.full = false
ssl.key.password = null
sasl.kerberos.min.time.before.relogin = 60000
connections.max.idle.ms = 540000
max.in.flight.requests.per.connection = 5
metrics.num.samples = 2
ssl.protocol = TLS
ssl.provider = null
ssl.enabled.protocols = [TLSv1.2, TLSv1.1, TLSv1]
batch.size = 16384
ssl.keystore.location = null
ssl.cipher.suites = null
security.protocol = PLAINTEXT
max.request.size = 1048576
value.serializer = class org.apache.kafka.common.serialization.StringSerializer
ssl.keymanager.algorithm = SunX509
metrics.sample.window.ms = 30000
partitioner.class = class org.apache.kafka.clients.producer.internals.DefaultPartitioner
linger.ms = 0
24/11/27 10:23:35 INFO producer.ProducerConfig: ProducerConfig values:
metric.reporters = []
metadata.max.age.ms = 300000
reconnect.backoff.ms = 50
sasl.kerberos.ticket.renew.window.factor = 0.8
bootstrap.servers = [sandbox.hortonworks.com:6667]
ssl.keystore.type = JKS
sasl.mechanism = GSSAPI
max.block.ms = 60000
interceptor.classes = null
ssl.truststore.password = null
client.id = producer-1
ssl.endpoint.identification.algorithm = null
request.timeout.ms = 30000
acks = 1
receive.buffer.bytes = 32768
ssl.truststore.type = JKS
retries = 0
ssl.truststore.location = null
29°C
Partly sunny
05:33 PM
27-11-2024
```

```
raw.githubusercontent.com/P... x Hortonworks Sandbox with HDI x | Ambari - Sandbox x root@sandbox-- Shell In A Box +  
← → C ⌂ 127.0.0.1:4200  
ssl.truststore.type = JKS  
retries = 0  
ssl.truststore.location = null  
ssl.keystore.password = null  
send.buffer.bytes = 131072  
compression.type = none  
metadata.fetch.timeout.ms = 60000  
retry.backoff.ms = 100  
sasl.kerberos.kinit.cmd = /usr/bin/kinit  
buffer.memory = 33554432  
timeout.ms = 30000  
key.serializer = class org.apache.kafka.common.serialization.StringSerializer  
sasl.kerberos.service.name = null  
sasl.kerberos.ticket.renew.jitter = 0.05  
ssl.trustmanager.algorithm = PKIX  
block.on.buffer.full = false  
ssl.key.password = null  
sasl.kerberos.min.time.before.relogin = 60000  
connections.max.idle.ms = 540000  
max.in.flight.requests.per.connection = 5  
metrics.num.samples = 2  
ssl.protocol = TLS  
ssl.provider = null  
ssl.enabled.protocols = [TLSv1.2, TLSv1.1, TLSv1]  
batch.size = 16384  
ssl.keystore.location = null  
ssl.cipher.suites = null  
security.protocol = PLAINTEXT  
max.request.size = 1048576  
value.serializer = class org.apache.kafka.common.serialization.StringSerializer  
ssl.keymanager.algorithm = SunX509  
metrics.sample.window.ms = 30000  
partitioner.class = class org.apache.kafka.clients.producer.internals.DefaultPartitioner  
linger.ms = 0  
  
24/11/27 10:23:35 WARN producer.ProducerConfig: The configuration key.deserializer = org.apache.kafka.common.serialization.StringDeserializer was supplied but isn't a known config.  
24/11/27 10:23:35 WARN producer.ProducerConfig: The configuration value.deserializer = org.apache.kafka.common.serialization.StringDeserializer was supplied but isn't a known config.  
24/11/27 10:23:35 WARN producer.ProducerConfig: The configuration hook.group.id = atlas was supplied but isn't a known config.  
0 29°C Partly sunny  05:34 PM 27-11-2024
```

```
root@sandbox:~$ curl -X POST -H "Content-Type: application/json" -d '{"records": [{"table": "ratings", "rating": 5, "user_id": 1, "movie_id": 1}, {"table": "ratings", "rating": 4, "user_id": 1, "movie_id": 2}, {"table": "ratings", "rating": 3, "user_id": 1, "movie_id": 3}, {"table": "ratings", "rating": 5, "user_id": 2, "movie_id": 1}, {"table": "ratings", "rating": 4, "user_id": 2, "movie_id": 2}, {"table": "ratings", "rating": 3, "user_id": 2, "movie_id": 3}, {"table": "ratings", "rating": 5, "user_id": 3, "movie_id": 1}, {"table": "ratings", "rating": 4, "user_id": 3, "movie_id": 2}, {"table": "ratings", "rating": 3, "user_id": 3, "movie_id": 3}], "operation": "insert"}' http://127.0.0.1:4200/_atlas/operations
```

