

MIMIC-III Big Data Pipeline User Manual: Healthcare Analytics with Hadoop, Hive, and MapReduce

Ahmed Mohamed Srouf

May 2025

Prepared by: Ahmed Mohamed Srouf

Contents

1	Overview	4
1.1	Purpose	4
1.2	Scope	4
1.3	Key Components	4
2	System Requirements	4
2.1	Hardware	4
2.2	Software	5
3	Installation and Setup	5
3.1	Prerequisites	5
3.2	Cloning the Repository	5
3.3	Starting Docker Containers	5
3.4	Configuring HDFS and Hive	6
4	Data Preprocessing	6
4.1	Downloading the MIMIC-III Demo Dataset	6
4.2	Data Cleaning and Conversion to Parquet	6
4.3	Loading Parquet Files into HDFS	7
5	Operating the Pipeline	7
5.1	Starting the Pipeline	7
5.2	Creating Hive Tables	8
5.3	Running Hive Analytics	9
5.4	Running MapReduce Jobs	9
5.4.1	Prerequisites	10
5.4.2	Installation and Setup	10
5.4.3	Creating the MapReduce Program	10
5.4.4	Compiling and Packaging the MapReduce Job	12
5.4.5	Transferring Data to HDFS	12
5.4.6	Running the MapReduce Job	12
5.4.7	Viewing the Output	12
6	Troubleshooting	13
6.1	Docker Issues	13
6.2	HDFS Issues	13
6.3	Hive Issues	13
6.4	Data Cleaning Issues	14
6.5	General Troubleshooting Tips	14
7	Technical Specifications	14
7.1	Architecture	14
7.2	Data Flow	15
7.3	Data Flow Diagram	15
8	Backup and Recovery	15
8.1	Data Backup Procedures	15

8.2	Data Recovery Procedures	15
8.3	Maintenance of Backup Process	16
9	Appendices	16
9.1	Glossary	16

1 Overview

This user manual provides comprehensive instructions for setting up, operating, and maintaining a big data pipeline for healthcare analytics using the MIMIC-III Clinical Database Demo v1.4. The pipeline processes key tables (PATIENTS, ADMISSIONS, ICUSTAYS, DIAGNOSES_ICD) to enable batch analytics, such as calculating average length of stay, ICU readmission distribution, and mortality rates by demographic groups. It leverages Hadoop for distributed storage, Hive for SQL-based analytics, MapReduce for custom processing, and Docker for containerized deployment.

1.1 Purpose

The pipeline automates the extraction, transformation, and loading (ETL) of healthcare data to:

- Store and manage MIMIC-III data efficiently in Hadoop Distributed File System (HDFS).
- Perform batch analytics using HiveQL and MapReduce for insights into patient outcomes.
- Enable scalable processing of structured healthcare data in Parquet format.

1.2 Scope

The pipeline processes a subset of the MIMIC-III demo dataset (100 patients) and supports:

- Data cleaning and conversion to Parquet for Hadoop and Hive compatibility.
- Distributed storage and analytics using Hadoop, Hive, and MapReduce.
- Containerized deployment for simplified setup and reproducibility.

1.3 Key Components

- **Hadoop:** Distributed storage and processing via HDFS and MapReduce.
- **Hive:** SQL-based querying for batch analytics.
- **MapReduce:** Custom processing for advanced analytics tasks.
- **Docker:** Containerized environment for Hadoop, Spark, and Hive.
- **Python:** Scripts for data preprocessing and cleaning.
- **Git:** Version control for code and configurations.

2 System Requirements

2.1 Hardware

- **CPU:** Minimum 4 cores (6+ recommended for faster processing).
- **RAM:** 16 GB (32 GB recommended for larger datasets).
- **Storage:** 50 GB free disk space for Docker images, HDFS, and data files.

