

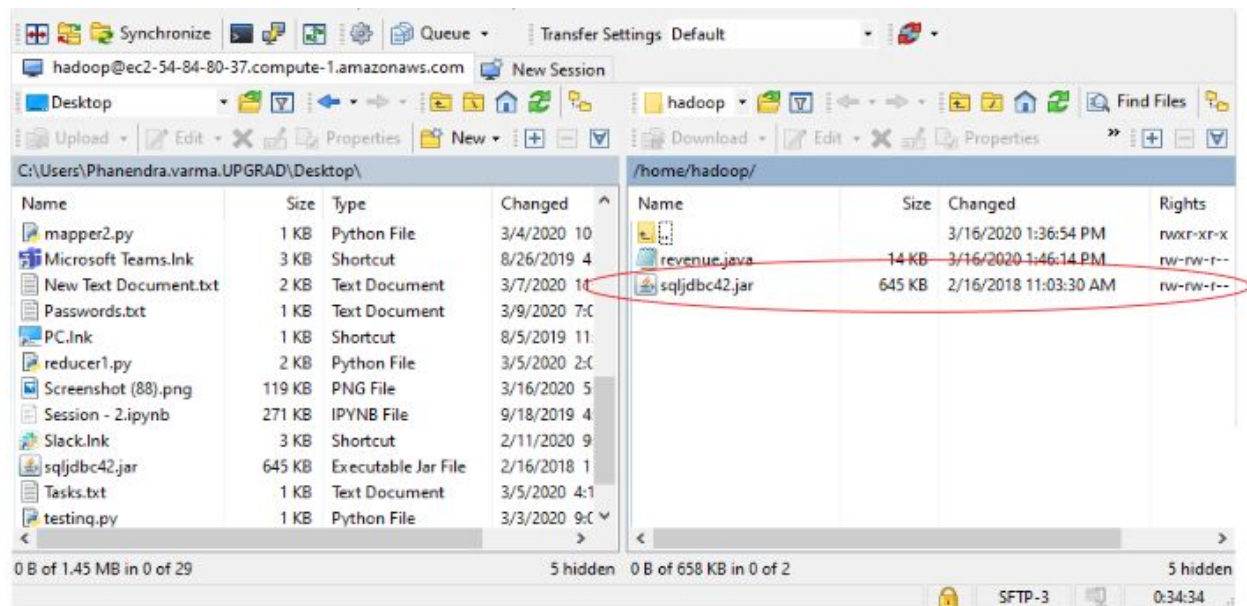
Things to ensure before running the sqoop job

1. Make sure the EMR version you are running is 5.24.0. The latest version of EMR isn't stable and creates problems while creating a workflow.
2. An RDS instance with necessary data loaded on the tables
3. Sample sqoop job running fine through CLI
4. Download the jdbc driver using the wget command and the link given below

OR

Download it locally and copy it to your EMR cluster using the WINSFTP tool -

<https://www.microsoft.com/en-us/download/details.aspx?id=54671>



5. Sample sqoop job running fine on your cluster

```
sqoop import --connect
jdbc:mysql://sample-database.cqsesz6h9yjjg.us-east-1.rds.amazonaws.com:3306/
telco --table revenue --target-dir /user/hadoop/telco/revenue/ --username
admin -P -m 1
```

sample-database.cqsesz6h9yjjg.us-east-1.rds.amazonaws.com - Replace this with the appropriate RDS instance address

telco - Name of the database

revenue - Name of the table

/user/hadoop/telco/revenue/ - Target directory

Now, steps to add a sqoop job in a oozie workflow

1. Launch the hue service by clicking on hue

Cluster: My cluster **Waiting** Cluster ready after last step completed.