```

raw.githubusercontent.com/PjX Hortonworks Sandbox with HD root@sandbox - Shell In A Box + 127.0.0.1:4200
← → ⌂ 127.0.0.1:4200
after 60000 ms.
    at org.apache.atlas.kafka.KafkaNotification.sendInternalToProducer(KafkaNotification.java:249)
    at org.apache.atlas.kafka.KafkaNotification.sendInternal(KafkaNotification.java:222)
    at org.apache.atlas.notification.AbstractNotification.send(AbstractNotification.java:84)
    at org.apache.atlas.hook.Atlasshook.notifyEntitiesInternal(Atlasshook.java:129)
    at org.apache.atlas.hook.Atlasshook.notifyEntities(Atlasshook.java:114)
    at org.apache.atlas.sqoop.hook.SqoopHook.publish(SqoopHook.java:177)
    at org.apache.atlas.sqoop.hook.SqoopHook.publish(SqoopHook.java:51)
    at org.apache.sqoop.mapreduce.PublishJobData.publishJobData(PublishJobData.java:52)
    at org.apache.sqoop.mapreduce.ImportJobBase.runImport(ImportJobBase.java:284)
    at org.apache.sqoop.manager.SqlManager.importTable(SqlManager.java:692)
    at org.apache.sqoop.manager.MySQLManager.importTable(MySQLManager.java:127)
    at org.apache.sqoop.tool.ImportTool.importTable(ImportTool.java:507)
    at org.apache.sqoop.tool.ImportTool.run(ImportTool.java:615)
    at org.apache.sqoop.Sqoop.run(Sqoop.java:147)
    at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:76)
    at org.apache.sqoop.Sqoop.runSqoop(Sqoop.java:183)
    at org.apache.sqoop.Sqoop.runTool(Sqoop.java:225)
    at org.apache.sqoop.Sqoop.runTool(Sqoop.java:234)
    at org.apache.sqoop.Sqoop.main(Sqoop.java:243)
Caused by: java.util.concurrent.ExecutionException: org.apache.kafka.common.errors.TimeoutException: Failed to update metadata after 60000 ms.
    at org.apache.kafka.clients.producer.KafkaProducer$FutureFailure.<init>(KafkaProducer.java:730)
    at org.apache.kafka.clients.producer.KafkaProducer.doSend(KafkaProducer.java:483)
    at org.apache.kafka.clients.producer.KafkaProducer.send(KafkaProducer.java:438)
    at org.apache.kafka.clients.producer.KafkaProducer.send(KafkaProducer.java:353)
    at org.apache.atlas.kafka.KafkaNotification.sendInternalToProducer(KafkaNotification.java:232)
    ... 18 more
Caused by: org.apache.kafka.common.errors.TimeoutException: Failed to update metadata after 60000 ms.
24/11/27 10:26:37 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `ratings` AS t LIMIT 1
24/11/27 10:26:37 INFO hive.HiveImport: Loading uploaded data into Hive
Logging initialized using configuration in jar:file:/usr/hdp/2.5.0.0-1245/hive/lib/hive-common-1.2.1000.2.5.0.0-1245.jar!/hive-log4j.properties
OK
Time taken: 15.693 seconds
Loading data to table default.ratings
Table default.ratings stats: [numFiles=1, numRows=0, totalSize=1979226, rawDataSize=0]
OK
Time taken: 3.356 seconds
[root@ sandbox ~]#
```

29°C Partly sunny 05:34 PM 27-11-2024

```

A Classwork for DSC 502 Big | ChatGPT | root@sandbox - Shell In | Hortonworks Sandbox with HD | Ambari - Sandbox | Big-Data/ratings.data at m | + 127.0.0.1:4200
← → ⌂ 127.0.0.1:4200
sandbox login: root
root@sandbox.hortonworks.com's password:
Last login: Mon Dec  2 18:33:45 2024 from 172.17.0.2
[root@sandbox ~]# sqoop import \ --connect jdbc:mysql://localhost/bigdata_db \ --username root \ --password hadoop \ --table ratings \ --warehouse-dir /user \ --driver com.mysql.jdbc.Driver \ -m 1
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:39:39 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Error parsing arguments for import:
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: --connect
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: jdbc:mysql://localhost/bigdata_db
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: --username
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: root
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: --password
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: hadoop
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: --table
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: ratings
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: --warehouse-dir
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: /user
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: --driver
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: com.mysql.jdbc.Driver
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: --m
24/12/02 18:39:39 ERROR tool.BaseSqoopTool: Unrecognized argument: 1

Try --help for usage instructions.
[root@sandbox ~]# sqoop import \ --connect jdbc:mysql://localhost/bigdata_db \ --username root \ --password hadoop \ --table ratings \ --warehouse-dir /user \ --driver com.mysql.jdbc.Driver \ -m 1
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:40:41 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Error parsing arguments for import:
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --connect
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: jdbc:mysql://localhost/bigdata_db
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --username
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: root
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --password
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: hadoop
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --table
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: ratings
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --warehouse-dir
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: /user
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --driver
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: com.mysql.jdbc.Driver
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --m
24/12/02 18:40:41 ERROR tool.BaseSqoopTool: Unrecognized argument: 1

[root@sandbox ~]#
```

Watchlist Ideas 12:30 AM 03-12-2024

```
Try --help for usage instructions.
[root@sandbox ~]# sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver -m 1
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:48:41 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Error parsing arguments for import:
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --connect
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: jdbc:mysql://localhost/bigdata_db
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --username
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: root
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --password
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: hadoop
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --table
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: ratings
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --warehouse-dir
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: /user
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: --driver
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: com.mysql.jdbc.Driver
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: -m
24/12/02 18:48:41 ERROR tool.BaseSqoopTool: Unrecognized argument: 1

Try --help for usage instructions.
[root@sandbox ~]# sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver -m 1
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:42:09 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:42:09 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
24/12/02 18:42:16 WARN sqoop.ConnFactory: Parameter --driver is set to an explicit driver however appropriate connection manager is not being set (via --connection-manager). Sqoop is going to fall back to org.apache.sqoop.manager.GenericJdbcManager. Please specify explicitly which connection manager should be used next time.
24/12/02 18:42:16 INFO manager.SqlManager: Using default fetchSize of 1000
24/12/02 18:42:16 INFO tool.CodeGenTool: Beginning code generation
24/12/02 18:42:11 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:42:11 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:42:11 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/hdp/2.5.0.0-1245/hadoop-mapreduce
Note: /tmp/sqoop-root/compile/ec02f936ac28461f45ed41c2b2d0cea5/ratings.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
24/12/02 18:42:17 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-root/compile/ec02f936ac28461f45ed41c2b2d0cea5/ratings.jar
24/12/02 18:42:17 INFO mapreduce.TwoPhaseJobBase: Beginning import of ratings
12:30 AM 03-12-2024
```

