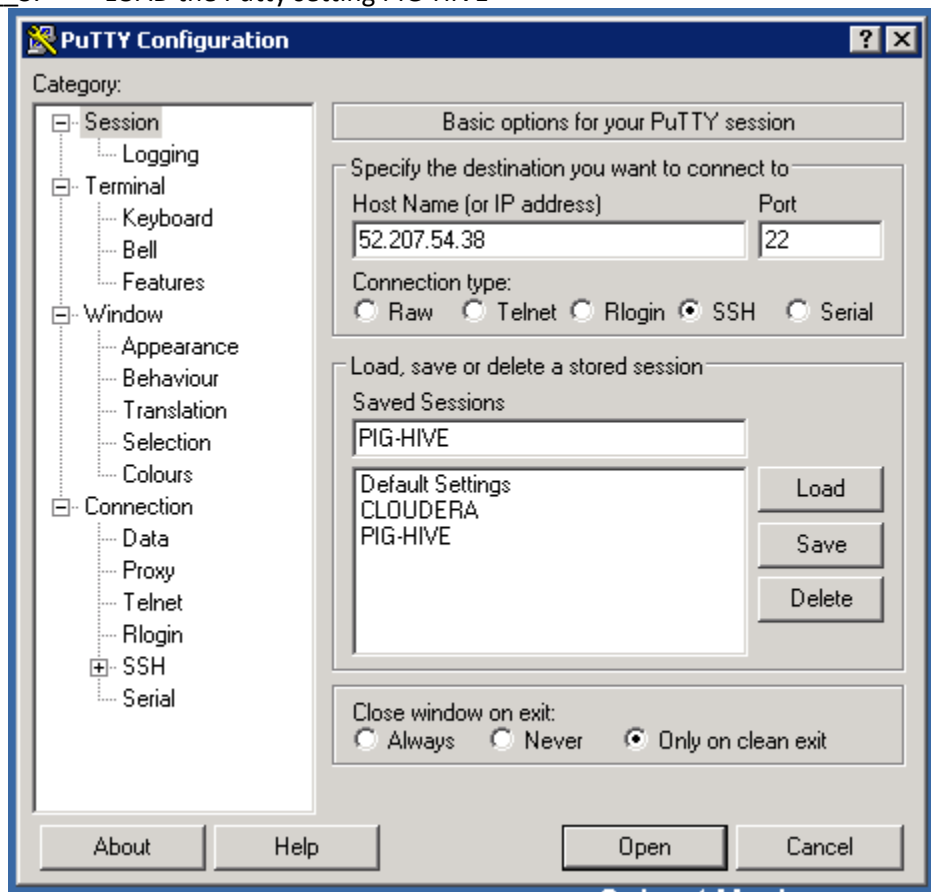


## Lab 2 – Introduction to Hive

In this lab we will demonstrate the use of Impala within the Linux shell and also within the HUE UI of Cloudera

- \_\_\_1. Go to the desktop of your Windows VM
- \_\_\_2. Double click on the Putty shortcut
- \_\_\_3. LOAD the Putty setting PIG-HIVE



- \_\_\_4. The password to login to the Linux environment : cloudera

```
cloudera@ip-172-31-88-237:~  
Using username "cloudera".  
cloudera@52.207.54.38's password:  
Last login: Tue Sep 12 18:48:49 2017 from 71-84-92-219.dhcp.rvsd.ca.charter.com  
[cloudera@ip-172-31-88-237 ~]$
```

- \_\_\_\_5. The dataset we are using here is from Flickserv. Flickserv is a Netflix Search Engine. The dataset is a simple text(Hadoop\_movies\_data.csv) file that lists movie names and its details like release year, rating, and runtime. To view the dataset within the hdfs, please type :
- \$ hadoop fs -ls

```
cloudera@ip-172-31-88-237:~  
Using username "cloudera".  
cloudera@52.207.54.38's password:  
Last login: Tue Sep 12 18:48:49 2017 from 71-84-92-219.dhcp.rvsd.ca.charter.com  
[cloudera@ip-172-31-88-237 ~]$ hadoop fs -ls  
Found 3 items  
drwxr-xr-x - cloudera cloudera 0 2017-09-12 18:42 .Trash  
-rw-r--r-- 1 cloudera cloudera 2940597 2017-09-12 18:00 hadoop_movies_data.csv  
drwxr-xr-x - cloudera cloudera 0 2017-09-12 19:14 oozie-oozi  
[cloudera@ip-172-31-88-237 ~]$
```

- \_\_\_\_6. To start hive, type : \$ hive

```

Details at logfile: /home/cloudera/pig_1505272914686.log
grunt> quit;
[cloudera@ip-172-31-88-237 ~]$ hive

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties
hive>

