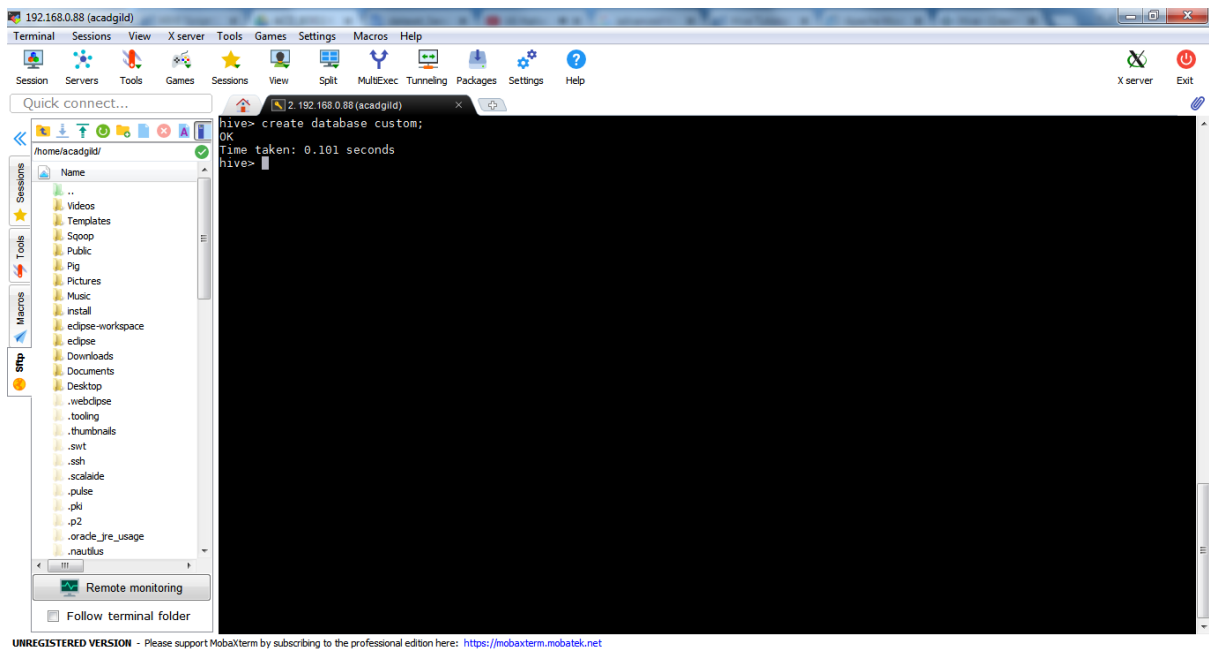# 5. Problem Statement

## Task 1

## Create a database named 'custom'.

Script-:

create database custom;

Create a table named temperature_data inside custom having below fields:

1. date (mm-dd-yyyy) format

2. zip code

3. temperature

The table will be loaded from comma-delimited file.