```
Try --help for usage instructions.
[root@sandbox ~]# sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver -m 1
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:42:09 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:42:09 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
24/12/02 18:42:16 WARN sqoop.ConnFactory: Parameter --driver is set to an explicit driver however appropriate connection manager is not being set (via --connection-manager). Sqoop is going to fall back to org.apache.sqoop.manager.GenericJdbcManager. Please specify explicitly which connection manager should be used next time.
24/12/02 18:42:16 INFO manager.SqlManager: Using default fetchSize of 1000
24/12/02 18:42:16 INFO tool.CodeGenTool: Beginning code generation
24/12/02 18:42:11 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:42:11 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:42:11 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/hdp/2.5.0.0-1245/hadoop-mapreduce
Note: /tmp/sqoop-root/compile/ec02f936ac28461f45ed41c2b2d0cea5/ratings.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
24/12/02 18:42:17 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-root/compile/ec02f936ac28461f45ed41c2b2d0cea5/ratings.jar
24/12/02 18:42:17 INFO mapreduce.ImportJobBase: Beginning import of ratings
24/12/02 18:42:18 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:42:21 INFO impl.TimelineClientImpl: Timeline service address: http://sandbox.hortonworks.com:8188/ws/v1/timeline/
24/12/02 18:42:21 INFO client.RMProxy: Connecting to ResourceManager at sandbox.hortonworks.com/172.17.0.2:8050
24/12/02 18:42:22 INFO client.AHSProxy: Connecting to Application History server at sandbox.hortonworks.com/172.17.0.2:10200
24/12/02 18:42:39 INFO db.DBInputFormat: Using read committed transaction isolation
24/12/02 18:42:39 INFO mapreduce.JobSubmitter: number of splits:1
24/12/02 18:42:40 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1733164307838_0001
24/12/02 18:42:45 INFO impl.YarnClientImpl: Submitted application application_1733164307838_0001
24/12/02 18:42:45 INFO mapreduce.Job: The url to track the job: http://sandbox.hortonworks.com:8088/proxy/application_1733164307838_0001/
24/12/02 18:42:45 INFO mapreduce.Job: Running job: job_1733164307838_0001
24/12/02 18:43:29 INFO mapreduce.Job: Job job_1733164307838_0001 running in uber mode : false
24/12/02 18:43:29 INFO mapreduce.Job: map 0% reduce 0%
24/12/02 18:43:48 INFO mapreduce.Job: map 100% reduce 0%
24/12/02 18:43:50 INFO mapreduce.Job: Job job_1733164307838_0001 completed successfully
24/12/02 18:43:51 INFO mapreduce.Job: Counters: 30
  File System Counters
    FILE: Number of bytes read=0
    FILE: Number of bytes written=162687
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=87
File Watchlist
Idea 12:30 AM 03-12-2024
```

```

24/12/02 18:43:29 INFO mapreduce.Job: Job job_1733164307838_0001 running in uber mode : false
24/12/02 18:43:29 INFO mapreduce.Job: map 0% reduce 0%
24/12/02 18:43:48 INFO mapreduce.Job: map 100% reduce 0%
24/12/02 18:43:50 INFO mapreduce.Job: Job job_1733164307838_0001 completed successfully
24/12/02 18:43:51 INFO mapreduce.Job: Counters: 30
  File System Counters
    FILE: Number of bytes read=0
    FILE: Number of bytes written=162687
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=87
    HDFS: Number of bytes written=1979226
    HDFS: Number of read operations=4
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Other local map tasks=1
    Total time spent by all maps in occupied slots (ms)=14626
    Total time spent by all reduces in occupied slots (ms)=0
    Total time spent by all map tasks (ms)=14626
    Total vcore-milliseconds taken by all map tasks=14626
    Total megabyte-milliseconds taken by all map tasks=3656500
  Map-Reduce Framework
    Map input records=100003
    Map output records=100003
    Input split bytes=87
    Spilled Records=0
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=153
    CPU time spent (ms)=7020
    Physical memory (bytes) snapshot=157761536
    Virtual memory (bytes) snapshot=1939927940
    Total committed heap usage (bytes)=60817408
  File Input Format Counters
    Bytes Read=0
  File Output Format Counters
    Bytes Written=1979226

```

```

24/12/02 18:43:51 INFO mapreduce.ImportJobBase: Transferred 1.8875 MB in 90.9698 seconds (21.247 KB/sec)
24/12/02 18:43:51 INFO mapreduce.ImportJobBase: Retrieved 100003 records.
[root@ sandbox ~]# hdfs dfs -ls /user/ratings
Found 2 items
-rw-r--r-- 1 root hdfs 1979226 2024-12-02 18:43 /user/ratings/_SUCCESS
-rw-r--r-- 1 root hdfs 1979226 2024-12-02 18:43 /user/ratings/part-m-00000
[root@ sandbox ~]# hdfs dfs -cat /user/ratings/part-m-00000 | head
0,50,5,881250949
0,172,5,881250949
0,133,1,881250949
196,242,3,881250949
186,302,3,891717742
22,377,1,878887116
244,51,2,880606923
166,346,1,886397596
298,474,4,884182806
115,265,2,881171488
cat: Unable to write to output stream.
[root@ sandbox ~]# hdfs dfs -count /user/ratings
      1          2        1979226 /user/ratings
[root@ sandbox ~]# sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver --m 1 --incremental append --check-column rating --last-value 5
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:47:45 INFO tool.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:47:45 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
24/12/02 18:47:45 WARN sqoop.ConnFactory: Parameter --driver is set to an explicit driver however appropriate connection manager is not being set (via --connection-manager). Sqoop is going to fall back to org.apache.sqoop.manager.GenericJdbcManager. Please specify explicitly which connection manager should be used next time.
24/12/02 18:47:45 INFO manager.SqlManager: Using default fetchSize of 1000
24/12/02 18:47:45 INFO tool.CodeGenTool: Beginning code generation
24/12/02 18:47:46 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:47:46 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0

```