Summary Application history Monitoring Hardware Configurations Events Steps Bootstrap actions

Connections: [Hue](#) [Resource Manager](#) ... (View All)

Master public DNS: ec2-54-84-80-37.compute-1.amazonaws.com [SSH](#)

Tags: -- [View All](#) / [Edit](#)

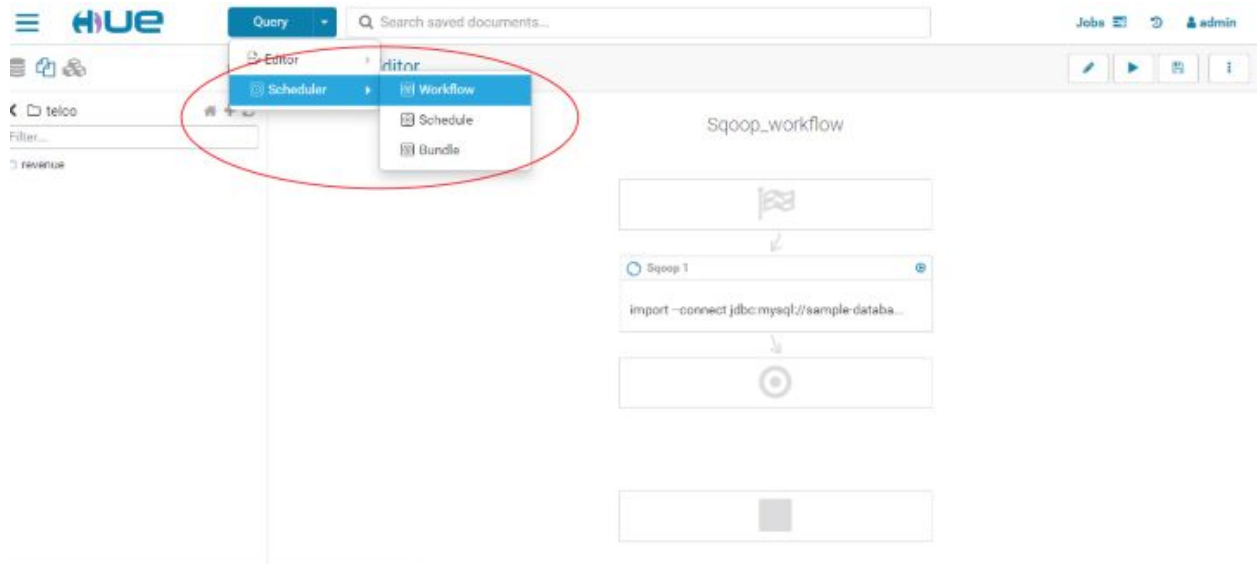
Summary	Configuration details
ID: j-10GMYH9WHQLFS	Release label: emr-5.24.0
Creation date: 2020-03-16 13:34 (UTC+5:30)	Hadoop distribution: Amazon 2.8.5
Elapsed time: 3 hours, 11 minutes	Applications: Hive 2.3.4, Pig 0.17.0, Hue 4.4.0, Sqoop 1.4.7, Oozie 5.1.0
After last step Cluster waits completes:	Log URI: s3://aws-logs-177300670946-us-east-1/elasticmapreduce/
Termination On Change protection:	EMRFS consistent view: Disabled
	Custom AMI ID: --

2. This will redirect to a new landing page where you are supposed to create a hue user.
3. On creating the hue user, the next step is to copy the driver to this user. This is to make sure that the jdbc driver is available with the hue user to run the sqoop job.