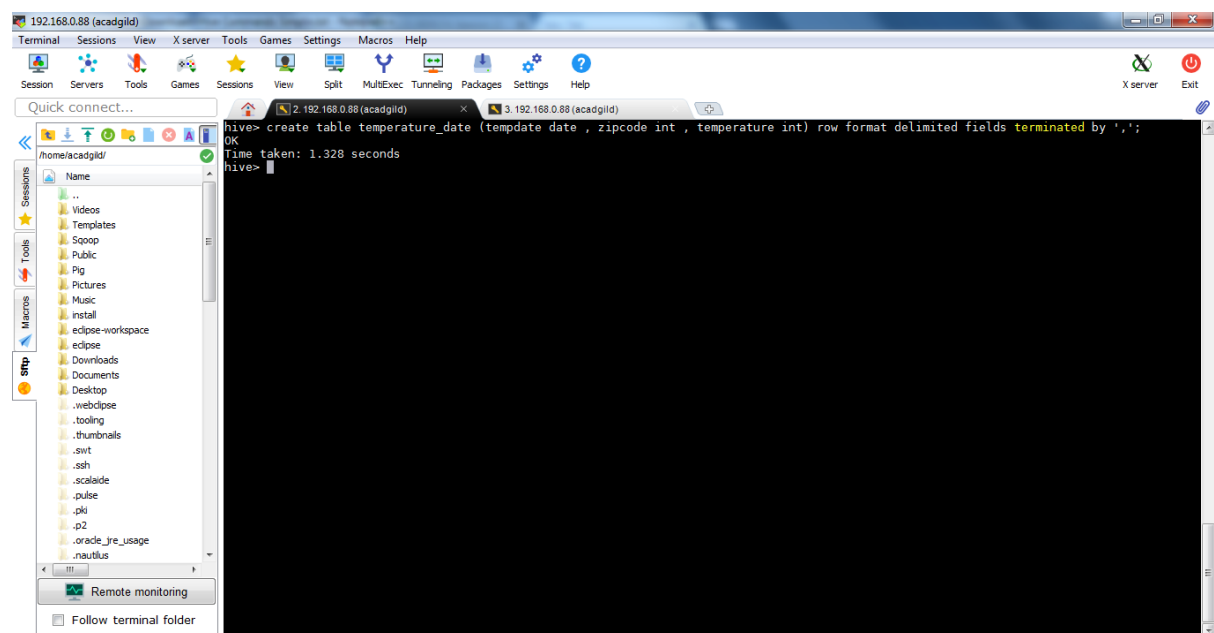Load the dataset.txt (which is ',' delimited) in the table.
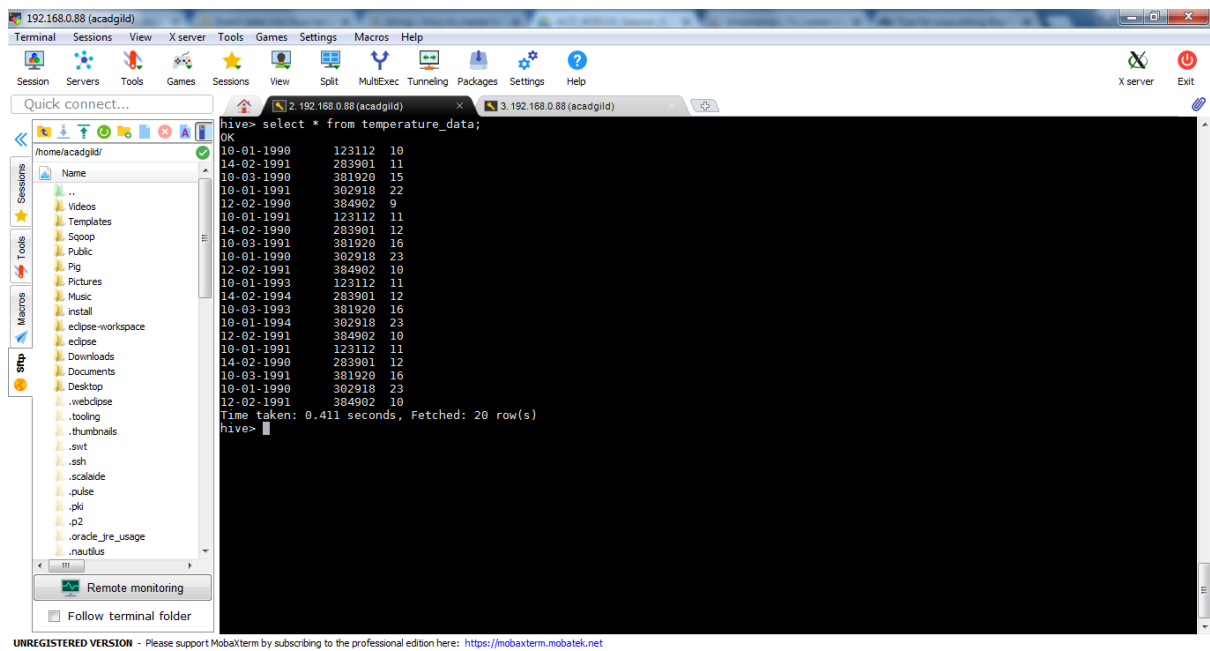
Script-:

Use custom;

create table temperature_data (tempdate date , zipcode int , temperature int) row format delimited fields terminated by ',';

LOAD DATA LOCAL INPATH '/home/acadgild/dataset_Session14.txt' into table temperature_data;
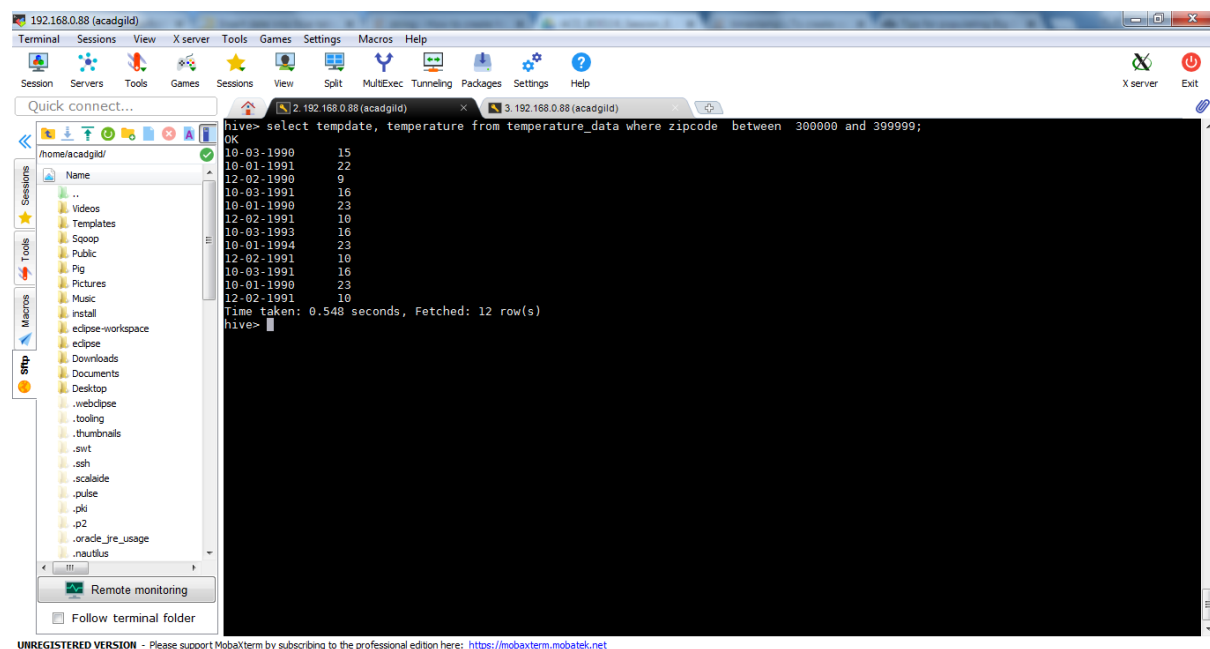
Output-:

## Task 2

● Fetch date and temperature from temperature_data where zip code is greater than

300000 and less than 399999.

Script-:

select tempdate, temperature from temperature_data where zipcode  between  300000 and 399999;

Output-:



● Calculate maximum temperature corresponding to every year from temperature_data table.

Script-:

select YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')),max(temperature) from temperature_data group by YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd'));

Output-:

● Calculate maximum temperature from temperature_data table corresponding to those years which have at least 2 entries in the table.

Script-:

select YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')),max(temperature) from temperature_data group by YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')) having count(YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')))>=2;

Output-:



● Create a view on the top of last query, name it temperature_data_vw.

Script-:

create view temperature_data_vw as select YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')),max(temperature) from temperature_data group by YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')) having count(YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')))>=2;

Output-:

● Export contents from temperature_data_vw to a file in local file system, such that each

file is '|' delimited.

Script-:

insert overwrite local directory '/home/acadgild/finaloutput.txt' row format delimited fields terminated by '|' select YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')),max(temperature) from temperature_data group by YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')) having count(YEAR(from_unixtime(unix_timestamp(tempdate ,'mm-dd-yyyy'), 'yyyy-MM-dd')))>=2;

Output-: