# Apache Sqoop

By Sumit Mittal

# Apache Sqoop Exercise 1

# IMPORTANT
## Copyright Infringement and Illegal Content Sharing Notice

# Sqoop Basics

## To enter into MySQL:

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

```
                                    cloudera@quickstart:~        _  □  ×
                           File  Edit  View  Search  Terminal  Help
                           [cloudera@quickstart ~]$ mysql -u root -p
                           Enter password: █
```

**Note:**

Enter password: cloudera

## MySQL root user:

root user has acess to all the databases.
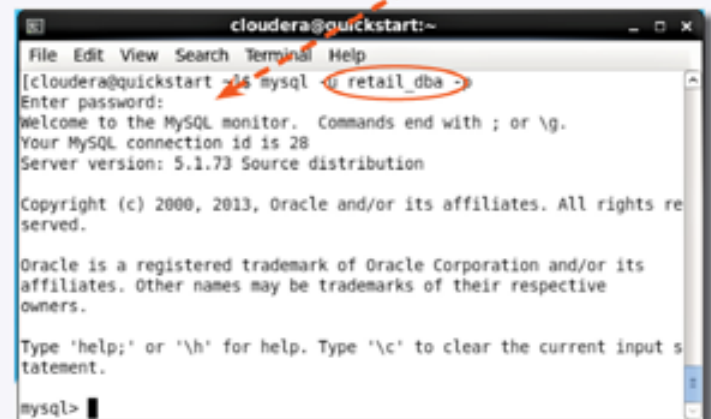
```
mysql -u root -p

(Enter password: cloudera)
```

```
                        cloudera@quickstart:~              _  □  ×
File  Edit  View  Search  Terminal  Help
[cloudera@quickstart ~]$ mysql -u root -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 20
Server version: 5.1.73 Source distribution

Copyright (c) 2000, 2013, 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> █
```

## MySQL reatil_dba user:

retail_dba user has acess to limited  databases.

```
mysql -u retail_dba -p

(Enter password: cloudera)
```

```
                        cloudera@quickstart:~              _  □  ×
File  Edit  View  Search  Terminal  Help
[cloudera@quickstart ~]$ mysql -u retail_dba -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 28
Server version: 5.1.73 Source distribution

Copyright (c) 2000, 2013, Oracle and/or its affiliates. All rights re
served.

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 s
tatement.

mysql> █
```
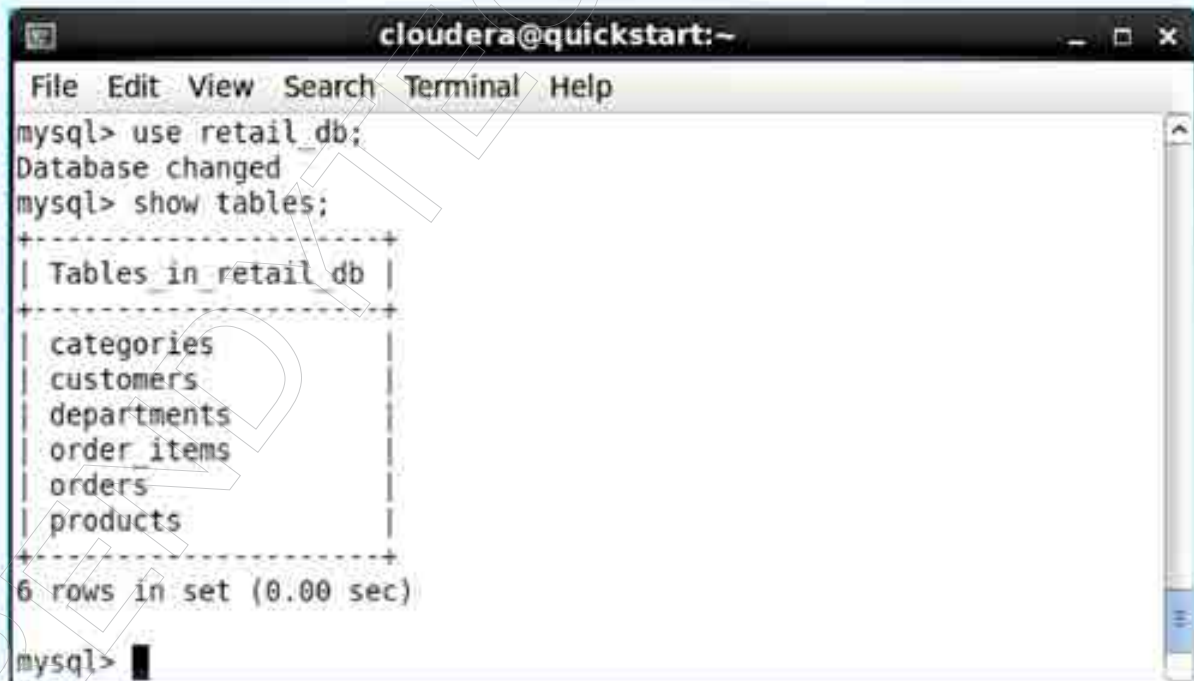
# To display databases in MySQL:

```
show databases;
```



