# Capstone Project: Instant Health Alert System - Final-Project Submission

## A script to extract patient info using Sqoop into hive table

Introduction:

This document provides an explanation of the sqoop script used to import patient details from remote database

## 1. Installation of MySql connector jar file.

```
wget https://de-mysql-connector.s3.amazonaws.com/mysql-connector-java-8.0.25.tar.gz
```

```
[hadoop@ip-172-31-12-140 ~]$ wget https://de-mysql-connector.s3.amazonaws.com/mysql-connector-java-8.0.25.tar.gz
--2024-11-11 07:27:39--  https://de-mysql-connector.s3.amazonaws.com/mysql-connector-java-8.0.25.tar.gz
Resolving de-mysql-connector.s3.amazonaws.com (de-mysql-connector.s3.amazonaws.com)... 52.217.199.41, 52.217.195.17, 3.5.23.166, ...
Connecting to de-mysql-connector.s3.amazonaws.com (de-mysql-connector.s3.amazonaws.com)|52.217.199.41|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 4079310 (3.9M) [application/x-gzip]
Saving to: 'mysql-connector-java-8.0.25.tar.gz'

100%[===================================>] 4,079,310   --.-K/s   in 0.06s

2024-11-11 07:27:39 (60.0 MB/s) - 'mysql-connector-java-8.0.25.tar.gz' saved [4079310/4079310]

[hadoop@ip-172-31-12-140 ~]$ tar -xvf mysql-connector-java-8.0.25.tar.gz
mysql-connector-java-8.0.25/
mysql-connector-java-8.0.25/src/
mysql-connector-java-8.0.25/src/build/
mysql-connector-java-8.0.25/src/build/java/
mysql-connector-java-8.0.25/src/build/java/documentation/
mysql-connector-java-8.0.25/src/build/java/instrumentation/
mysql-connector-java-8.0.25/src/build/misc/
mysql-connector-java-8.0.25/src/build/misc/debian.in/
mysql-connector-java-8.0.25/src/build/misc/debian.in/source/
mysql-connector-java-8.0.25/src/demo/
mysql-connector-java-8.0.25/src/demo/java/
mysql-connector-java-8.0.25/src/demo/java/demo/
mysql-connector-java-8.0.25/src/demo/java/demo/x/
mysql-connector-java-8.0.25/src/demo/java/demo/x/devapi/
mysql-connector-java-8.0.25/src/generated/
mysql-connector-java-8.0.25/src/generated/java/
mysql-connector-java-8.0.25/src/generated/java/com/
mysql-connector-java-8.0.25/src/generated/java/com/mysql/
mysql-connector-java-8.0.25/src/generated/java/com/mysql/cj/
mysql-connector-java-8.0.25/src/generated/java/com/mysql/cj/x/
mysql-connector-java-8.0.25/src/generated/java/com/mysql/cj/x/protobuf/
mysql-connector-java-8.0.25/src/legacy/
mysql-connector-java-8.0.25/src/legacy/java/
mysql-connector-java-8.0.25/src/legacy/java/com/
```

## 2. Extract of MySql connector tar file.

This section creates an external Hive table named `threshold_reference_table`. External tables are used when data is stored outside of Hive's own data warehouse; in this case, the table is linked to an HBase table, which facilitates real-time access and modifications to the data.

```
tar -xvf mysql-connector-java-8.0.25.tar.gz
```

## 3. Copy the MySQL Connector directory to the Sqoop library to complete the installation.

```
cd mysql-connector-java-8.0.25/
sudo cp mysql-connector-java-8.0.25.jar /usr/lib/sqoop/lib/
```



## 4. Install MySql connector in EMR cluster using the extracted folder

```
mysql_secure_installation
```

Enter current password for root (enter for none): ENTER
Set root password [Y/n] Y
New password: 123
Re-enter password: 123
Remove anonymous users [Y/n] Y
Disallow root login remotely [Y/n] n
Remove test database and access to it [Y/n] Y
Reload privilege tables now [Y/n] Y

## 5. Accessing the MySql shell

After this a prompt for password will be shown hit 123 and grant root privileges and restart mariaDB

```
mysql -u root -p
GRANT ALL PRIVILEGES ON *.* TO 'root'@'%' identified by 'root123' WITH GRANT OPTION;
flush privileges;
exit;
sudo service mariadb restart
```

# Sqoop Commands

## 1. Importing patients information into HDFS

sqoop import --connect
jdbc:mysql://upgraddetest.cyaielc9bmnf.us-east-1.rds.amazonaws.com/testdatabase --table
patients_information --username student --password STUDENT123 --target-dir
/user/hadoop/health-alert/patients-contact-info -m 1

```
                Other local map tasks=1
                Total time spent by all maps in occupied slots (ms)=161328
                Total time spent by all reduces in occupied slots (ms)=0
                Total time spent by all map tasks (ms)=3361
                Total vcore-milliseconds taken by all map tasks=3361
                Total megabyte-milliseconds taken by all map tasks=5162496
        Map-Reduce Framework
                Map input records=5
                Map output records=5
                Input split bytes=87
                Spilled Records=0
                Failed Shuffles=0
                Merged Map outputs=0
                GC time elapsed (ms)=67
                CPU time spent (ms)=1890
                Physical memory (bytes) snapshot=261730304
                Virtual memory (bytes) snapshot=3281002496
                Total committed heap usage (bytes)=247463936
        File Input Format Counters
                Bytes Read=0
        File Output Format Counters
                Bytes Written=230
23/03/25 07:11:39 INFO mapreduce.ImportJobBase: Transferred 230 bytes in 20.9571
 seconds (10.9748 bytes/sec)
23/03/25 07:11:39 INFO mapreduce.ImportJobBase: Retrieved 5 records.
```

## 2. Creating Hive Table `patients_contact_info` for the imported data

2.1 Create a database health

```
create database health;
use health;
```



```
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> create database health;
OK
Time taken: 0.48 seconds
hive> use health;
OK
Time taken: 0.02 seconds
hive> 
```

### 2.2 Create a table **patients_contact_info**

```
CREATE EXTERNAL TABLE IF NOT EXISTS health.patients_contact_info (
    patientid int,
    patientname string,
    patientaddress string,
    phone_number string,
    admitted_ward int,
    age int,
    other_details string
)
row format delimited
fields terminated by ','
lines terminated by '\n'
location '/user/hadoop/health-alert/patients-contact-info';
```

### 2.3 View the imported data

```
Select * from  patients_contact_info
```

```
hive> select * from Patients_Contact_Info;
OK
patients_contact_info.patientid patients_contact_info.patientname       patients
_contact_info.patientaddress    patients_contact_info.phone_number      patients
_contact_info.admitted_ward     patients_contact_info.age       patients_contact
_info.other_details
1       Alex S  XDC test Address        8982739282      1       23      null
2       Sammy A New Building Address    2382739282      2       45      null
3       Karan C Aws Address     8923739282      3       56      null
4       Dara M  India Address   2182739282      4       67      null
5       Pam     ABC test Address        4982739282      5       72      null
Time taken: 1.541 seconds, Fetched: 5 row(s)
```