# Analyze a Dataset in Hive
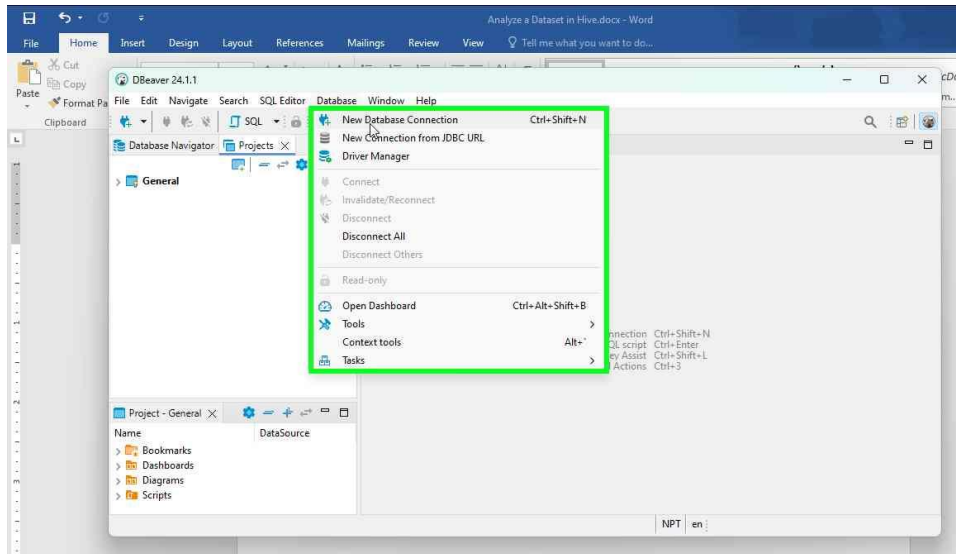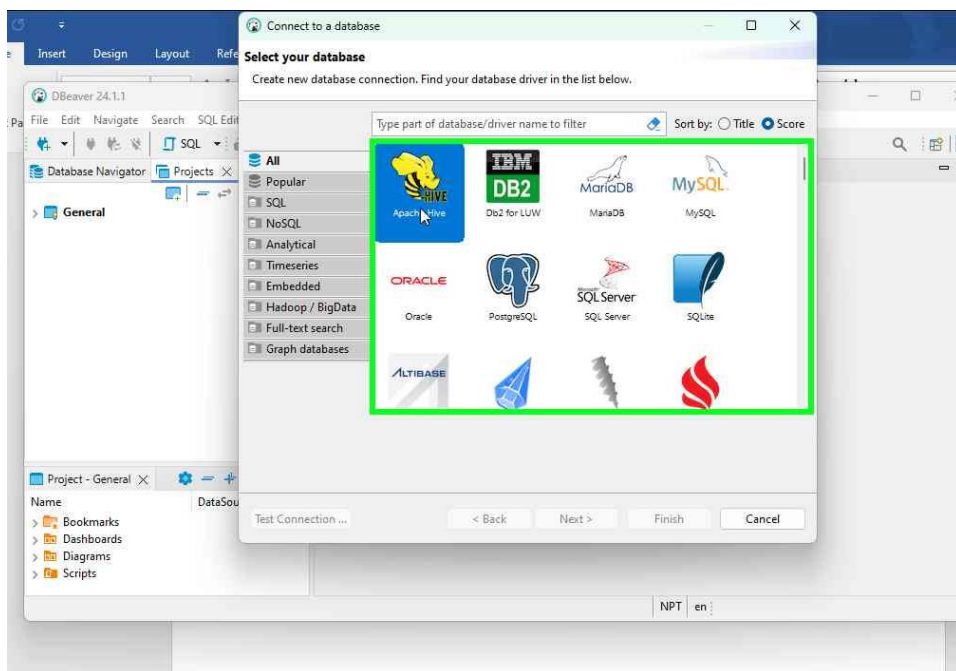
***

Note: This manual builds up from the previous manual of *Hive Installation*. Please refer to it before proceeding

▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪

**1. Download DBeaver from site: https://dbeaver.io/download/**

**2. Open the DBeaver and navigate to Database and click on "New Database Connection"**
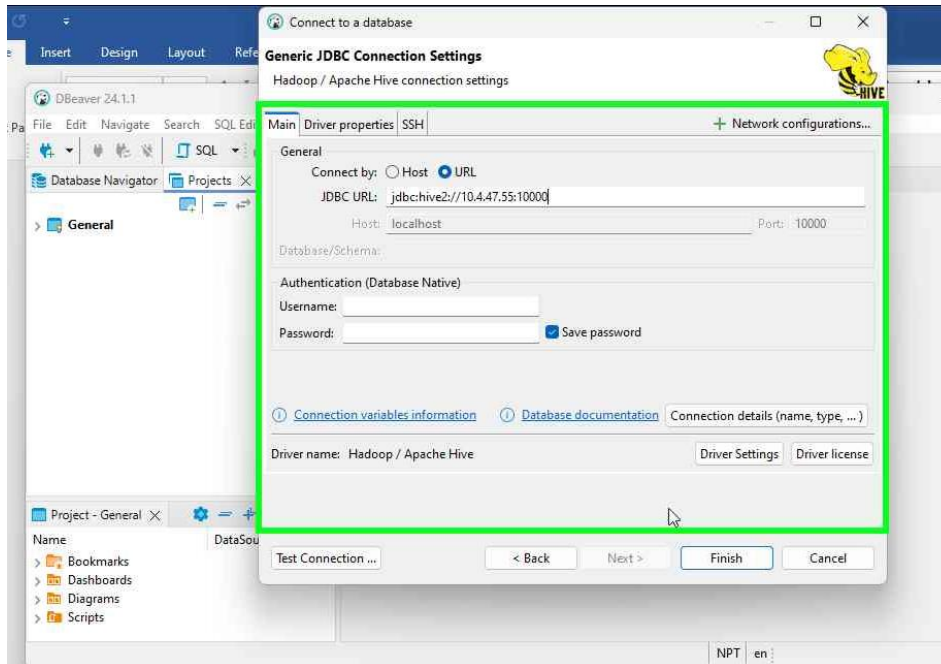


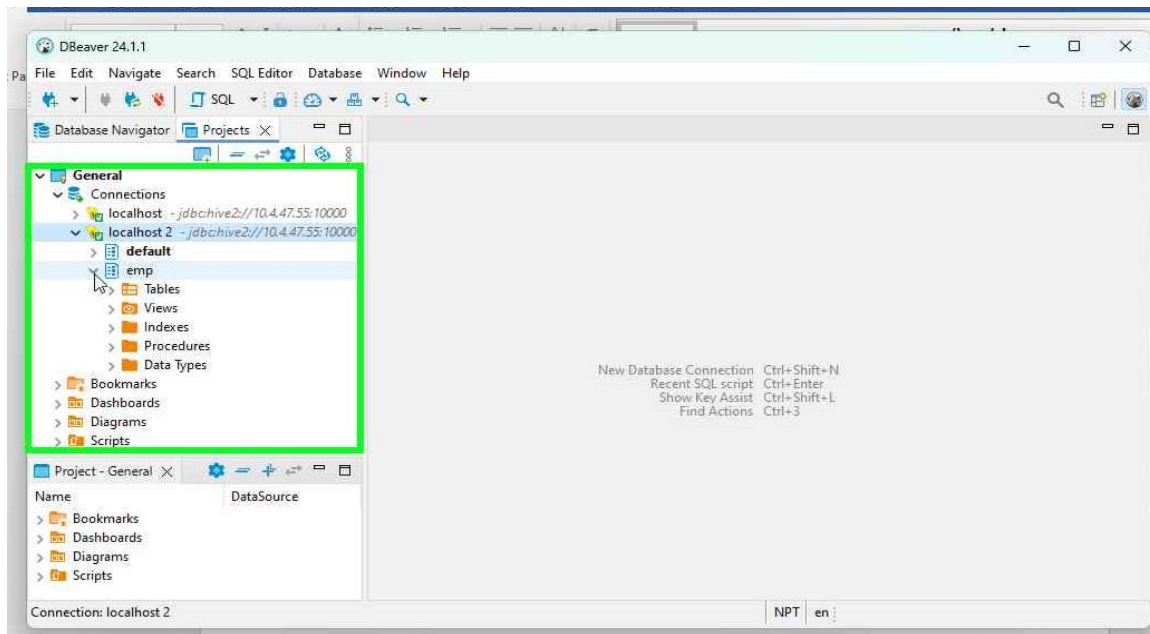**3. Select the Apache Hive as the Database type:**