# Use databases and display tables:
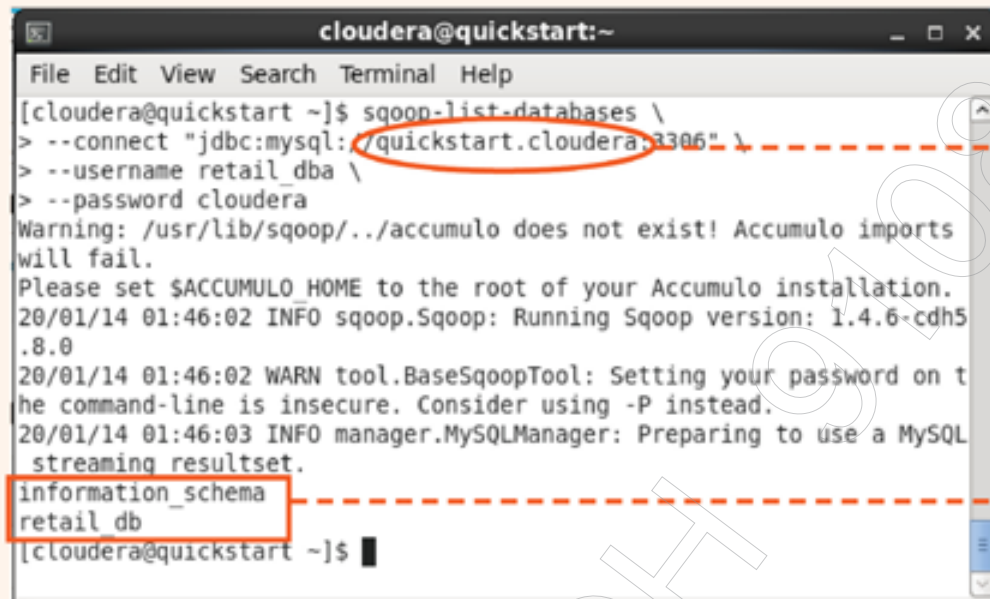
```
use retail_db;


show tables;
```

## Acessenig MySQL databses from Hadoop using Sqoop:

```
sqoop-list-databases \
--connect "jdbc:mysql://quickstart.cloudera:3306" \
--username retail_dba \
--password cloudera
```

**Space with backslash (\) Indicates continuation of line**


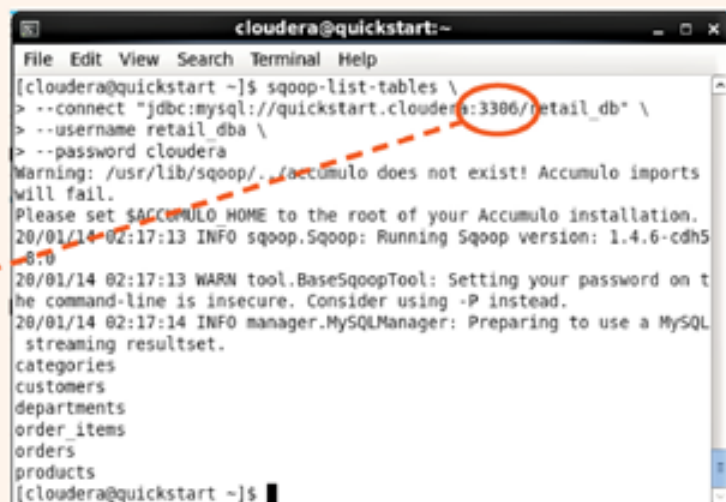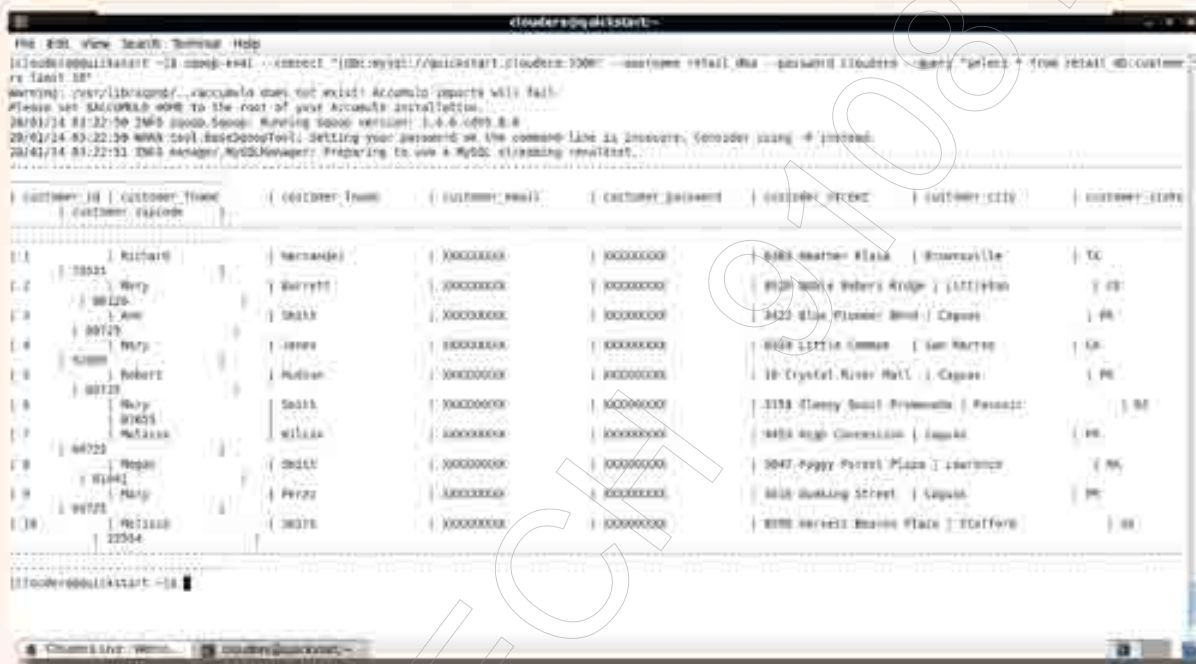
**Local host name**

**List of databese**

## Acessenig MySQL tables using the *root* user:

```
sqoop-list-tables \
--connect "jdbc:mysql://quickstart.cloudera:3306/retail_db" \
--username retail_dba \
--password cloudera
```

**Local port no. where MySQL runs**

TRENDYTECH

## Displaying table data using sqoop-eval:

```
sqoop-eval \
--connect "jdbc:mysql://quickstart.cloudera:3306" \
--username retail_dba \
--password cloudera \
--query "select * from retail_db.customers limit 10"
```



## Create and use a database in MySQL:

```
CREATE database trendytech;



USE trendytech;
```

## Create a table in MySQL:

```
CREATE TABLE people
(
PersonID int,
LastName varchar(255),
FirstName varchar(255),
Address varchar(255),
City varchar(255)
);
```
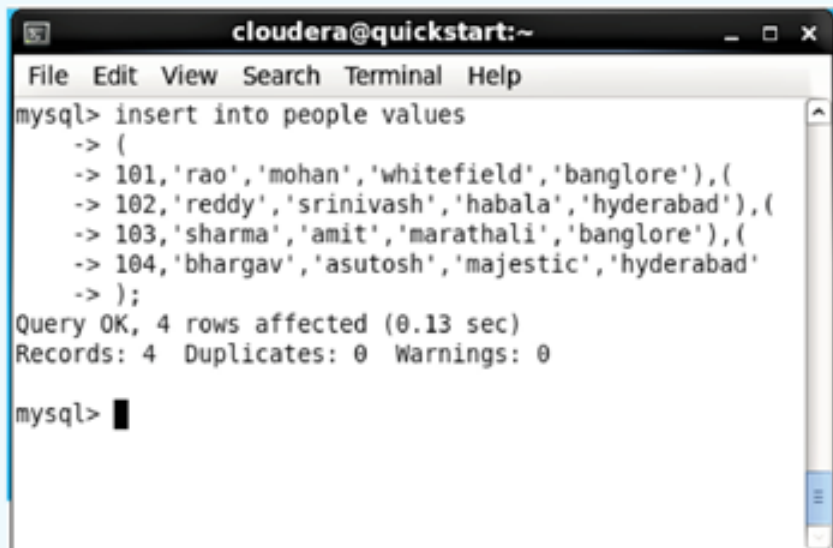


```
mysql> CREATE TABLE people (
    -> PersonID int,
    -> LastName varchar(255),
    -> FirstName varchar(255),
    -> Address varchar(255),
    -> City varchar(255)
    -> );
Query OK, 0 rows affected (0.35 sec)

mysql>
```

## Insert records into the people table:

```
insert into people values
(
101,'rao','mohan','whitefield','banglore'),(
102,'reddy','srinivash','habala','hyderabad'),(
103,'sharma','amit','marathali','banglore'),(
104,'bhargav','asutosh','majestic','hyderabad'
);



commit;
```



```
mysql> insert into people values
    -> (
    -> 101,'rao','mohan','whitefield','banglore'),(
    -> 102,'reddy','srinivash','habala','hyderabad'),(
    -> 103,'sharma','amit','marathali','banglore'),(
    -> 104,'bhargav','asutosh','majestic','hyderabad'
    -> );
Query OK, 4 rows affected (0.13 sec)
Records: 4  Duplicates: 0  Warnings: 0

mysql>
```
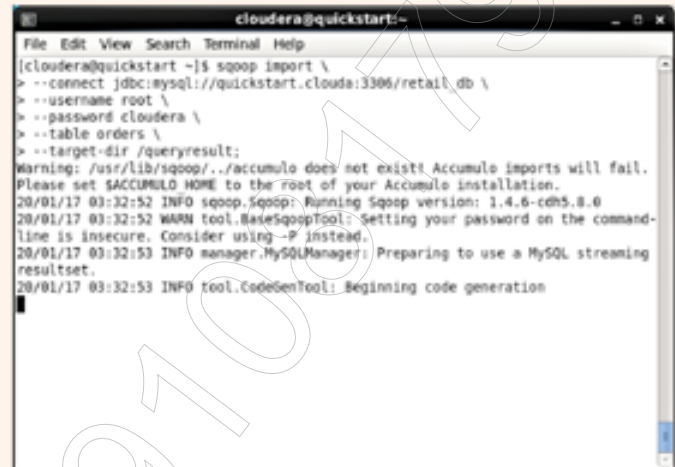
TRENDY T ECH

# Import data from MySQL to Sqoop:

```
sqoop import \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username root \
--password cloudera \
--table orders \
--target-dir /queryresult
```
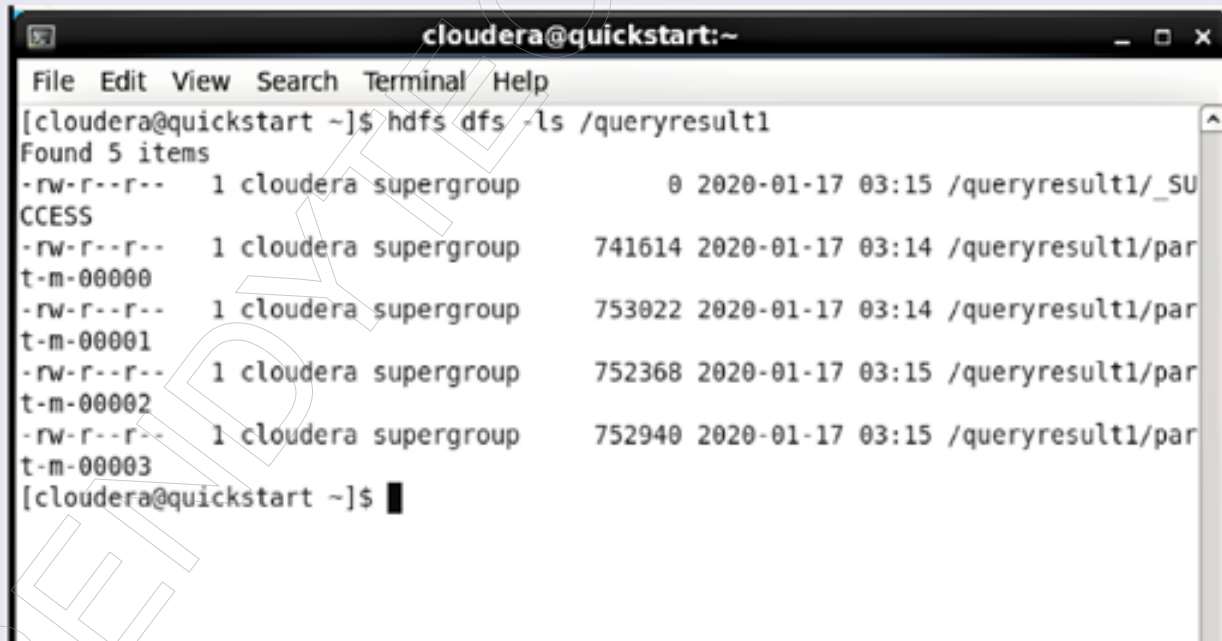


**Note**: If table don't have primary key than it will not import.

# To display contents of queryresult directory in HDFS (use terminal):

```
hadoop fs -ls /queryresult
```



**Note**: By default the number of mappers are 4, so 4 output files are created.

**Instructions** ► Import the **people** table (which we have created earlier in MySQL) with same command as we did above.

## To import people table from MySQL to HDFS:

```
sqoop import \
--connect jdbc:mysql://quickstart.cloudera:3306/trendytech \
--username root --password cloudera --table people \
--target-dir /peopleresult
```
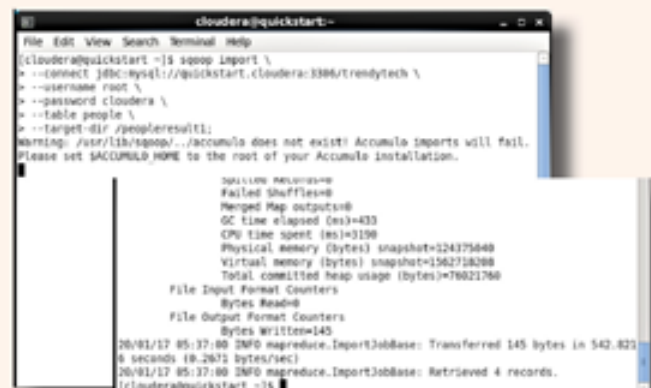


**NOTE**: it will throw error. Becasue **people** table doesn't have primary key.

**Instructions** ► Now, run the above command with mapper (**-m 1**):

## To import people table from MySQL to HDFS with one Mapper:

```
sqoop import \
--connect jdbc:mysql://quickstart.cloudera:3306/trendytech \
--username root \
--password cloudera \
--table people \
-m 1 \
--target-dir /peopleresult1
```

# To display people table from HDFS:

```
hadoop fs -ls /peopleresult1

hadoop fs -cat /peopleresult1/*
```



**Note:** You will find one mapper file only (**part-m-00000**).

# To import all tables from "MySQL" database:

```
sqoop-import-all-tables \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username retail_dba \
--password cloudera \
--as-sequencefile \     <------------------------- File Format
-m 4 \
--warehouse-dir /user/cloudera/sqoopdir
```

**Note:** Here no of mappers are 4 that means we will get 4 files

We can also mention **file format** while importing data as mentioned above.

Sqoop supports **4 types** of file formats:

- Text file format
- Sequence file format
- Avro file format
- Parquet file format

**Note**: If you do not mention any file format, by default it will be text file format.

By default Sqoop provides 4 mappers - so we can skip the above **–m  4** command and still get the same result.

Difference between Sqoop **target** directory & **warehouse** directory.

The difference is that:

–target-dir is a full directory path and the data files will be created directly inside the specified folder.

–warehouse-dir is used to specify a base directory within hdfs where SQOOP will create a sub folder inside with the name of the source table, and import the data files into that folder.

Directory structure for **retail_db** will be:

/user/cloudera/sqoopdir/employee
/user/cloudera/sqoopdir/customer
/user/cloudera/sqoopdir/table3
/user/cloudera/sqoopdir/tablw4

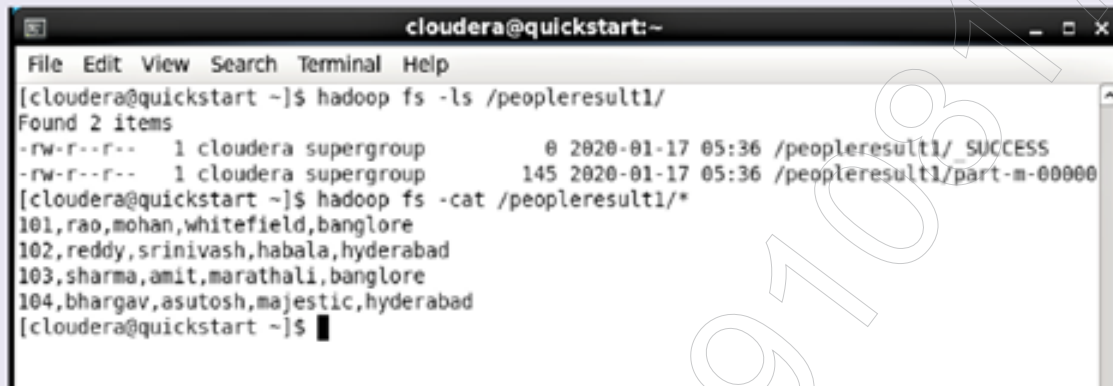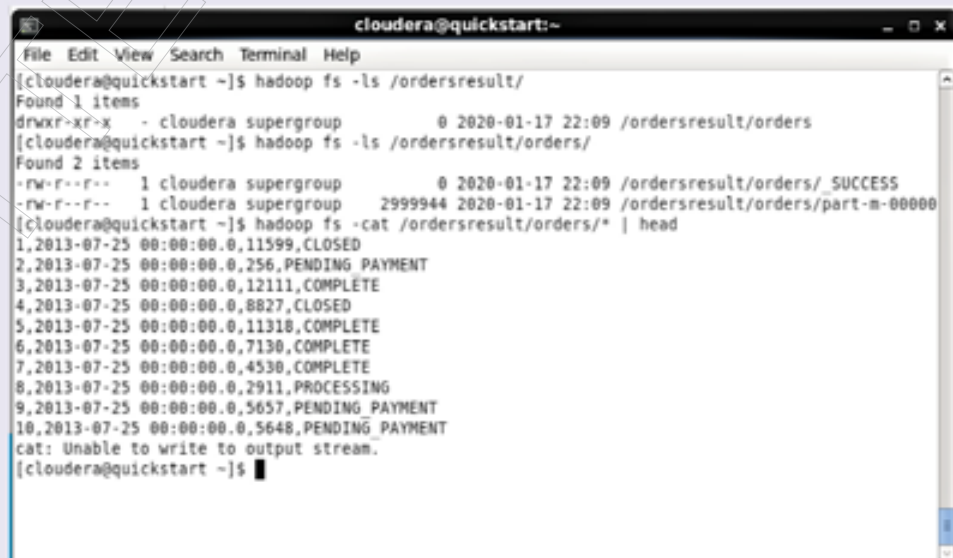## Now try to run following code to import the *orders* table with `--warehouse-dir` path:

```
sqoop import \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username root \
--password cloudera \
--table orders \
--warehouse-dir /ordersresult
```