2.2 Software

- **Operating System:** Linux (Ubuntu 20.04+), macOS, or Windows 10/11 with WSL2.
- **Docker:** Docker Desktop or Engine (version 20.10+).
- **Docker Compose:** Version 1.29+.
- **Git:** Version 2.25+ for repository cloning.
- **Python:** Version 3.8+ with `pandas` and `pyarrow` for data preprocessing.
- **Java:** OpenJDK 8 or 11 for Hadoop, Hive, and MapReduce.

3 Installation and Setup

3.1 Prerequisites

- Install **Docker**: Refer to the [Docker Get Started Guide](#).
- Install **Docker Compose**: See [Compose Installation](#).
- Install **Git**: Follow [Git Installation](#).
- Install **Python**: Ensure Python 3.8+ is installed with `pip`.
- Install **Java**: Ensure OpenJDK 8 or 11 is installed.

3.2 Cloning the Repository

Clone the Dockerized Hadoop, Spark, and Hive environment:

```
1 git clone https://github.com/Marcel-Jan/docker-hadoop-spark.git
2 cd docker-hadoop-spark
```

3.3 Starting Docker Containers

Launch the containers using Docker Compose:

```
1 docker-compose up -d
```

Verify containers are running (e.g., `namenode`, `datanode`, `hive-server`):

```
1 docker ps
```

Note: The environment requires approximately 5-6 GB of storage. An internet connection is needed for initial setup.

3.4 Configuring HDFS and Hive

- **HDFS:** Access the Hadoop NameNode container:

```
1 docker exec -it namenode bash
```

- **Hive:** Access the Hive server and connect to Beeline:

```
1 docker exec -it hive-server bash
2 beeline -u jdbc:hive2://localhost:10000
```

4 Data Preprocessing

4.1 Downloading the MIMIC-III Demo Dataset

Download the MIMIC-III Clinical Database Demo v1.4 from [PhysioNet](#). Extract the following CSV files:

- PATIENTS.csv
- ADMISSIONS.csv
- ICUSTAYS.csv
- DIAGNOSES_ICD.csv

4.2 Data Cleaning and Conversion to Parquet

Clean the dataset to ensure compatibility with Parquet and Hive, including data type conversions and timestamp alignment. Use the following Python script:

```
1 import pandas as pd
2 import pyarrow.parquet as pq
3 import pyarrow as pa
4 from datetime import datetime
5
6 # Define input CSV files
7 csv_files = ['PATIENTS.csv', 'ADMISSIONS.csv', 'ICUSTAYS.csv', '
    DIAGNOSES_ICD.csv']
8
9 for file in csv_files:
10     # Read CSV
11     df = pd.read_csv(file)
12     # Data cleaning
13     if file == 'PATIENTS.csv':
14         # Convert date columns to datetime
15         for col in ['DOB', 'DOD', 'DOD_HOSP', 'DOD_SSN']:
16             if col in df.columns:
17                 df[col] = pd.to_datetime(df[col], errors='coerce')
18         # Ensure EXPIRE_FLAG is integer
19         df['EXPIRE_FLAG'] = df['EXPIRE_FLAG'].astype('int32')
20     elif file == 'ADMISSIONS.csv':
21         # Convert date columns to datetime
```

```

22     for col in ['ADMITTIME', 'DISCHTIME', 'DEATHTIME', 'EDREGTIME', '
23         EDOUTTIME']:
24         if col in df.columns:
25             df[col] = pd.to_datetime(df[col], errors='coerce')
26     # Ensure IDs are integers
27     df['HADM_ID'] = df['HADM_ID'].astype('int32')
28     df['SUBJECT_ID'] = df['SUBJECT_ID'].astype('int32')
29 elif file == 'ICUSTAYS.csv':
30     # Convert date columns to datetime
31     for col in ['INTIME', 'OUTTIME']:
32         if col in df.columns:
33             df[col] = pd.to_datetime(df[col], errors='coerce')
34     # Ensure IDs are integers
35     for col in ['ICUSTAY_ID', 'HADM_ID', 'SUBJECT_ID']:
36         df[col] = df[col].astype('int32')
37 elif file == 'DIAGNOSES_ICD.csv':
38     # Ensure ICD9_CODE is string and handle missing values
39     df['ICD9_CODE'] = df['ICD9_CODE'].fillna('UNKNOWN').astype(str)
40     # Ensure IDs and SEQ_NUM are integers
41     for col in ['HADM_ID', 'SUBJECT_ID', 'SEQ_NUM']:
42         df[col] = df[col].astype('int32')
43 # Convert to Parquet
44 table = pa.Table.from_pandas(df)
pq.write_table(table, file.replace('.csv', '.parquet'))