```

[root@sandbox ~]# sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver --incremental append --check-column rating --last-value 5
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:47:45 INFO tool.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:47:45 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
24/12/02 18:47:45 WARN tool.SqoopConnectionFactory: Parameter --driver is set to an explicit driver however appropriate connection manager is not being set (via --connection-manager). Sqoop is going to fall back to org.apache.sqoop.manager.GenericJdbcManager. Please specify explicitly which connection manager should be used next time.
24/12/02 18:47:46 INFO manager.SqlManager: Using default fetchSize of 1000
24/12/02 18:47:46 INFO manager.SqlManager: Beginning code generation
24/12/02 18:47:46 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:47:46 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:47:46 INFO manager.SqlManager: Executing SQL statement: SELECT MAX(rating) FROM ratings
Note: /tmp/sqoop-root/tmp/_sqlfile_278bfid3ae223abc91fc8da2704cfa47d/ratings.java uses or overrides a deprecated API.
Note: /tmp/sqoop-root/tmp/_sqlfile_278bfid3ae223abc91fc8da2704cfa47d/ratings.java uses or overrides a deprecated API.
[root@sandbox ~]# sqoop export --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --export-dir /user/ratings/part-m-00000 --driver com.mysql.jdbc.Driver --m 1
Warning: /usr/hdp/2.5.0.0-1245/accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
24/12/02 18:47:51 INFO tool.Sqoop: Running Sqoop version: 1.4.6.2.5.0.0-1245
24/12/02 18:47:51 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
24/12/02 18:47:51 WARN tool.SqoopConnectionFactory: Parameter --driver is set to an explicit driver however appropriate connection manager is not being set (via --connection-manager). Sqoop is going to fall back to org.apache.sqoop.manager.GenericJdbcManager. Please specify explicitly which connection manager should be used next time.
24/12/02 18:47:51 INFO manager.SqlManager: Using default fetchSize of 1000
24/12/02 18:47:51 INFO tool.CodeGenTool: Beginning code generation
24/12/02 18:47:51 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:47:51 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:47:51 INFO manager.SqlManager: Executing SQL statement: HADOOP_MAPRED_HOME is /usr/hdp/2.5.0.0-1245/hadoop-mapreduce
Note: /tmp/sqoop-root/tmp/_sqlfile_278bfid3ae223abc91fc8da2704cfa47d/ratings.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
24/12/02 18:47:52 INFO manager.SqlManager: Writing jar file: /tmp/sqoop-root/compile/78bfid3ae223abc91fc8da2704cfa47d/ratings.jar
24/12/02 18:47:52 INFO tool.ImportTool: Maximal id query for Free form Incremental import: SELECT MAX(rating) FROM ratings
24/12/02 18:47:52 INFO tool.ImportTool: incremental import based on column rating
24/12/02 18:47:52 INFO tool.ImportTool: No new rows detected since last import.
[root@sandbox ~]#

```

```

[root@sandbox ~]# sqoop import --connect jdbc:mysql://localhost/bigdata_db --username root --password hadoop --table ratings --warehouse-dir /user --driver com.mysql.jdbc.Driver
24/12/02 18:52:22 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/hdp/2.5.0.0-1245/hadoop-mapreduce
Note: /tmp/sqoop-root/compile/a6ef27e9f7fd934fcdb1d2c6a7ca949/ratings.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
24/12/02 18:52:25 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-root/compile/a6ef27e9f7fd934fcdb1d2c6a7ca949/ratings.jar
24/12/02 18:52:25 INFO mapreduce.ExportJobBase: Beginning export of ratings
24/12/02 18:52:29 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:52:29 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:52:29 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM ratings AS t WHERE 1=0
24/12/02 18:52:30 INFO impl.TimelineClientImpl: Timeline service address: http://sandbox.hortonworks.com:8188/ws/v1/timeline/
24/12/02 18:52:30 INFO client.RMProxy: Connecting to ResourceManager at sandbox.hortonworks.com/172.17.0.2:8080
24/12/02 18:52:43 INFO input.FileInputFormat: Total input paths to process : 1
24/12/02 18:52:43 INFO input.FileInputFormat: Total input paths to process : 1
24/12/02 18:52:43 INFO lzo.GPLNativeCodeLoader: Loaded native gpl library
24/12/02 18:52:43 INFO lzo.LzoCodec: Successfully loaded & initialized native-lzo library [hadoop-lzo rev 7a4b57bedce694048432dd5bf5b90a6c8ccdba80]
24/12/02 18:52:43 INFO mapreduce.JobSubmitter: number of splits:1
24/12/02 18:52:44 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1733164307838_0002
24/12/02 18:52:45 INFO impl.YarnClientImpl: Submitted application application_1733164307838_0002
24/12/02 18:52:46 INFO mapreduce.Job: The url to track the job: http://sandbox.hortonworks.com:8088/proxy/application_1733164307838_0002/
24/12/02 18:52:46 INFO mapreduce.Job: Job: Running job: job_1733164307838_0002
24/12/02 18:53:09 INFO mapreduce.Job: Job job_1733164307838_0002 running in uber mode : false
24/12/02 18:53:09 INFO mapreduce.Job: map 0% reduce 0%
24/12/02 18:53:27 INFO mapreduce.Job: Job map 100% reduce 0%
24/12/02 18:53:29 INFO mapreduce.Job: Job job_1733164307838_0002 completed successfully
24/12/02 18:53:29 INFO mapreduce.Job: Counters: 0
File System Counters
FILE: Number of bytes read=0
FILE: Number of bytes written=162542
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=1979371
HDFS: Number of bytes written=0
HDFS: Number of read operations=4
HDFS: Number of large read operations=0
HDFS: Number of write operations=0
Job Counters
Launched map tasks=1
Data-local map tasks=1
Total time spent by all maps in occupied slots (ms)=14211
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=14211
Total vcore-milliseconds taken by all map tasks=14211

```