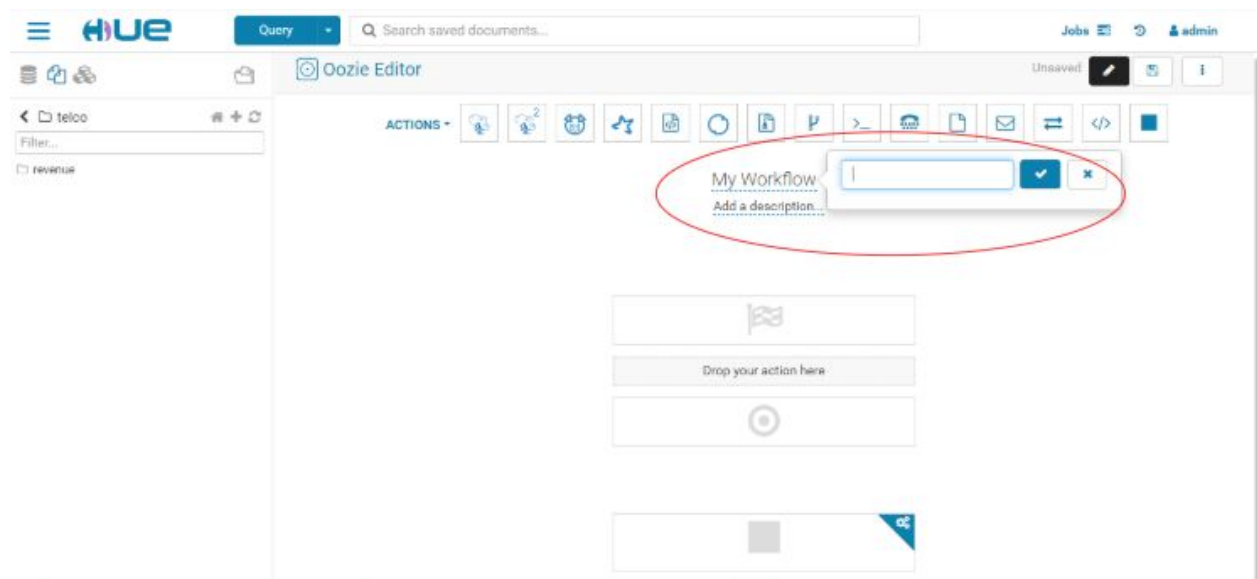
On the CLI run the following commands to copy the file to the hue user.

hadoop fs -put sqljdbc42.jar /user/admin/

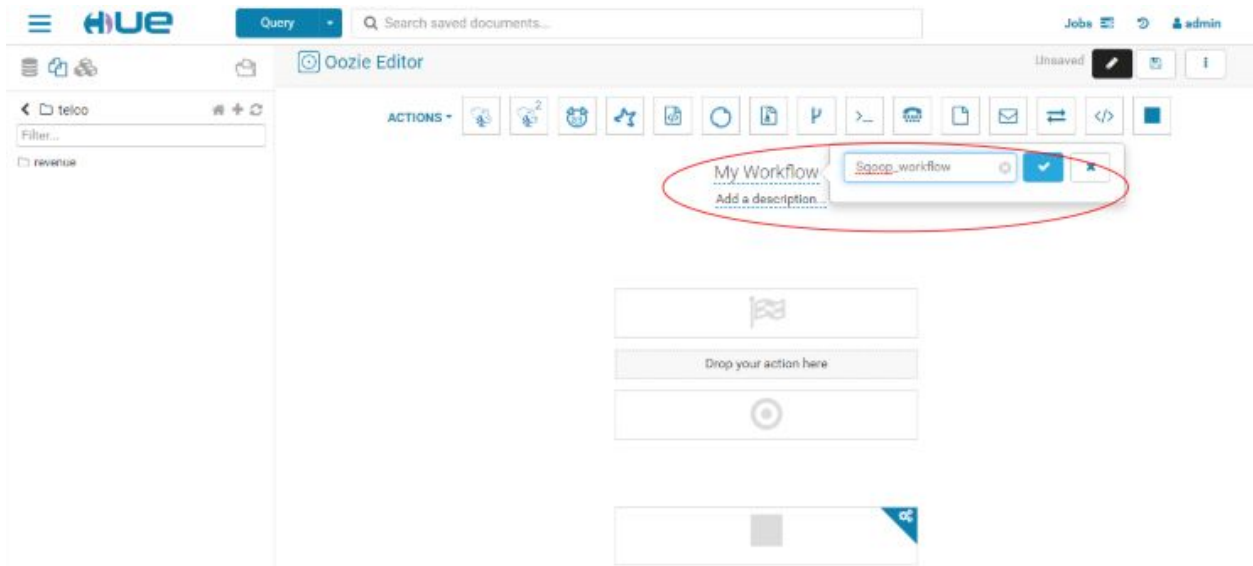
4. Now that you have the driver with the hue user the next step is to create a workflow. For this click on **schedule** and select **workflow**.