```

Run the script:

```

1 pip install pandas pyarrow
2 python convert_to_parquet.py

```

4.3 Loading Parquet Files into HDFS

Copy the Parquet files to HDFS:

```

1 docker exec -it namenode bash
2 hdfs dfs -mkdir /mimic_data
3 hdfs dfs -put /path/to/parquet/*.parquet /mimic_data/

```

Verify files in HDFS:

```

1 hdfs dfs -ls /mimic_data

```

5 Operating the Pipeline

5.1 Starting the Pipeline

Ensure Docker containers are running:

```

1 docker-compose up -d

```

5.2 Creating Hive Tables

Access the Hive server and create tables for the Parquet files:

```
1 docker exec -it hive-server bash
2 beeline -u jdbc:hive2://localhost:10000
```

Example Hive table creation scripts:

- PATIENTS:

```
1 CREATE EXTERNAL TABLE patients (
2     subject_id INT,
3     gender STRING,
4     dob TIMESTAMP,
5     dod TIMESTAMP,
6     dod_hosp TIMESTAMP,
7     dod_ssn TIMESTAMP,
8     expire_flag INT
9 )
10 STORED AS PARQUET
11 LOCATION '/mimic_data/PATIENTS.parquet';
```

- ADMISSIONS:

```
1 CREATE EXTERNAL TABLE admissions (
2     subject_id INT,
3     hadm_id INT,
4     admittime TIMESTAMP,
5     dischtime TIMESTAMP,
6     deathtime TIMESTAMP,
7     admission_type STRING,
8     admission_location STRING,
9     discharge_location STRING,
10    insurance STRING,
11    language STRING,
12    religion STRING,
13    marital_status STRING,
14    ethnicity STRING,
15    edregtime TIMESTAMP,
16    edouttime TIMESTAMP,
17    diagnosis STRING,
18    hospital_expire_flag INT
19 )
20 STORED AS PARQUET
21 LOCATION '/mimic_data/ADMISSIONS.parquet';
```

- ICUSTAYS:

```
1 CREATE EXTERNAL TABLE icustays (
2     subject_id INT,
3     hadm_id INT,
4     icustay_id INT,
5     dbsource STRING,
6     first_careunit STRING,
```



```

7      last_careunit STRING,
8      first_wardid INT,
9      last_wardid INT,
10     intime TIMESTAMP,
11     outtime TIMESTAMP,
12     los DOUBLE
13 )
14 STORED AS PARQUET
15 LOCATION '/mimic_data/ICUSTAYS.parquet';

```

- **DIAGNOSES_ICD:**

```

1 CREATE EXTERNAL TABLE diagnoses_icd (
2     subject_id INT,
3     hadm_id INT,
4     seq_num INT,
5     icd9_code STRING
6 )
7 STORED AS PARQUET
8 LOCATION '/mimic_data/DIAGNOSES_ICD.parquet';

```

5.3 Running Hive Analytics

Run HiveQL queries for batch analytics. Example queries:

- **Average Length of Stay per Diagnosis:**

```

1 SELECT d.icd9_code, AVG(DATEDIFF(a.disctime, a.admittime)) AS avg_los
2 FROM admissions a
3 JOIN diagnoses_icd d ON a.hadm_id = d.hadm_id
4 GROUP BY d.icd9_code;