**4. Select the "Connect by" option to: URL and enter the URL as shown in the image.**
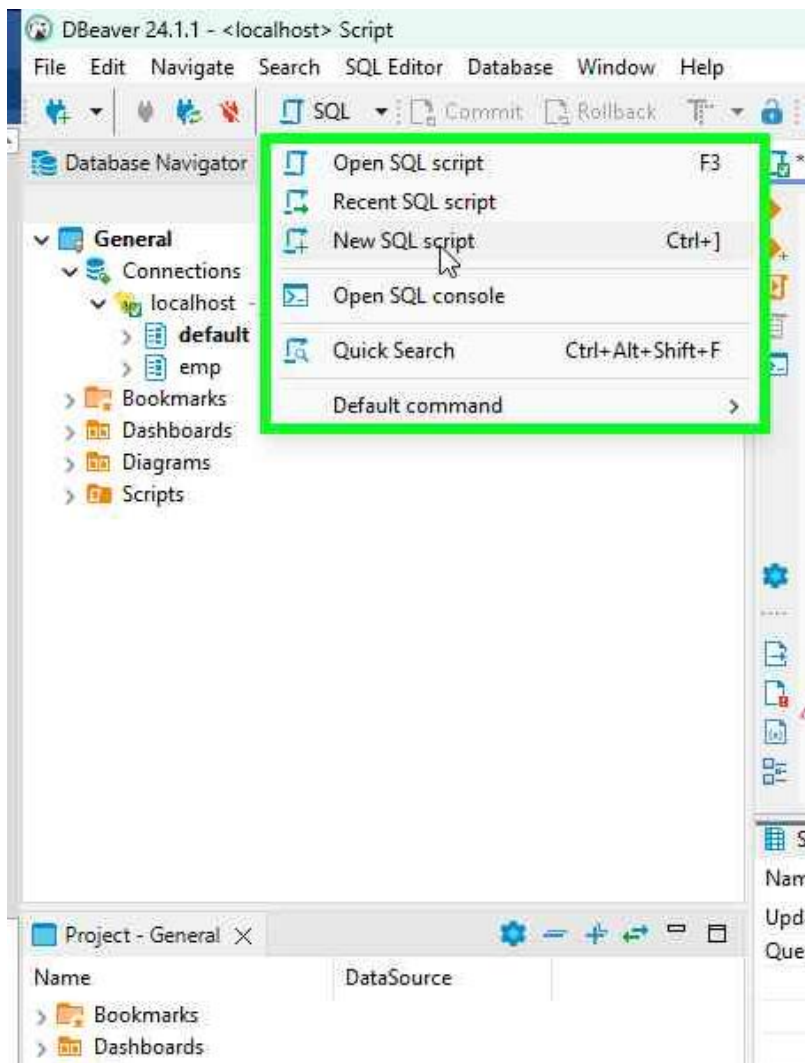
**Remember to put the IP of HIVE SERVER instead of the below shown IP 10.4.47.55. The port remains same.**



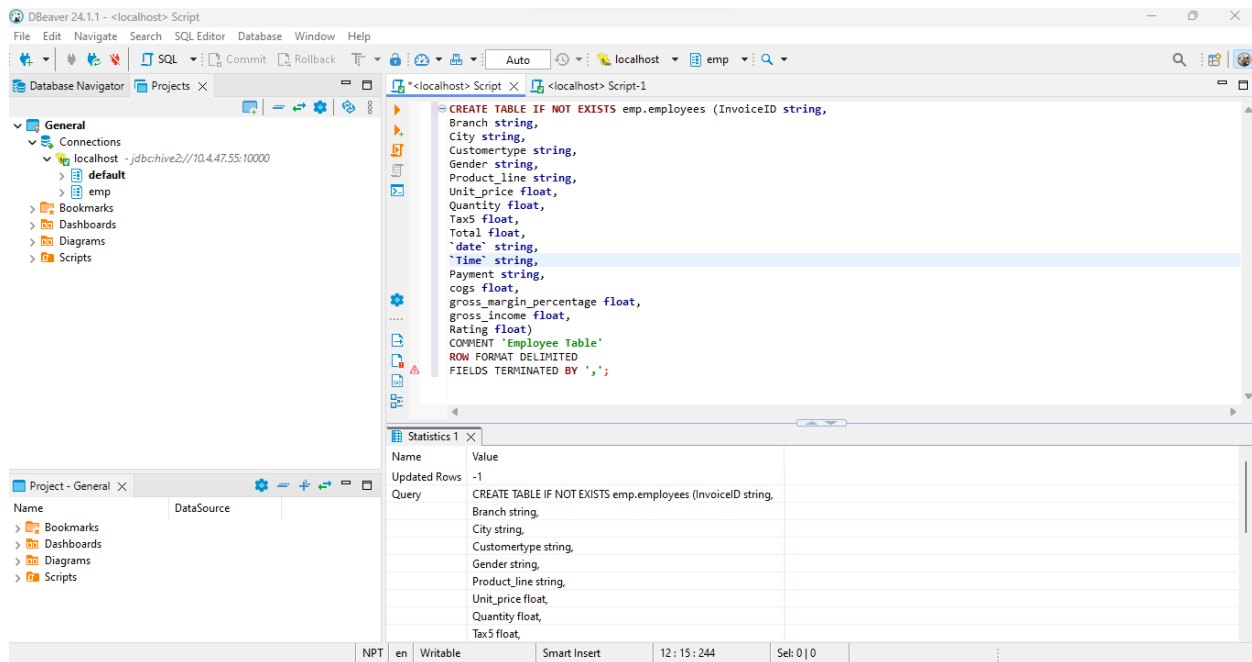**5. You should get the Hive database under your connections.**