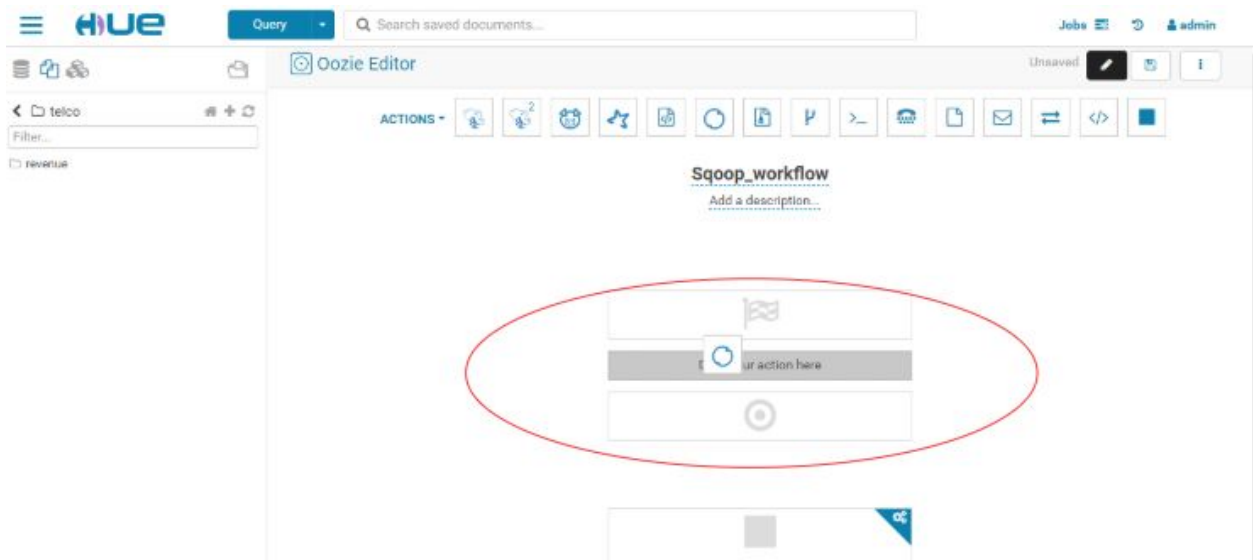
5. Give the name to the workflow as shown in the image below