```

- **Distribution of ICU Readmissions:**

```

1 SELECT subject_id, COUNT(DISTINCT hadm_id) AS readmission_count
2 FROM icustays
3 GROUP BY subject_id
4 HAVING readmission_count > 1;

```

- **Mortality Rates by Demographic Groups:**

```

1 SELECT p.gender,
2     COUNT(CASE WHEN p.dod IS NOT NULL THEN 1 END) / CAST(COUNT(*) AS
3     DOUBLE) AS mortality_rate
4 FROM patients p
5 JOIN admissions a ON p.subject_id = a.subject_id
6 GROUP BY p.gender;

```

5.4 Running MapReduce Jobs

This section describes how to implement and run a MapReduce job to calculate the average age of patients in the MIMIC-III dataset using the PATIENTS.csv file.

5.4.1 Prerequisites

- Ensure Java (OpenJDK 8 or 11) is installed in the **namenode** container.
- Verify that the **PATIENTS.csv** file is available for processing.

5.4.2 Installation and Setup

1. Access the **namenode** container:

```
1 docker exec -it namenode bash
```

2. Update the package list and install Java:

```
1 apt update
2 apt install -y openjdk-8-jdk
```

3. Verify Java installation:

```
1 java -version
```

4. Install **nano** if not already present:

```
1 apt update
2 echo "deb http://archive.debian.org/debian stretch main" > /etc/apt/
  sources.list
3 echo "deb http://archive.debian.org/debian-security stretch/updates
  main" >> /etc/apt/sources.list
4 apt-get update
5 apt-get install nano
```

5.4.3 Creating the MapReduce Program

Create a Java program to calculate the average age of patients:

```
1 nano AverageAge.java
```

Paste the following script into **AverageAge.java**:

```
1 import java.io.IOException;
2 import java.text.SimpleDateFormat;
3 import java.util.Date;
4 import java.util.Locale;
5 import java.util.StringTokenizer;
6 import org.apache.hadoop.conf.Configuration;
7 import org.apache.hadoop.fs.Path;
8 import org.apache.hadoop.io.*;
9 import org.apache.hadoop.mapreduce.*;
10 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
11 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
12
13 public class AverageAge {
14
15     public static class AgeMapper extends Mapper<LongWritable, Text, Text,
        IntWritable> {
```

```

16     private final static Text word = new Text("age");
17     private final static SimpleDateFormat format = new
        SimpleDateFormat("yyyy-MM-dd", Locale.ENGLISH);
18
19     public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {
20         String[] fields = value.toString().split(",");
21         if (fields[0].equals("row_id") || fields.length < 5) return;
            // skip header or bad lines
22
23         try {
24             Date dob = format.parse(fields[3]);
25             Date dod = fields[4].isEmpty() ? format.parse("2200-01-01"
                ) : format.parse(fields[4]);
26             long age = (dod.getTime() - dob.getTime()) / (1000L * 60 *
                60 * 24 * 365);
27             if (age > 0 && age < 150) // sanity check
                context.write(word, new IntWritable((int) age));
28         } catch (Exception e) {
29             // Ignore malformed dates
30         }
31     }
32 }
33
34
35 public static class AvgReducer extends Reducer<Text, IntWritable, Text
    , DoubleWritable> {
36     public void reduce(Text key, Iterable<IntWritable> values, Context
        context) throws IOException, InterruptedException {
37         int sum = 0, count = 0;
38         for (IntWritable val : values) {
39             sum += val.get();
40             count++;
41         }
42         if (count != 0)
43             context.write(new Text("Average Age"), new DoubleWritable
                ((double) sum / count));
44     }
45 }
46
47 public static void main(String[] args) throws Exception {
48     Configuration conf = new Configuration();
49     Job job = Job.getInstance(conf, "average age");
50     job.setJarByClass(AverageAge.class);
51     job.setMapperClass(AgeMapper.class);
52     job.setReducerClass(AvgReducer.class);
53     job.setOutputKeyClass(Text.class);
54     job.setOutputValueClass(IntWritable.class);
55     FileInputFormat.addInputPath(job, new Path(args[0]));
56     FileOutputFormat.setOutputPath(job, new Path(args[1]));
57     System.exit(job.waitForCompletion(true) ? 0 : 1);
58 }
59 }

```

Save and exit nano.