```
127.0.0.1:4200
Job Counters
Launched map tasks=1
Data-local map tasks=1
Total time spent by all maps in occupied slots (ms)=14211
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=14211
Total vcore-milliseconds taken by all map tasks=14211
Total megabyte-milliseconds taken by all map tasks=3552750

Map-Reduce Framework
Map input records=100003
Map output records=100003
Input split bytes=142
Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=1352
CPU time spent (ms)=5158
Physical memory (bytes) snapshot=141717504
Virtual memory (bytes) snapshot=1938128896
Total committed heap usage (bytes)=37748736

File Input Format Counters
Bytes Read=0
File Output Format Counters
Bytes Written=0
24/12/02 18:53:29 INFO mapreduce.ExportJobBase: Transferred 1.8877 MB in 60.2378 seconds (32.0891 KB/sec)
24/12/02 18:53:29 INFO mapreduce.ExportJobBase: Exported 100003 records.
[root@sandbox ~]# mysql -u root -p hadoop
Enter password:
ERROR 1049 (42000): Unknown database 'hadoop'
[root@sandbox ~]# mysql -u root -p
Enter password:
Welcome to the MySQL monitor. Commands end with ; or \g.
Your MySQL connection id is 59
Server version: 5.6.34 MySQL Community Server (GPL)

Copyright (c) 2000, 2016, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

25C Mostly cloudy 12:30 AM 03-12-2024
```

```
127.0.0.1:4200
File Output Format Counters
Bytes Written=0
24/12/02 18:53:29 INFO mapreduce.ExportJobBase: Transferred 1.8877 MB in 60.2378 seconds (32.0891 KB/sec)
24/12/02 18:53:29 INFO mapreduce.ExportJobBase: Exported 100003 records.
[root@sandbox ~]# mysql -u root -p hadoop
Enter password:
ERROR 1049 (42000): Unknown database 'hadoop'
[root@sandbox ~]# mysql -u root -p
Enter password:
Welcome to the MySQL monitor. Commands end with ; or \g.
Your MySQL connection id is 59
Server version: 5.6.34 MySQL Community Server (GPL)

Copyright (c) 2000, 2016, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> USE bigdata_db;
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> SHOW TABLES;
+-----+
| Tables_in_bigdata_db |
+-----+
| ratings              |
+-----+
1 row in set (0.00 sec)

mysql> SELECT COUNT(*) FROM ratings;
+-----+
| COUNT(*) |
+-----+
| 200006 |
+-----+
1 row in set (0.00 sec)

25C Mostly cloudy 12:31 AM 03-12-2024
```

4. Hive Queries

Objective:

To perform analytical queries on the imported data in Hive to derive insights.

Queries:

1. Basic Analysis:

1. Fetch the first 10 rows to verify data accuracy.

Query Used:

```
SELECT * FROM ratings LIMIT 10;
```

Output:

The screenshot shows the Hortonworks Sandbox with HDFS interface. The top navigation bar includes links for Ambari, Sandbox, Dags, and Alerts. The main window has tabs for Hive, Query, Saved Queries, History, UDFs, and Upload Table. On the left, a Database Explorer panel shows databases like default, Ratings, Example 07, Example 08, foodmart, and academo. The central area is a Query Editor with a Worksheet tab containing the SQL query: `1 SELECT * FROM Ratings LIMIT 10;`. Below the editor are buttons for Execute, Explain, and Save as... A status bar at the bottom indicates "Query Process Results (Status: SUCCEEDED)". The bottom section displays the results of the query as a table:

	ratings.user_id	ratings.movie_id	ratings.rating	ratings.timestamp
0	50	5	881250949	
0	172	5	881250949	
0	133	1	881250949	
196	242	3	881250949	
186	302	3	891717742	
22	377	1	878887116	
244	51	2	880606923	
166	346	1	886397596	
298	474	4	884182806	
115	265	2	881171488	

2. Determine the maximum and minimum ratings.

Query Used:

```
SELECT MAX(rating) AS max_rating, MIN(rating) AS min_rating FROM ratings;
```

Output:

The screenshot shows a browser window with a tab titled '127.0.0.1:8080/#/main/views/HIVE/1.5.0/AUTO_HIVE_INSTANCE'. The main content area is a 'Query Editor' with a 'Worksheet' pane containing the SQL query:

```
1 SELECT MAX(rating) AS max_rating, MIN(rating) AS min_rating FROM Ratings;
```

Below the worksheet are buttons for 'Execute', 'Explain', and 'Save as...'. To the right is a sidebar with icons for 'SQL', 'TEZ', and a mail icon with a red notification dot. The bottom section is a 'Query Process Results' table with a status of 'SUCCEEDED'. It has tabs for 'Logs' and 'Results', with the 'Results' tab selected. The results show two columns: 'max_rating' and 'min_rating', with values 5 and 1 respectively. The browser's taskbar at the bottom shows various open applications, and the system tray indicates the date and time as 27-11-2024 04:29 PM.

3. Calculate the average rating for each movie.

Query Used:

```
SELECT movie_id, AVG(rating) AS avg_rating  
FROM ratings  
GROUP BY movie_id;
```

Output:

The screenshot shows a web-based interface for running Hive queries. The top navigation bar includes tabs for 'Hive', 'Query', 'Saved Queries', 'History', 'UDFs', and 'Upload Table'. On the left, a 'Database Explorer' panel shows a 'default' database with tables like 'ratings', 'sample_07', 'sample_08', 'foodmart', and 'addemo'. The main area is a 'Query Editor' with a 'Worksheet' pane containing the following SQL code:

```
1 SELECT movie_id, AVG(rating) AS avg_rating  
2 FROM Ratings  
3 GROUP BY movie_id;
```

Below the worksheet are buttons for 'Execute', 'Explain', and 'Save as...'. A progress bar indicates the query is at 100% completion. At the bottom, a 'Logs' tab is selected in the 'Query Process Results' section, which displays the status as 'SUCCEEDED'. The results tab shows the output of the query:

movie_id	avg_rating
1	3.8783185840707963
2	3.2061068702290076
3	3.033333333333333
4	3.550239234449761
5	3.302325581395349
6	3.576923076923077
7	3.798469387755102
8	3.9954337899543377
9	3.8963210702341136
10	3.831460674157303
11	3.847457627118644
12	4.385767790262173
13	3.4184782608695654
14	3.9672131147540983
15	3.7781569965870307

4. Identify the highest-rated movie based on average ratings.

Query Used:

```
SELECT movie_id, AVG(rating) AS avg_rating  
FROM ratings  
GROUP BY movie_id  
ORDER BY avg_rating DESC  
LIMIT 1;
```

Output:

The screenshot shows a browser window with several tabs open. The active tab is titled '127.0.0.1:8080/#/main/views/HIVE/1.5.0/AUTO_HIVE_INSTANCE'. The main content area is a 'Worksheet' interface. On the left, there's a sidebar with a dropdown set to 'default', a search bar for tables, and a 'Databases' section listing 'default', 'ratings', 'sample 07', 'sample 08', 'foodmart', and 'xdademo'. The central area contains the SQL query:

```
1 SELECT movie_id, AVG(rating) AS avg_rating  
2 FROM Ratings  
3 GROUP BY movie_id  
4 ORDER BY avg_rating DESC  
5 LIMIT 1;
```

Below the query are buttons for 'Execute', 'Explain', and 'Save as...'. To the right of the worksheet is a vertical toolbar with icons for 'SQL', 'TEZ', and other options. At the bottom, a 'Logs' tab is visible, showing the status 'SUCCEEDED'. The results table has columns 'movie_id' and 'avg_rating', with one row: '814' and '5.0'. The browser's taskbar at the bottom shows various application icons and the date/time '27-11-2024 04:35 PM'.

5. Identify the highest-rated movie based on average ratings.

Query Used:

```
SELECT movie_id, AVG(rating) AS avg_rating  
FROM ratings  
GROUP BY movie_id  
ORDER BY avg_rating DESC  
LIMIT 1;
```

Output:

The screenshot shows a browser window with several tabs open. The active tab is a worksheet in a SQL editor, displaying the following code:

```
1 SELECT COUNT(DISTINCT user_id) FROM Ratings;
```

Below the code, there are buttons for "Execute", "Explain", and "Save as...". The status bar indicates "100%".

After executing the query, a new window titled "Query Process Results (Status: SUCCEEDED)" appears. It has tabs for "Logs" and "Results". The "Results" tab shows a single row of data:

c0
944

The browser's taskbar at the bottom shows various application icons, and the system tray indicates it's 04:40 PM on 27-11-2024.

2. Complex Queries

1. Find movies with more than 500 ratings.

Query Used:

```
SELECT movie_id, COUNT(*) AS num_ratings
FROM ratings
GROUP BY movie_id
HAVING num_ratings > 500;
```

Output:

The screenshot shows the Hortonworks Sandbox interface with two browser windows. Both windows have the URL `127.0.0.1:8080/#/main/views/HIVE/1.5.0/AUTO_HIVE_INSTANCE`.

Top Window (Query Editor):

- Database Explorer:** Shows databases: default, Ratings, sample_07, sample_08, foodmart, academo.
- Query Editor:** Worksheet pane contains the following SQL query:

```
1 SELECT movie_id, COUNT(*) AS num_ratings
2 FROM Ratings
3 GROUP BY movie_id
4 HAVING num_ratings > 500;
```
- Buttons:** Execute, Explain, Save as.., New Worksheet.

Bottom Window (Query Process Results):

- Logs:** Status: SUCCEEDED.
- Results:** A table showing movie IDs and their ratings counts:

movie_id	num_ratings
50	584
100	508
181	507
268	509

Licensed under the Apache License, Version 2.0.
See third-party tools/resources that Ambari uses and their respective authors

2. Get the average rating for each movie, along with the number of ratings:

Query Used:

```
SELECT movie_id, COUNT(*) AS num_ratings, AVG(rating) AS avg_rating
FROM Ratings
GROUP BY movie_id
ORDER BY avg_rating DESC;
```

Output:

The screenshot shows the Hortonworks Sandbox interface. In the top navigation bar, there are tabs for 'Hive', 'Query', 'Saved Queries', 'History', 'UDFs', and 'Upload Table'. The 'Query' tab is selected. Below the tabs is a 'Database Explorer' panel showing databases: default, Ratings, Sample_07, Sample_08, Foodmart, and Ademo. The main area is the 'Query Editor' with a 'Worksheet' tab containing the SQL query:

```
1 SELECT movie_id, COUNT(*) AS num_ratings, AVG(rating) AS avg_rating
2 FROM Ratings
3 GROUP BY movie_id
4 ORDER BY avg_rating DESC;
```

Below the worksheet are buttons for 'Execute', 'Explain', 'Save as...', and 'New Worksheet'. The status bar at the bottom indicates 'Status: SUCCEEDED'. The 'Results' tab in the 'Query Process Results' section shows a table with columns: movie_id, num_ratings, and avg_rating. The table data is as follows:

movie_id	num_ratings	avg_rating
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254

This screenshot shows the same Hortonworks Sandbox interface as the previous one, but the 'Results' tab is more populated with data. The table now includes many more rows, showing the average rating for a large number of movies. The columns remain the same: movie_id, num_ratings, and avg_rating.

movie_id	num_ratings	avg_rating
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318	298	4.466442953020135
169	118	4.466101694915254
1653	1	5.0
1293	3	5.0
1467	2	5.0
814	1	5.0
1500	2	5.0
1201	1	5.0
1122	1	5.0
1189	3	5.0
1599	1	5.0
1536	1	5.0
1449	8	4.625
1594	2	4.5
119	4	4.5
1398	2	4.5
1642	2	4.5
408	112	4.491071428571429
318		

3. Find the user with the most ratings

Query Used:

```
SELECT user_id, COUNT(*) AS num_ratings
FROM Ratings
GROUP BY user_id
ORDER BY num_ratings DESC
LIMIT 1;
```

Output:

The screenshot shows a browser window with multiple tabs. The active tab is titled '127.0.0.1:8080/#/main/views/HIVE/1.5.0/AUTO_HIVE_INSTANCE'. The main content area is a 'Worksheet' interface. On the left, there's a sidebar with a dropdown set to 'default', a search bar, and a 'Databases' section listing 'default', 'Ratings', 'sample_07', 'sample_08', 'foodmart', and 'gademo'. The main workspace contains the following SQL query:

```
1 SELECT user_id, COUNT(*) AS num_ratings
2 FROM Ratings
3 GROUP BY user_id
4 ORDER BY num_ratings DESC
5 LIMIT 1;
```

Below the query are three buttons: 'Execute', 'Explain', and 'Save as...'. To the right of the workspace is a vertical toolbar with icons for 'SQL', 'Hive', 'Tez', and 'Email'. A red notification badge with the number '11' is visible next to the 'TEZ' icon. At the bottom of the worksheet is a 'New Worksheet' button.

Below the worksheet is a 'Query Process Results' section with a status of 'SUCCEEDED'. It has tabs for 'Logs' (selected) and 'Results'. The results table shows one row:

user_id	num_ratings
405	737

At the very bottom of the screen, a Windows taskbar is visible, showing various application icons and the system clock indicating '04:56 PM 27-11-2024'.

4. Find the most rated movies (with more than 500 ratings)

Query Used:

```
SELECT movie_id, COUNT(*) AS total_ratings
FROM Ratings
GROUP BY movie_id
HAVING total_ratings > 500
ORDER BY total_ratings DESC
LIMIT 5;
```

Output:

The screenshot shows two instances of the Ambari Hive Query Editor interface. Both instances have the same URL: `127.0.0.1:8080/#/main/views/HIVE/1.5.0/AUTO_HIVE_INSTANCE`.

Query Editor (Top Window):

- Database Explorer:** Shows the default database selected.
- Query Editor:** Displays the SQL query:

```
1 SELECT movie_id, COUNT(*) AS total_ratings
2 FROM Ratings
3 GROUP BY movie_id
4 HAVING total_ratings > 500
5 ORDER BY total_ratings DESC
6 LIMIT 5;
```
- Buttons:** Execute, Explain, Save as...

Query Process Results (Bottom Window):

- Logs / Results:** Status: SUCCEEDED
- Results:** A table showing the top 5 most rated movies.

movie_id	total_ratings
50	584
258	509
100	508
181	507

Licensed under the Apache License, Version 2.0.
See third-party tools/resources that Ambari uses and their respective authors

5. Get the top 3 users who have given the highest average ratings

Query Used:

```
SELECT user_id, AVG(rating) AS avg_rating
FROM Ratings
GROUP BY user_id
ORDER BY avg_rating DESC
LIMIT 3;
```

Output:

The screenshot shows a Jupyter Notebook interface with a single worksheet. The code cell contains the following SQL query:

```
1 SELECT user_id, AVG(rating) AS avg_rating
2 FROM Ratings
3 GROUP BY user_id
4 ORDER BY avg_rating DESC
5 LIMIT 3;
```

The results section displays the query process status as "SUCCEEDED" and the resulting data:

user_id	avg_rating
849	4.869565217391305
688	4.833333333333333
507	4.724137931034483

6. Find movies that have ratings from the most number of distinct users

Query Used:

```
SELECT movie_id, COUNT(DISTINCT user_id) AS unique_users
FROM Ratings
GROUP BY movie_id
ORDER BY unique_users DESC
LIMIT 10;
```

Output:

The screenshot shows two windows of the Hortonworks Sandbox interface. The top window is the 'Query Editor' showing the executed SQL query:

```
1 SELECT movie_id, COUNT(DISTINCT user_id) AS unique_users
2 FROM Ratings
3 GROUP BY movie_id
4 ORDER BY unique_users DESC
5 LIMIT 10;
```

The bottom window is the 'Query Process Results' window, which displays the following table:

movie_id	unique_users
50	584
258	509
100	508
181	507
294	485
286	481
288	478
1	452
300	431
121	429

7. Get the average rating for each movie, with user count, and order by movie ID

Query Used:

```
SELECT movie_id, AVG(rating) AS avg_rating, COUNT(*) AS num_ratings
FROM Ratings
GROUP BY movie_id
ORDER BY movie_id ASC;
```

Output:

The screenshot shows two browser windows displaying the same interface. The interface includes a top navigation bar with tabs for 'Hive', 'Query', 'Saved Queries', 'History', 'UDFs', and 'Upload Table'. On the left is a 'Database Explorer' pane showing databases like 'default', 'Ratings', 'sample_07', 'sample_08', 'foodmart', and 'gademo'. The main area is a 'Query Editor' with a 'Worksheet' tab containing the following SQL code:

```
1 SELECT movie_id, AVG(rating) AS avg_rating, COUNT(*) AS num_ratings
2 FROM Ratings
3 GROUP BY movie_id
4 ORDER BY movie_id ASC;
```

Below the worksheet are buttons for 'Execute', 'Explain', and 'Save as...'. To the right of the editor is a vertical toolbar with icons for 'SQL', 'Hive', 'Logs', 'Results', and 'TEZ'. The status bar at the bottom indicates '05:09 PM 27-11-2024'.

The second window shows the 'Query Process Results' with a status of 'SUCCEEDED'. It has tabs for 'Logs' and 'Results'. The 'Results' tab displays a table with three columns: 'movie_id', 'avg_rating', and 'num_ratings'. The data is as follows:

movie_id	avg_rating	num_ratings
1	3.8783185840707963	452
2	3.2061068702290076	131
3	3.033333333333333	90
4	3.550238234449761	209
5	3.302325581955349	86
6	3.576923076923077	26
7	3.798469387755102	392
8	3.9954337899543377	219
9	3.8963210702341136	269
10	3.831460674157303	89
11	3.847457627118644	236
12	4.385767790262173	267
13	3.4184782608695654	184
14	3.9672131147540983	183
15	3.7781569965870307	293
16	3.2051282051282053	39

❖ Conclusion

This project provides a comprehensive end-to-end demonstration of how different tools within the Hadoop ecosystem can work together to manage, store, and analyze large datasets efficiently. By employing Hadoop Distributed File System (HDFS), MySQL, Apache Sqoop, and Apache Hive, it highlights the powerful capabilities of big data technologies in solving complex data challenges, especially in the context of user ratings for movies.

This implementation serves as an excellent demonstration of how the tools in the Hadoop ecosystem can be leveraged in tandem to manage, analyze, and derive insights from large-scale datasets. The ability to seamlessly transfer data across platforms, from GitHub to HDFS, then into MySQL and finally into Hive, underscores the flexibility and scalability of the Hadoop ecosystem.

Key Takeaways:

1. Data Movement and Integration: The use of tools like Sqoop facilitates smooth data movement across different platforms (relational to big data), ensuring that data from various sources can be integrated efficiently and stored in the most suitable environment for processing.
2. Scalable Storage Solutions: By using HDFS for scalable storage and MySQL for structured data management, the project highlights how both systems complement each other. While MySQL handles transactional and structured data, HDFS offers fault tolerance and scalability for larger datasets.
3. Insightful Data Analysis: Hive's SQL-like interface allows for efficient querying of large datasets, enabling data exploration and analysis without the need for complex coding. Its ability to handle petabytes of data across distributed systems makes it an ideal tool for big data analysis.
4. Practical Application: The ability to generate insights from large datasets, such as identifying trends, the highest-rated movies, and active users, demonstrates the practical value of big data tools in real-world applications such as movie recommendations, user profiling, and content popularity analysis.

Final Conclusion:

This project effectively showcases the potential of big data tools to solve real-world data challenges. It not only highlights the technical capabilities of each component (HDFS, MySQL, Sqoop, Hive) but also illustrates the collaborative power of the Hadoop ecosystem. From data ingestion to analysis, the project demonstrates how organizations can leverage big data technologies to handle vast amounts of information, streamline operations, and generate actionable insights for better decision-making. The approach is flexible, scalable, and efficient, making it a robust solution for handling large-scale data analysis in today's data-driven world.

❖ References

<https://www.proquest.com/openview/dcb49cc3706d2057b05f1285783cd77d/1?pq-origsite=gscholar&cbl=2035897>

<https://www.proquest.com/openview/dcb49cc3706d2057b05f1285783cd77d/1?pq-origsite=gscholar&cbl=2035897>

https://books.google.co.in/books?hl=en&lr=&id=ufjJDAAAQBAJ&oi=fnd&pg=PP1&dq=papers+on+hdfs+sqoop+apache+hive&ots=WkwpkyGrBT&sig=KOMT_a2oB6BOXDxsbfEwQc6_Hhc&redir_esc=y#v=onepage&q=papers%20on%20hdfs%20sqoop%20apache%20hive&f=false

<https://dl.acm.org/doi/abs/10.1145/3299869.3314045>

<https://ieeexplore.ieee.org/abstract/document/7453400>

https://books.google.co.in/books?hl=en&lr=&id=IMNiDwAAQBAJ&oi=fnd&pg=PP1&dq=papers+on+apache+hive&ots=BHRdelfJh4&sig=dbYZ2yJxF9h3d2V98oMHDN09DIk&redir_esc=y#v=onepage&q=papers%20on%20apache%20hive&f=false