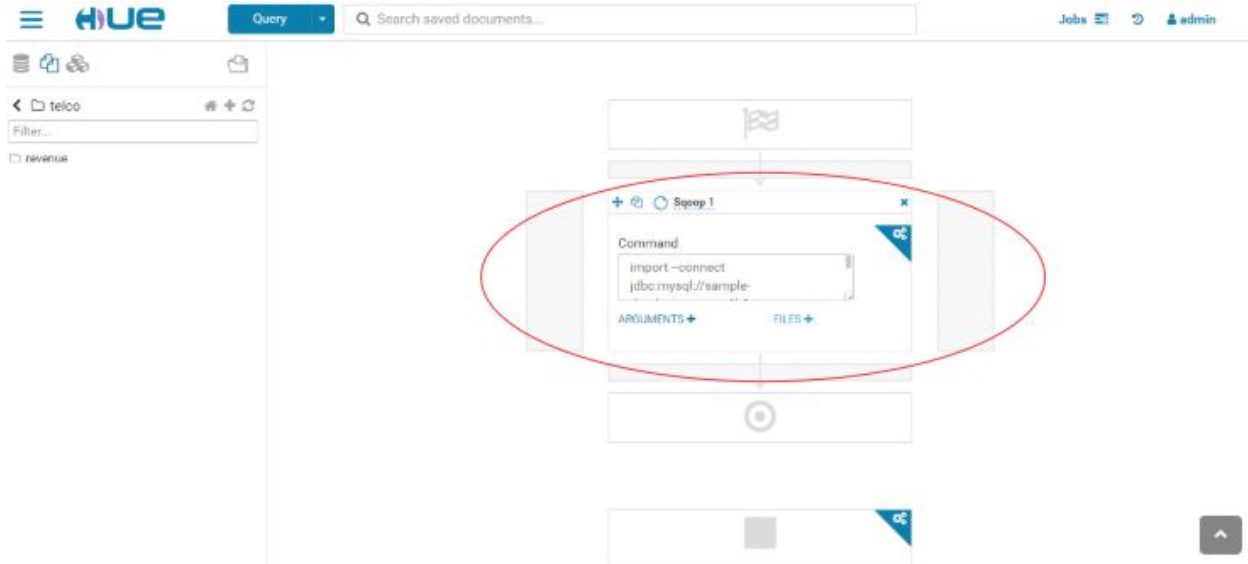
For example you can name it as **sqoop_workflow**



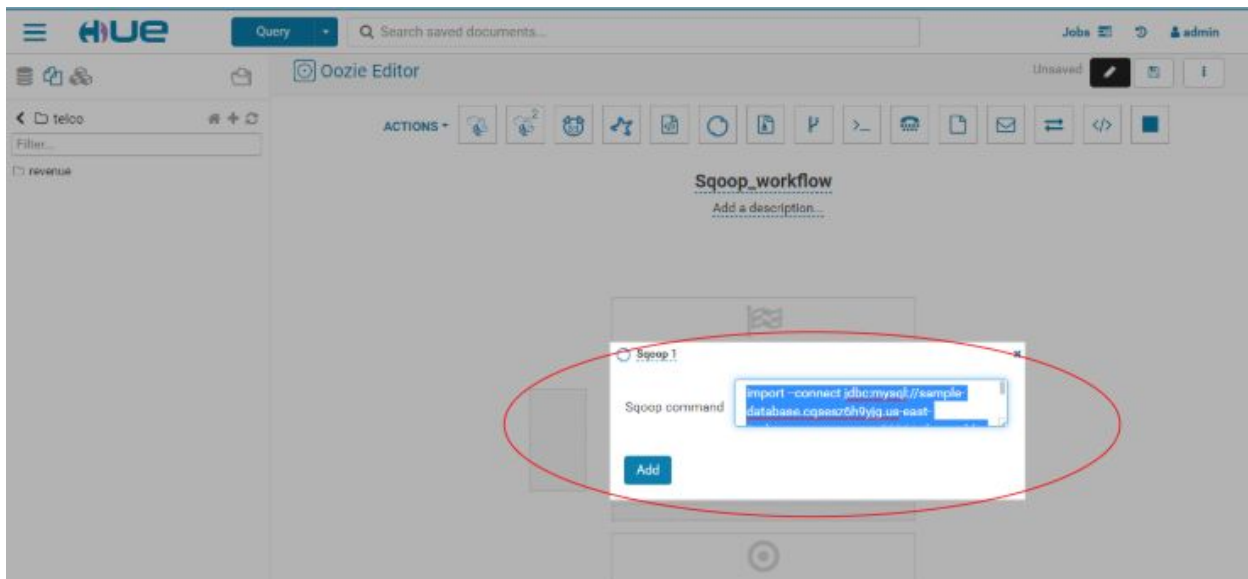
6. Next based on the desired action drag and drop the desired service. In this case drop the sqoop service as shown in the figure



Dropping the sqoop service enables you to enter the necessary command



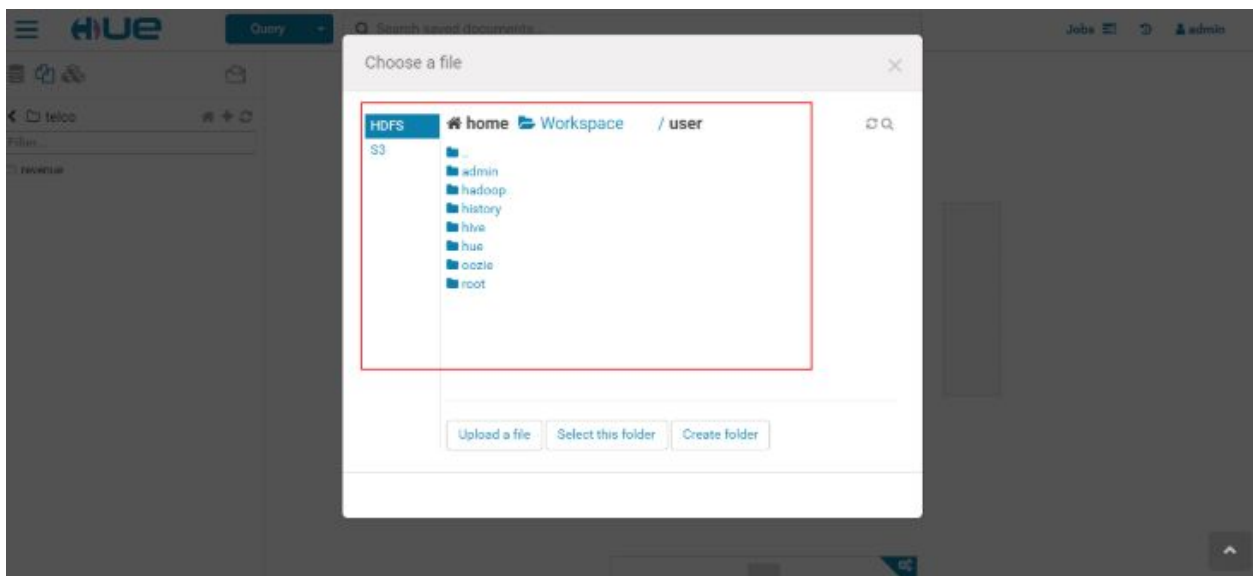
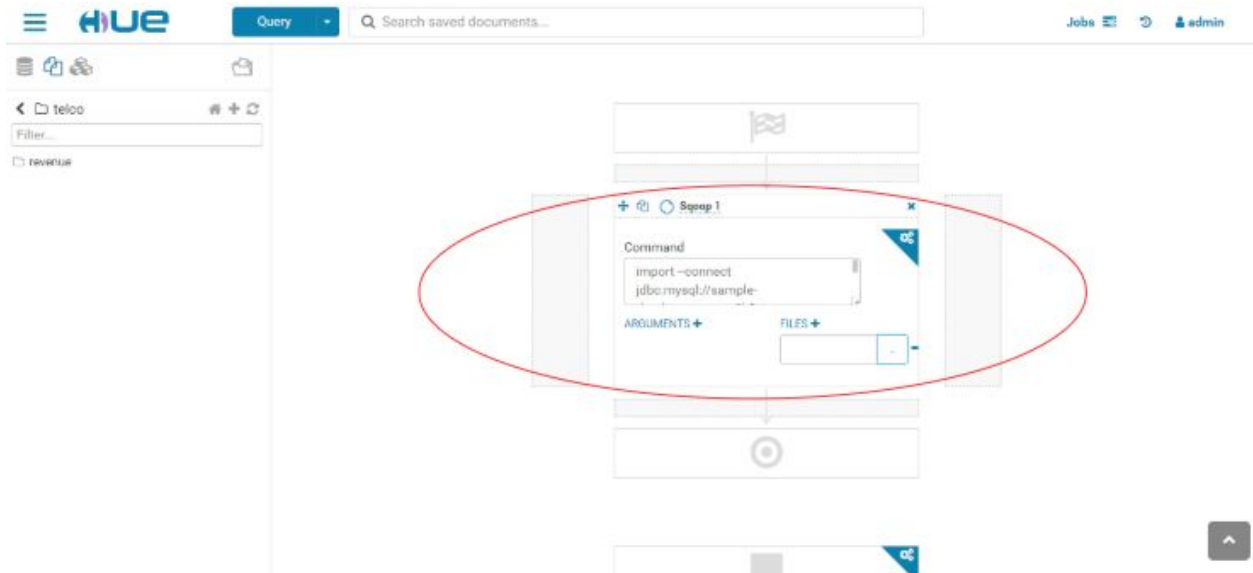
7. Under the command field give the following command

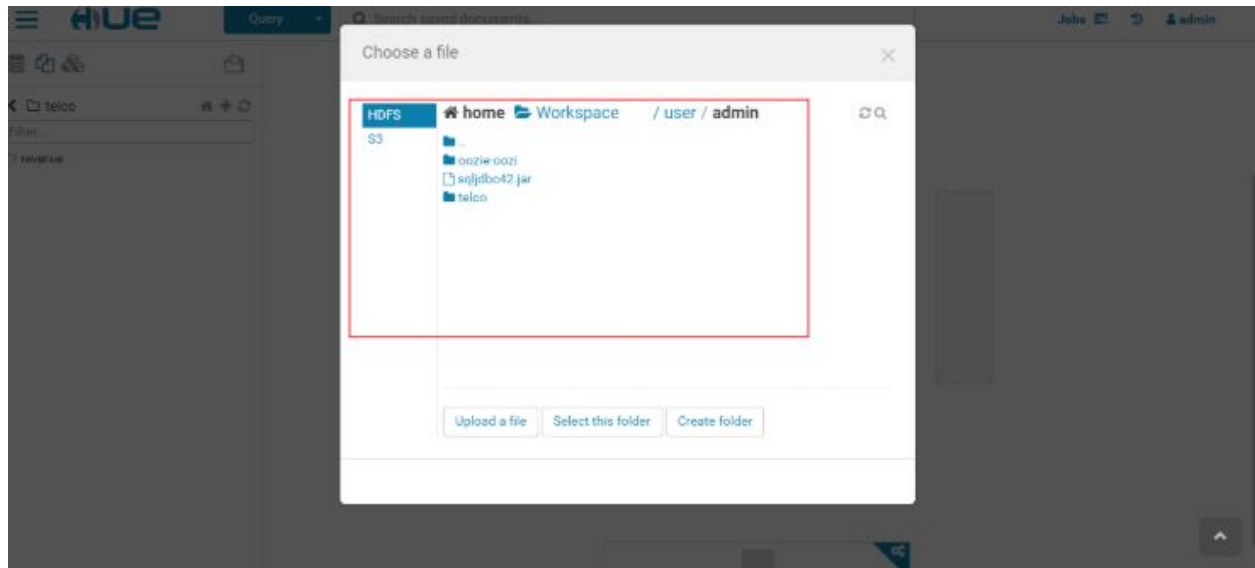


```
import --connect
jdbc:mysql://sample-database.cqsesz6h9yjq.us-east-1.rds.amazonaws.com:3306/telco
--table revenue --target-dir /user/admin/telco/revenue/ --username admin --password
admin123 --m 1 --driver com.microsoft.sqlserver.jdbc.SQLServerDriver
```