**6. Click on new SQL script to open a SQL window.**



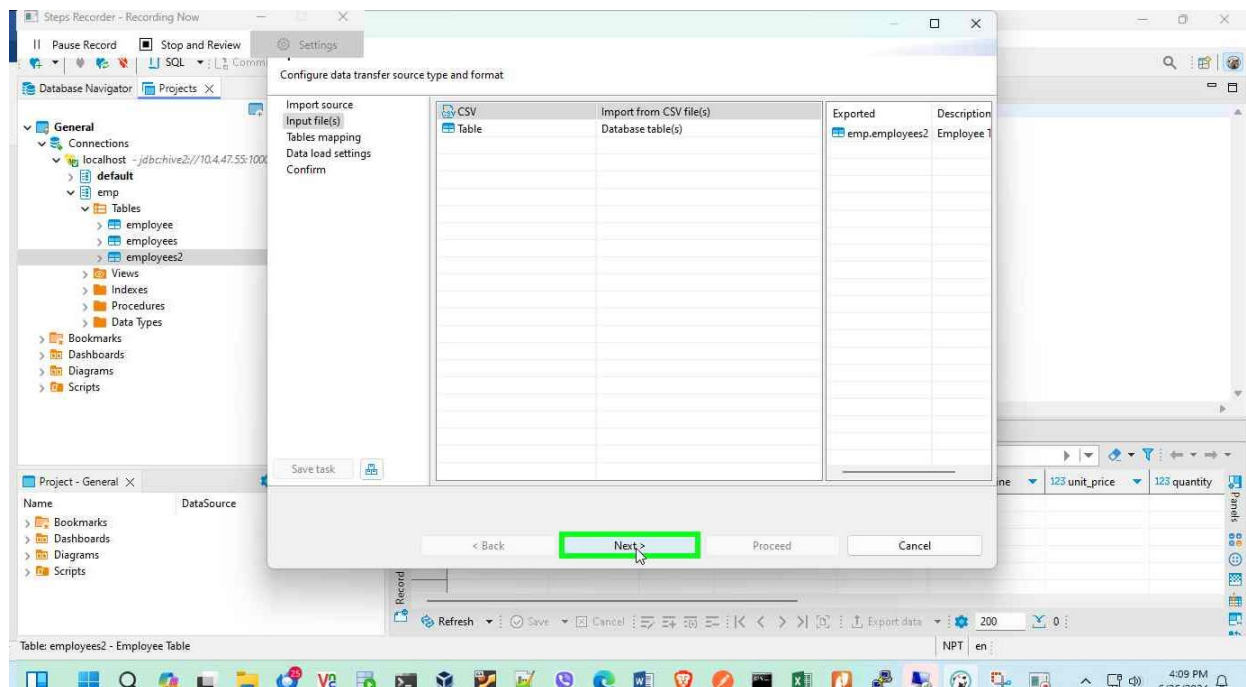**7. Create a table using following syntax. hit Ctrl+Enter to execute.**

```
CREATE TABLE IF NOT EXISTS emp.employees (InvoiceID string,
Branch string,
City string,
Customertype string,
Gender string,
Product_line string,
Unit_price float,
Quantity float,
Tax5 float,
Total float,
`date` string,
`Time` string,
Payment string,
cogs float,
gross_margin_percentage float,
gross_income float,
Rating float)
COMMENT 'Employee Table'
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ',';
```
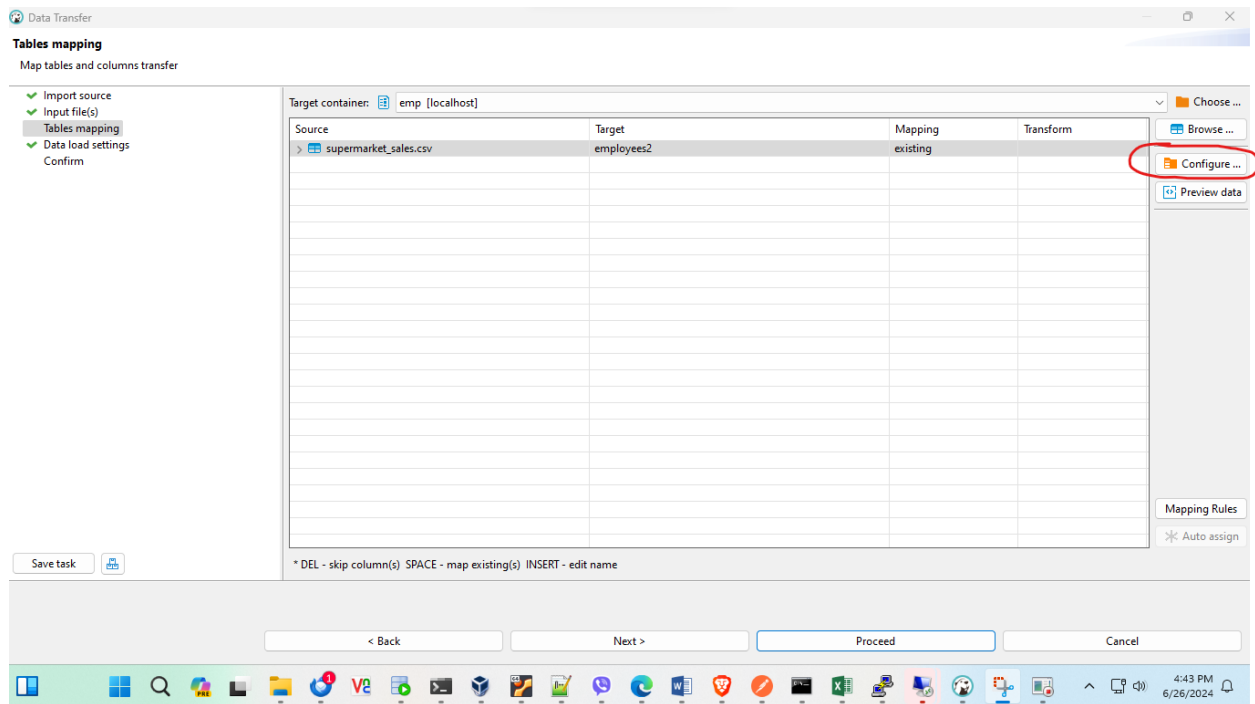
**8. Let's download a CSV file which will be inserted to the table.**
https://github.com/sersun/supermarket-sales-analysis/blob/main/supermarket_sales.csv
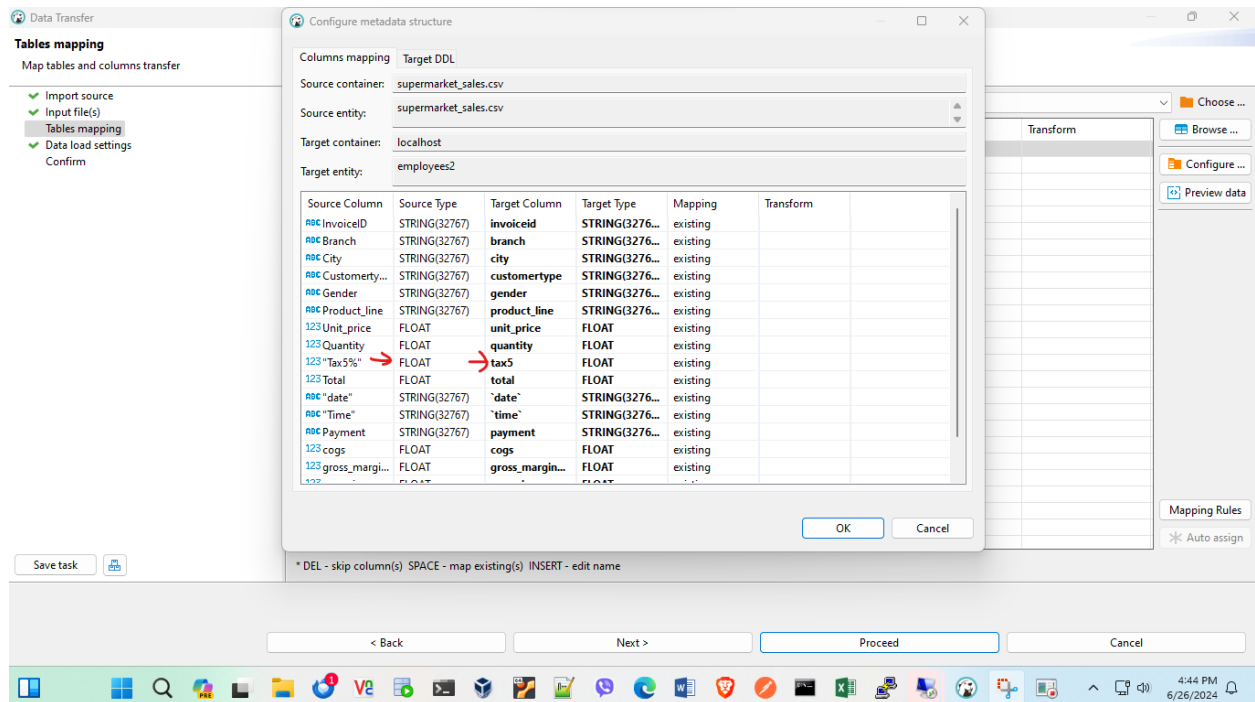
**9. Right click on the database table and select import data.**