## To check the file structure in HDFS:

```
hadoop fs -ls /ordersresult/
```

/user/cloudera/warehouse/ordersresult/orders

part-m-00000_0

# To display a list of all available tools:

```
sqoop help
```

```
                                cloudera@quickstart:~                          _  □ ×

 File  Edit  View  Search  Terminal  Help

[cloudera@quickstart ~]$ sqoop help
Warning: /usr/lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
20/01/17 22:49:12 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.8.0
Usage: sqoop COMMAND [ARGS]

Available commands:
  codegen            Generate code to interact with database records
  create-hive-table  Import a table definition into Hive
  eval               Evaluate a SQL statement and display the results
  export             Export an HDFS directory to a database table
  help               List available commands
  import             Import a table from a database to HDFS
  import-all-tables  Import tables from a database to HDFS
  import-mainframe   Import datasets from a mainframe server to HDFS
  job                Work with saved jobs
  list-databases     List available databases on a server
  list-tables        List available tables in a database
  merge              Merge results of incremental imports
  metastore          Run a standalone Sqoop metastore
  version            Display version information

See 'sqoop help COMMAND' for information on a specific command.
[cloudera@quickstart ~]$ █
```

# To know sqoop version:

```
sqoop version
```

```
                                cloudera@quickstart:~                          _  □ ×

 File  Edit  View  Search  Terminal  Help

[cloudera@quickstart ~]$ sqoop version
Warning: /usr/lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
20/01/17 23:01:50 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.8.0
Sqoop 1.4.6-cdh5.8.0
git commit id
Compiled by jenkins on Thu Jun 16 12:25:21 PDT 2016
[cloudera@quickstart ~]$ █
```

# Sqoop help with command Aliases

```
sqoop help eval
```



```
sqoop help import
```

The argument **--password** takes authentication password in plain text.

```
sqoop-list-databases \
--connect jdbc:mysql://quickstart.cloudera:3306 \
--username retail_dba \
--password cloudera
```



While the argument **-P** read password from console.

```
sqoop-list-databases \
--connect jdbc:mysql://quickstart.cloudera:3306 \
--username retail_dba \
-P
```

The argument **--query** can be replaced with **-e**.

```
sqoop-eval \
--connect jdbc:mysql://quickstart.cloudera:3306 \
--username retail_dba \
--password cloudera \
--query "select * from retail_db.customers limit 10"
```

OR

```
sqoop-eval \
--connect "jdbc:mysql://quickstart.cloudera:3306" \
--username retail_dba \
--password cloudera \
-e "select * from retail_db.customers limit 10"
```

Similarly **-m** and **--num-mappers** are same.

```
sqoop import \
--connect jdbc:mysql://quickstart.cloudera:3306/trendytech \
--username root \
--password cloudera \
--table people -m 1 \
--target-dir /peopleresult1
```

OR

```
sqoop import \
--connect jdbc:mysql://quickstart.cloudera:3306/trendytech \
--username root \
--password cloudera \
--table people --num-mappers 1 \
--target-dir /peopleresult1
```

## Redirecting logs:

```
sqoop import \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username root \
--password cloudera \
--table orders \
--warehouse-dir /queryresult4 1>query.out 2>query.err
```

```
cloudera@quickstart:~                                    _ □ ✗

File  Edit  View  Search  Terminal  Help

[cloudera@quickstart ~]$ sqoop import \
> --connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
> --username root \
> --password cloudera \
> --table orders \
> --warehouse-dir /queryresult4 1>query.out 2>query.err
[cloudera@quickstart ~]$ ▮
```

## To check the content of the queryresult4:

```
hadoop fs -ls /queryresult4/orders/
```

```
cloudera@quickstart:~                                    _ □ ✗

File  Edit  View  Search  Terminal  Help

[cloudera@quickstart ~]$ hadoop fs -ls /queryresult4/
Found 1 items
drwxr-xr-x   - cloudera supergroup          0 2020-01-20 20:51 /queryresult4/orders
[cloudera@quickstart ~]$ hadoop fs -ls /queryresult4/orders/
Found 5 items
-rw-r--r--   1 cloudera supergroup          0 2020-01-20 20:51 /queryresult4/orders/_SUCCESS
-rw-r--r--   1 cloudera supergroup     741614 2020-01-20 20:50 /queryresult4/orders/part-m-00000
-rw-r--r--   1 cloudera supergroup     753022 2020-01-20 20:50 /queryresult4/orders/part-m-00001
-rw-r--r--   1 cloudera supergroup     752368 2020-01-20 20:51 /queryresult4/orders/part-m-00002
-rw-r--r--   1 cloudera supergroup     752940 2020-01-20 20:51 /queryresult4/orders/part-m-00003
[cloudera@quickstart ~]$ ▮
```

# To check the contents of log files:

```
cat query.out
```

```
cloudera@quickstart:~
File  Edit  View  Search  Terminal  Help
[cloudera@quickstart ~]$ cat query.out
Warning: /usr/lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
[cloudera@quickstart ~]$
```

```
cat query.err
```

```
cloudera@quickstart:~
File  Edit  View  Search  Terminal  Help
[cloudera@quickstart ~]$ cat query.err
20/01/20 20:49:05 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.13.0
20/01/20 20:49:05 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P
 instead.
20/01/20 20:49
20/01/20 20:49          cloudera@quickstart:~
20/01/20 20:49  File  Edit  View  Search  Terminal  Help
20/01/20 20:49          HDFS: Number of bytes written=2999944
20/01/20 20:49          HDFS: Number of read operations=10
Note: /tmp/sqo          HDFS: Number of large read operations=0
Note: Recompil          HDFS: Number of write operations=0
20/01/20 20:49     Job Counters
875a75c8e6/ord          Launched map tasks=4
20/01/20 20:49          Other local map tasks=4
20/01/20 20:49          Total time spent by all maps in occupied slots (ms)=195315
20/01/20 20:49          Total time spent by all reduces in occupied slots (ms)=0
20/01/20 20:49          Total time spent by all map tasks (ms)=195315
20/01/20 20:49          Total vcore-milliseconds taken by all map tasks=195315
20/01/20 20:49          Total megabyte-milliseconds taken by all map tasks=200002560
er.address         Map Reduce Framework
20/01/20 20:49          Map input records=68883
20/01/20 20:49          Map output records=68883
20/01/20 20:49          Input split bytes=469
20/01/20 20:49          Spilled Records=0
20/01/20 20:49          Failed Shuffles=0
'orders'                Merged Map outputs=0
20/01/20 20:49          GC time elapsed (ms)=694
20/01/20 20:49          CPU time spent (ms)=24730
20/01/20 20:49          Physical memory (bytes) snapshot=620457472
20/01/20 20:49          Virtual memory (bytes) snapshot=6300893184
20/01/20 20:49          Total committed heap usage (bytes)=392691712
                   File Input Format Counters
                        Bytes Read=0
                   File Output Format Counters
                        Bytes Written=2999944
20/01/20 20:51:08 INFO mapreduce.ImportJobBase: Transferred 2.861 MB in 102.9298 seconds (28.4624 KB/sec)
20/01/20 20:51:08 INFO mapreduce.ImportJobBase: Retrieved 68883 records.
[cloudera@quickstart ~]$
```

# Sqoop import execution flow

How Mappers devide their work when a query fired:

- Selects 1 record and by using that it gets the metadata and builds the java file



- Using above java file it builds the jar file



- **BoundingValsQuery** based on min and max on primary key



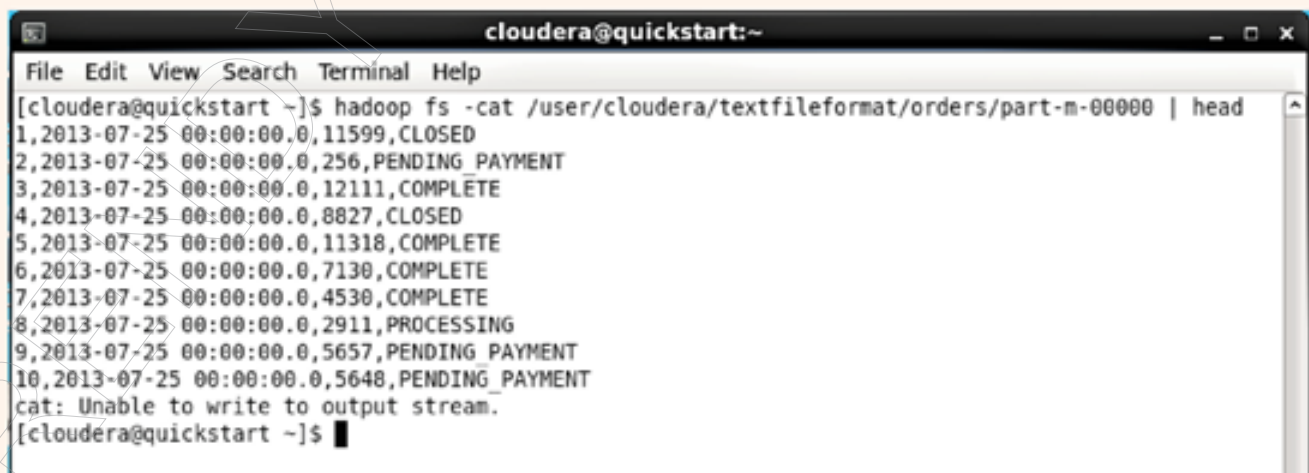- Calculates (max - min)/4 and it gets the split size.

## File formats:

Sqoop import supports following file formats:

1. Text file format - command argument **--as-textfile**
2. Sequence file format - command argument **--as-sequencefile**
3. Avro file format - command argument **--as-avrodatafile**
4. Parquet file format - command argument **--as-parquetfile**

**Note**:  If you are not mentioning any file format, by default  sqoop
uses **--as-textfile**

## Text file format:

```
sqoop-import \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username retail_dba \
--password cloudera \
--table orders \
--as-textfile \
-m 4 \
--warehouse-dir /user/cloudera/textfileformat
```

```
cloudera@quickstart:~                                    _ □ x

File  Edit  View  Search  Terminal  Help
[cloudera@quickstart ~]$ hadoop fs -cat /user/cloudera/textfileformat/orders/part-m-00000 | head
1,2013-07-25 00:00:00.0,11599,CLOSED
2,2013-07-25 00:00:00.0,256,PENDING_PAYMENT
3,2013-07-25 00:00:00.0,12111,COMPLETE
4,2013-07-25 00:00:00.0,8827,CLOSED
5,2013-07-25 00:00:00.0,11318,COMPLETE
6,2013-07-25 00:00:00.0,7130,COMPLETE
7,2013-07-25 00:00:00.0,4530,COMPLETE
8,2013-07-25 00:00:00.0,2911,PROCESSING
9,2013-07-25 00:00:00.0,5657,PENDING_PAYMENT
10,2013-07-25 00:00:00.0,5648,PENDING_PAYMENT
cat: Unable to write to output stream.
[cloudera@quickstart ~]$ █
```

## Sequence file format:

```
sqoop-import \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username retail_dba \
--password cloudera \
--table orders \
--as-sequencefile \
-m 4 \
--warehouse-dir /user/cloudera/sequencefileformat
```
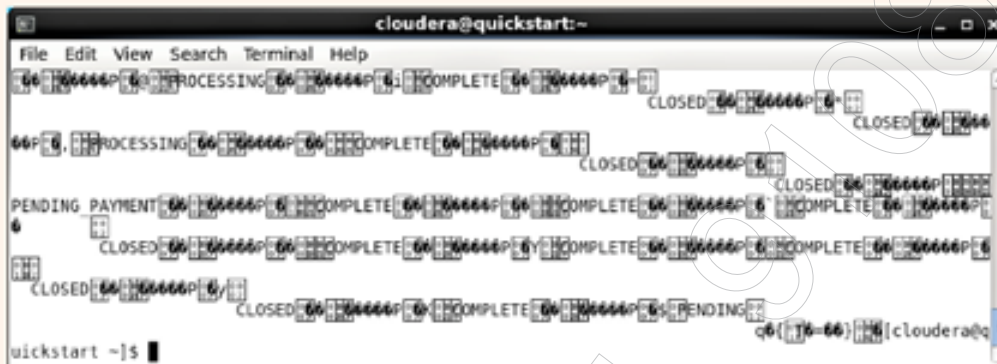


**Note**: SequenceFiles are a binary format that store individual records in custom record-specific data types. These data types are manifested as Java classes.

TRENDY**T**ECH

## Avro file format:

```
sqoop-import \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username retail_dba --password cloudera \
--table orders \
--as-avrodatafile -m 4 \
--warehouse-dir /user/cloudera/avrofileformat
```



## Parquet file format:

```
sqoop-import \
--connect jdbc:mysql://quickstart.cloudera:3306/retail_db \
--username retail_dba --password cloudera \
--table orders \
--as-parquetfile -m 4 \
--warehouse-dir /user/cloudera/parquetfileformat
```

# Follow US

| | |
|---|---|
| **Trainer** | **Mr. Sumit Mittal** |
| **LinkedIn** | https://www.linkedin.com/in/bigdatabysumit/ |
| **Website** | https://trendytech.in/courses/big-data-online-training/ |
| **Phone** | 9108179578 |
| **Email** | trendytech.sumit@gmail.com |
| **Youtube** | TrendyTech |
| **Twitter** | @BigdataBySumit |
| **Instagram** | bigdatabysumit |
| **Facebook** | https://www.facebook.com/trendytech.in/ |

TRENDYTECH
UPLIFT YOUR CAREER!