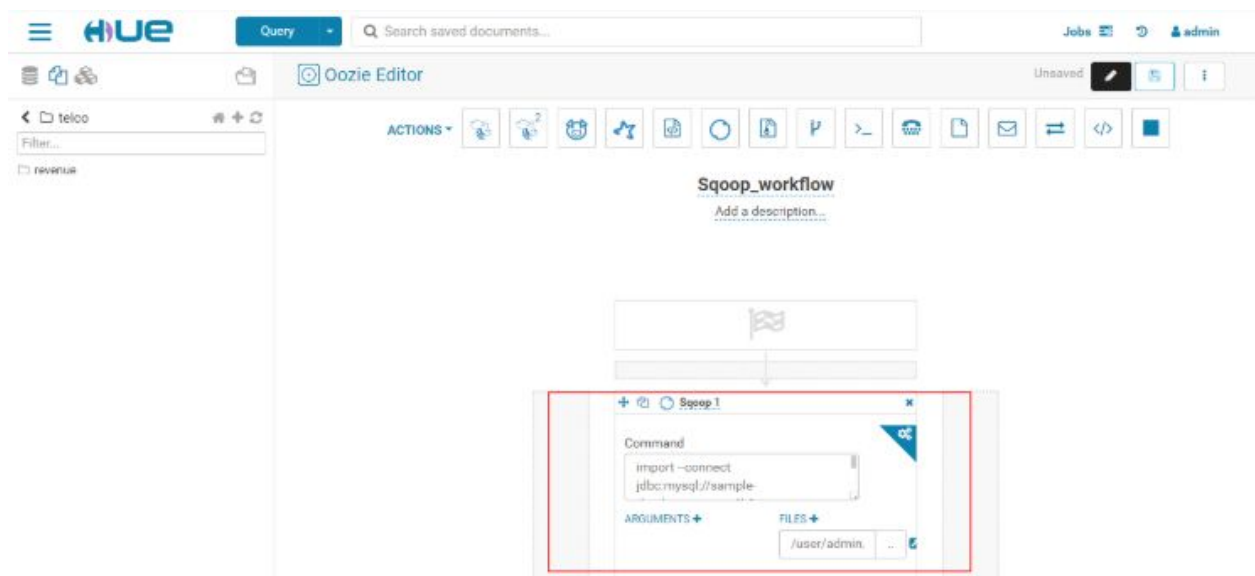
The username and password here are that of the RDS instance and not that of the hue service

8. In the command section you can find the option to choose the necessary files while running the sqoop job.





9. Finally after selecting the file save and run the job



10. Once you are job is submitted you can see the status under the **job browser section** and by clicking on the job you can see the details of the different tasks involved in the job.

The top screenshot shows the Hue Job Browser interface. The 'Jobs' tab is selected, and the 'Running' section is active. A table lists the running jobs, including their names, users, types, statuses, progress, groups, start times, durations, and IDs. The 'Completed' section also lists jobs that have finished, including their names, users, types, statuses, progress, groups, start times, durations, and IDs.

The bottom screenshot shows the detailed view of a specific job. The job name is 'oozie:action:T=sqoop:W=oozie_sqoop:A=sqoop-3f30:ID=0000004-200316081254187-oozie-oozi-W'. The job is in the 'SUCCEEDED' status, with a progress of 100%. The job type is 'MAPREDUCE'. The user is 'admin'. The job duration is 5s. The job is currently running on 100% of the available resources (1/1). The job is currently running on 0% of the available resources (0/0). The job is currently running on 5s of the available resources (5s).

The job details section shows the following information:

- ID: job_15843463899...
- NAME: oozie:action:T=sqoop...
- TYPE: MAPREDUCE
- STATUS: SUCCEEDED
- USER: admin
- PROGRESS: 100%
- MAP: 100% 1 / 1
- REDUCE: 0% 0 / 0
- DURATION: 5s

The job details section also shows the following information:

- Log Type: stdout
- Log Upload Time: Mon Mar 16 12:02:49 +0000 2020
- Log Length: 0

11. After successful execution of the sqoop job, you can see the files created at the destination specified

< admin

Filter...

oozie-oozi
sqljdbc42.jar
telco

< telco

Filter...

revenue

< revenue

Filter...

_SUCCESS
part-m-00000