5.4.4 Compiling and Packaging the MapReduce Job

1. Create a directory for compiled classes:

```
1 mkdir -p avg_classes
```

2. Set the Hadoop classpath:

```
1 export HADOOP_CLASSPATH=$(hadoop classpath)
```

3. Compile the Java file:

```
1 javac -classpath $HADOOP_CLASSPATH -d avg_classes AverageAge.java
```

4. Package the compiled classes into a JAR file:

```
1 jar -cvf avg.jar -C avg_classes/ .
```

5.4.5 Transferring Data to HDFS

1. Copy the PATIENTS.csv file to the namenode container (from the host machine):

```
1 docker cp /path/to/Patients.csv namenode:/root/
```

2. Create a directory in HDFS:

```
1 hdfs dfs -mkdir -p /user/root/mimic
```

3. Upload the CSV file to HDFS:

```
1 hdfs dfs -put /root/Patients.csv /user/root/mimic/
```

4. Verify the file in HDFS:

```
1 hdfs dfs -ls /user/root/mimic/
```

5.4.6 Running the MapReduce Job

Run the MapReduce job to calculate the average age:

```
1 hadoop jar avg.jar AverageAge /user/root/mimic/Patients.csv /user/root/output_avg
```

Note: If you rerun the job, delete the previous output directory:

```
1 hdfs dfs -rm -r /user/root/output_avg
```

5.4.7 Viewing the Output

Display the result:

```
1 hdfs dfs -cat /user/root/output_avg/part-r-00000
```

6 Troubleshooting

6.1 Docker Issues

- **Containers Fail to Start:**
 - **Error:** Containers exit immediately or fail to initialize.
 - **Solution:**
 - * Check container status: `docker ps -a`.
 - * View logs: `docker logs <container-name>`.
 - * Ensure sufficient memory (16 GB minimum) and no port conflicts (e.g., HDFS ports 9000, 9870; Hive port 10000).
 - * Restart containers: `docker-compose down && docker-compose up -d`.
- **Cannot Access Hive Server:**
 - **Error:** Connection refused on `jdbc:hive2://localhost:10000`.
 - **Solution:**
 - * Verify Hive server is running: `docker ps | grep hive-server`.
 - * Check logs: `docker logs hive-server`.
 - * Ensure network connectivity: `docker network ls`.

6.2 HDFS Issues

- **Failed to Upload Parquet Files:**
 - **Error:** `hdfs dfs -put` fails with permission or connection errors.
 - **Solution:**
 - * Verify HDFS is running: `hdfs dfsadmin -report`.
 - * Check permissions: `hdfs dfs -chmod -R 777 /mimic_data.EnsureNameNodeisaccessible : telnet namenode 9000`.
- **Corrupted Parquet Files:**
 - **Error:** Hive queries fail with schema mismatch or corruption errors.
 - **Solution:**
 - * Regenerate Parquet files using `convert_to_parquet.py.Ensure datatypes match Hive table schema` (e.g.,
 - * Inspect CSV files for malformed data before conversion.

6.3 Hive Issues

- **Query Fails with Schema Mismatch:**
 - **Error:** Column type mismatch or invalid column type.
 - **Solution:**
 - * Ensure Hive table schema matches Parquet file schema.
 - * Regenerate Parquet files with correct data types.

* Drop and recreate table: `DROP TABLE patients; CREATE EXTERNAL TABLE patients....`

- **Slow Query Performance:**

- **Error:** Queries take too long.

- **Solution:**

- * Partition tables by `subject_id, dorhadm_id`. *Increase Hive memory in `hive-site.xml` (`hive.tez.container...`*

6.4 Data Cleaning Issues

- **Timestamp Parsing Errors:**

- **Error:** Invalid timestamp format during Parquet conversion.

- **Solution:**

- * Use `errors='coerce'` in `pd.to_datetime(convert_to_parquet.py)`. *Standardized date formats in CSV files*

- **Data Type Mismatch in Hive:**

- **Error:** Queries fail due to type mismatch (e.g., STRING vs. INT).

- **Solution:**

- * Verify data types in `convert_to_parquet.py` (e.g., `int32` for IDs, `string` for ICD9_CDE). *Update Hive tables*

6.5 General Troubleshooting Tips

- **Check Logs:** Use `docker logs <container-name>` for detailed errors.

- **Restart Services:** Run `docker-compose restart` for transient issues.

- **Update Images:** Run `docker-compose pull` for the latest Docker images.

- **Contact Support:** For unresolved issues, create an issue in the private Git repository or contact the project owner.

7 Technical Specifications

7.1 Architecture

The pipeline uses a containerized architecture with:

- **Hadoop:** Stores MIMIC-III data in HDFS and supports MapReduce.

- **Hive:** Executes SQL-based analytics on Parquet files.

- **MapReduce:** Performs custom analytics, such as calculating average patient age.

- **Docker:** Runs Hadoop, Hive, and supporting services in isolated containers.

7.2 Data Flow

1. **Extraction:** Download MIMIC-III demo CSV files from PhysioNet.
2. **Transformation:** Clean data (data types, timestamps) and convert to Parquet using Python.
3. **Loading:** Upload Parquet files to HDFS.
4. **Analytics:** Create Hive tables, run HiveQL queries, or execute MapReduce jobs.

7.3 Data Flow Diagram

1. MIMIC-III CSV Files → (Python Cleaning and Conversion) →
2. Parquet Files → (HDFS Upload) →
3. HDFS Storage → (Hive Table Creation or MapReduce Job) →
4. Analytics (e.g., length of stay, readmissions, average age).

8 Backup and Recovery

8.1 Data Backup Procedures

- **HDFS Backup:**

- Export HDFS data to local storage:

```
1 hdfs dfs -get /mimic_data /backup/mimic_data_$(date +%F)
```

- Schedule backups using a cron job or manual script.

- **Hive Metadata Backup:**

- Export Hive table definitions:

```
1 docker exec hive-server hive -e 'SHOW CREATE TABLE patients;' > /  
  backup/patients_ddl.sql
```

- Repeat for other tables.

8.2 Data Recovery Procedures

- **HDFS Recovery:**

- Restore data to HDFS:

```
1 hdfs dfs -put /backup/mimic_data_<date> /mimic_data
```

- **Hive Recovery:**

- Recreate tables using backed-up DDL scripts:

```
1 docker exec hive-server hive -f /backup/patients_ddl.sql
```

8.3 Maintenance of Backup Process

- Schedule regular HDFS backups (e.g., weekly).
- Verify backups by restoring to a test environment.
- Define a retention policy (e.g., keep backups for 30 days).

9 Appendices

9.1 Glossary

- **Hadoop**: Open-source framework for distributed storage and processing.
- **HDFS**: Hadoop Distributed File System for large datasets.
- **Hive**: Data warehouse infrastructure for SQL-based querying on Hadoop.
- **MapReduce**: Programming model for processing large datasets in parallel.
- **Parquet**: Columnar storage format optimized for big data.
- **HiveQL**: SQL-like query language for Hive.
- **Docker**: Platform for containerizing applications.
- **MIMIC-III**: Freely accessible critical care database.