```

\_\_\_\_7. Please type in:

```
CREATE TABLE IF NOT EXISTS hadoop_movies_data
```

```
( id int, name String, year int, rating double, duration int)
```

```
COMMENT 'Movies details' ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LINES TERMINATED BY
'\n' STORED AS TEXTFILE;
```

```

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties
hive> CREATE TABLE IF NOT EXISTS hadoop_movies_data ( id int, name String, year int, rating d
ouble, duration int)
> COMMENT 'Movies details' ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LINES TERMINATED
BY '\n' ;
OK
Time taken: 1.184 seconds
hive>

```

\_\_\_\_8. Generally, after creating a table in SQL, we can insert data using the Insert statement. But in Hive, we can insert data using the LOAD DATA statement.

While inserting data into Hive, it is better to use LOAD DATA to store bulk records. There are two ways to load data: one is from local file system and second is from Hadoop file system.

## Syntax

The syntax for load data is as follows:

```
LOAD DATA [LOCAL] INPATH 'filepath' [OVERWRITE] INTO TABLE tablename
[PARTITION (partcol1=val1, partcol2=val2 ...)]
```

- LOCAL is identifier to specify the local path. It is optional.
- OVERWRITE is optional to overwrite the data in the table.
- PARTITION is optional.

Check your File Browser in HUE to see if `hadoop_movies_data.csv` is located in `/user/cloudera`  
If not you can upload the file, in the `C:\movies` folder on the Windows VM

The following query loads the given text into the table.

```
hive> LOAD DATA INPATH 'hadoop_movies_data.csv' OVERWRITE INTO TABLE  
Hadoop_movies_data;
```

Now the location of the file is located in the hdfs under the `/user/hive/warehouse/`

To verify exit out of hive with

```
hive> quit;
```

```
$ hadoop fs -ls /user/hive/warehouse
```

You should see the `hadoop_movies_data` folder

\_\_\_\_9. Now we can re-enter hive

```
$hive
```

The metadata tables remain there in hive

\_\_\_\_10. Let's run some queries

\_\_\_\_11. The following query is executed to retrieve the employee details using the above table:

```
SELECT * FROM hadoop_movies_data WHERE rating > 4.0;
```

\_\_\_\_12. *List the movies that were released between 1950 and 1960*

```
hive> SELECT * FROM hadoop_movies_data WHERE year > 1950 and year < 1960;
```

\_\_\_13. Enter the following in a web browser:

http://<ip-address supplied on desktop>:8888

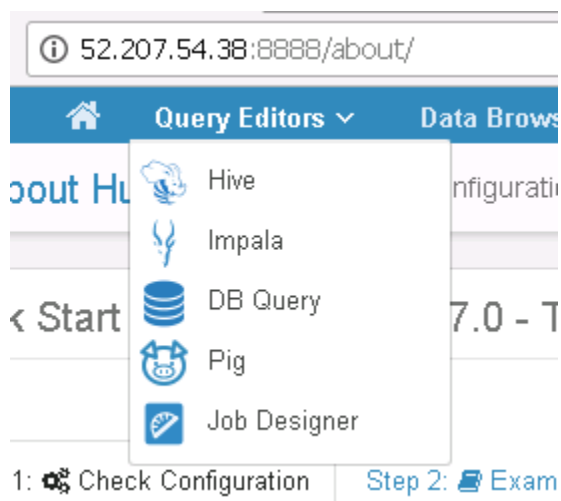
username: cloudera

password: cloudera

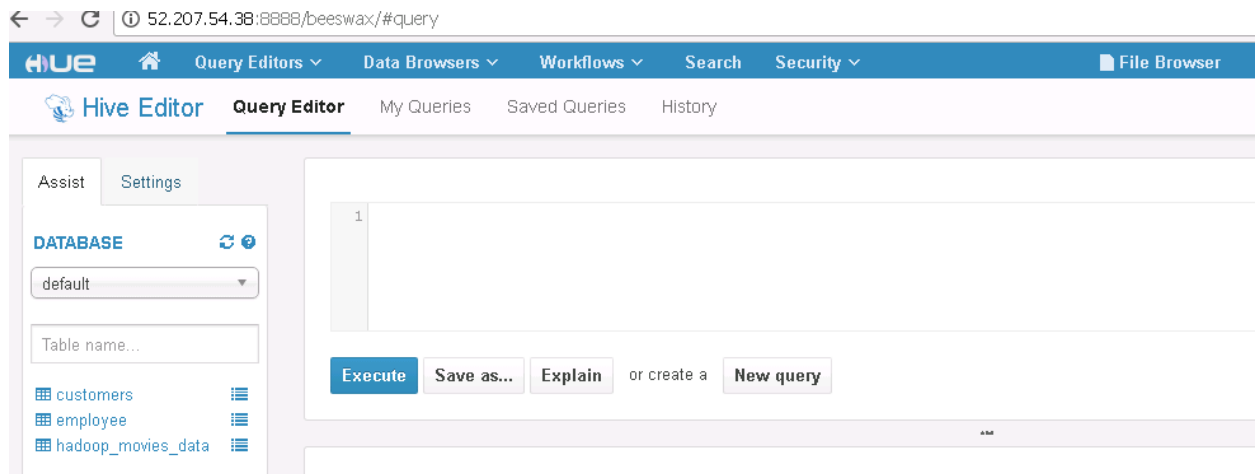


Sign in to continue to Hue

\_\_\_14. Go to Query Editors -> Hive



\_\_\_15. You are now in the beeswax

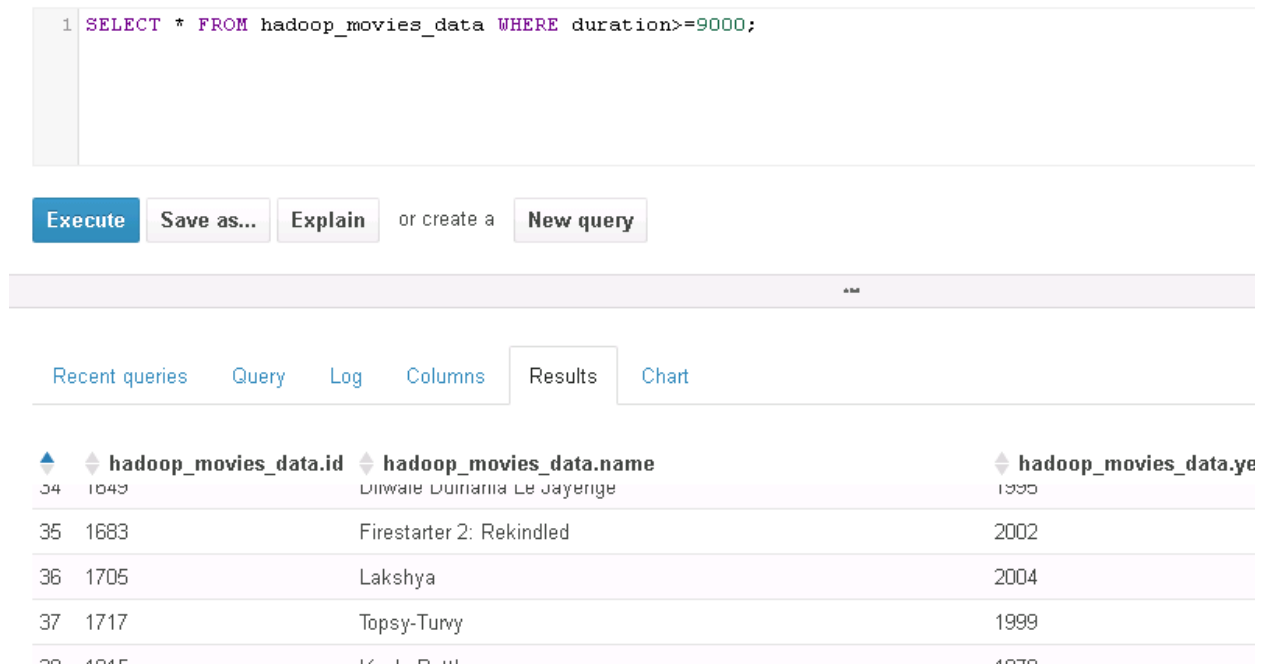


Take a look at the tables on the default database.

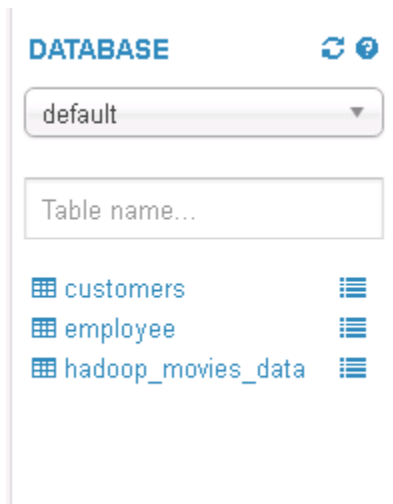
\_\_\_16. Now enter a query within the Query Editor

**SELECT \* FROM hadoop\_movies\_data WHERE duration>=9000;**

While processing, click on the Results Tab to see the results



\_\_\_17. You can click on the tables on the left to see details



\_\_\_18. Save your query as Query-movies, and then view your saved query by clicking on My Queries. You can also view the history there too.