## 10. Under the "Table Mapping" option, click on configure.



## 11. Make sure the special characters such as % are removed in the target Column.

**12. Click OK, Next, Next and finally Proceed. This will take some time to complete as the data is being transferred from host computer to Virtual Box VM to Hive DB.**

*13. To do this fast, we can download the CSV file in our LINUX machine first. And then load it to HIVE. Remember you should either import from DBeaver or do this step 13. Don't do both.*

**a. First open the hive shell in CLI.**

*cd $HIVE_HOME*

*bin/beeline -u jdbc:hive2://10.4.47.55:10000 hadoop*

**b. Load data into the HIVE.**

*LOAD DATA LOCAL INPATH '/home/hadoop/data.csv' INTO TABLE emp.employees;*

in above command put the appropriate location of your CSV file. in my case, the I am uploading the file data.csv which is in location /home/hadoop.

Now you can continue over to DBeaver.

**14. Examples:**

**a. Sort by unit_price**

```
select * from emp.employees order by unit_price desc;
```

**b. Selecting only `invoiceid`, `city`, `product_line`, `tax5`, `payment` columns**

```sql
select invoiceid, city, product_line, tax5,payment  from emp.employees;
```



**c. Count the total entries**

```sql
select count(*) from emp.employees;
```

**d. Group by `gender, product_line, quantity, total`**

```
select   gender, product_line, quantity, total
from emp.employees where quantity > 5
group by gender,product_line, quantity, total;
```



**e. Select using multiple criteria's**

```
select   gender, product_line, quantity, total
from emp.employees where quantity > 5
and total > 200 and product_line ='Electronic accessories';
```



**\*\*\* END OF MANUAL\*\*\***