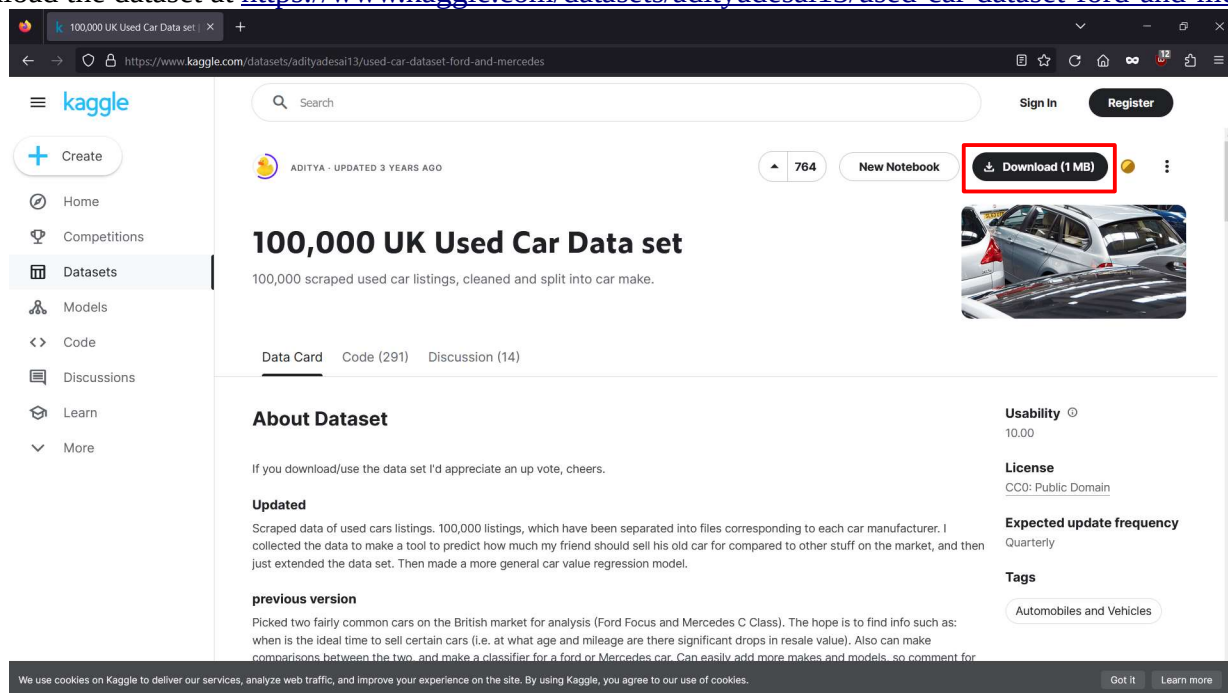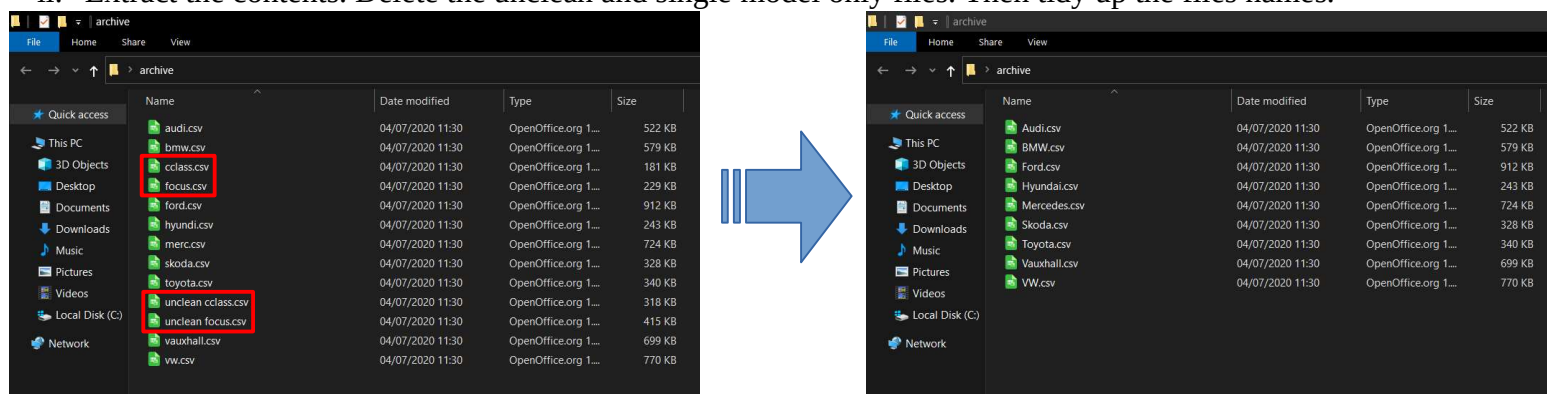**ETL to AWS RDS Project by Daniel Lee**

Objective:
1. Create a python script to merge csv datasets that are pre-cleansed & have a common schema into a single file inclusive of some light automated transformation.
2. Host the dataset in AWS RDS MySQL database.
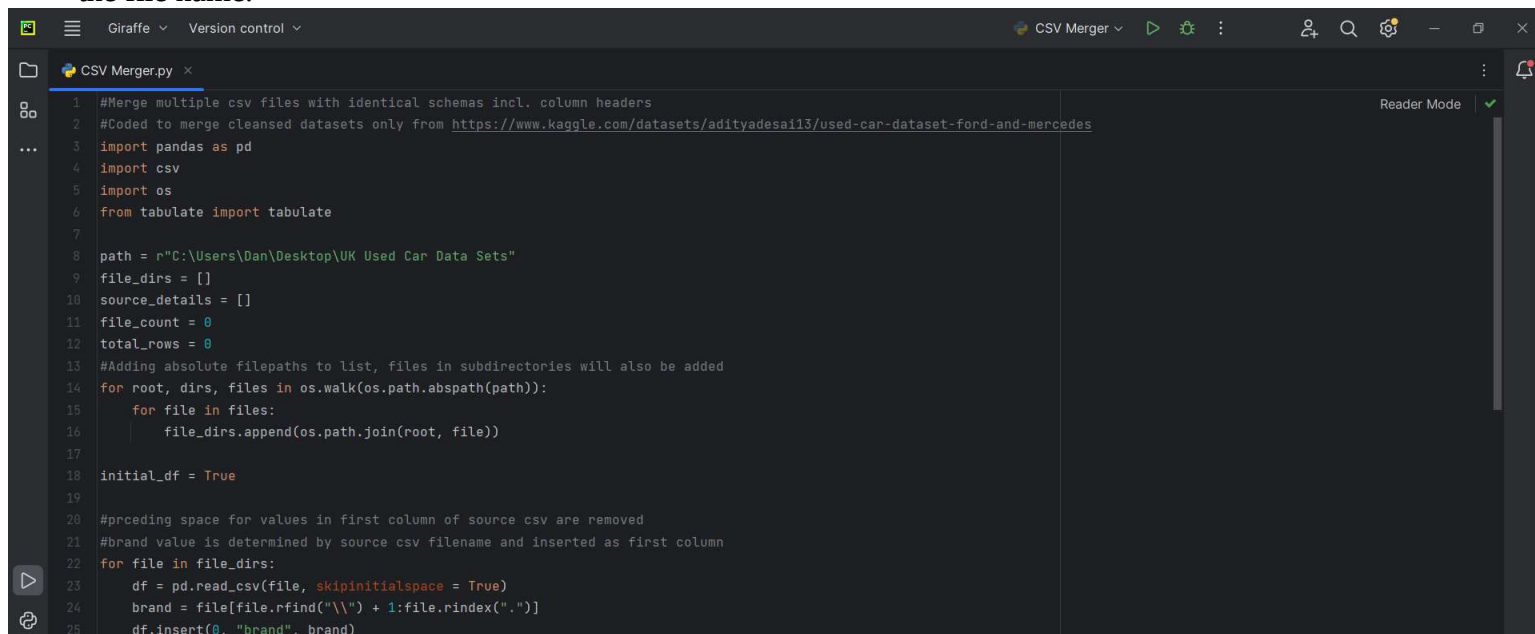3. Using MySQL Workbench to load the data to the RDS.

Process:
i. Download the dataset at https://www.kaggle.com/datasets/adityadesai13/used-car-dataset-ford-and-mercedes



ii. Extract the contents. Delete the unclean and single model only files. Then tidy up the files names.



iii. Create & run the python script to merge the files. This will also insert a new column for the brand which references the file name.

```python
    if initial_df == True:
        df.to_csv(r"C:\Users\Dan\Desktop\All Cars.csv", header=True, index=False, encoding="utf-8")
        initial_df = False
    else:
        df.to_csv(r"C:\Users\Dan\Desktop\All Cars.csv", mode="a", header=False, index=False, encoding="utf-8")

    filename = file[file.rfind("\\") + 1:]
    rows = len(df)
    source_details.append((filename, rows))

    file_count += 1
    total_rows += rows

source_details_col_names = ["Filename", "Rows"]
print(tabulate(source_details, headers=source_details_col_names))

print("")

summary = []
summary.append((file_count, total_rows))
summary_col_names = ["Total Files Read", "Total Rows Read"]
print(tabulate(summary, headers=summary_col_names))

print("")

output_rows = pd.read_csv(r"C:\Users\Dan\Desktop\All Cars.csv")
print("Rows in merged file: ", len(output_rows))
```

Console output:



```
C:\Users\Dan\PycharmProjects\Giraffe\venv\Scripts\python.exe "C:\Users\Dan\PycharmProjects\Giraffe\venv\CSV Merger.py"
Filename        Rows
------------    ------
Audi.csv        10668
BMW.csv         10781
Ford.csv        17965
Hyundai.csv      4860
Mercedes.csv    13119
Skoda.csv        6267
Toyota.csv       6738
Vauxhall.csv    13632
VW.csv          15157

  Total Files Read    Total Rows Read
------------------  -----------------
                 9              99187

Rows in merged file:  99187

Process finished with exit code 0
```

Schema transformation:



iv. Go to the AWS RDS management console and start creating the database



v. In the configuration, select **Easy Create**, **MySQL** & **Free Tier**

vi. Set a database name, then the master username and password to *admin* & *password* respectively. Select the **Create Database** button at the bottom.



vii. Ignore the error message at the top in red. Wait for the status to change to available then click on the database name



viii.    You should now see an Endpoint. Copy and paste into Notepad for future reference. Click on Modify

ix. Scroll down to **Connectivity**, and expand **Additional Configuration**. Select **Publicly accessible**

x.  On the next screen select Apply Immediately and Modify DB Instance



xi.  Now click on the VPC security group
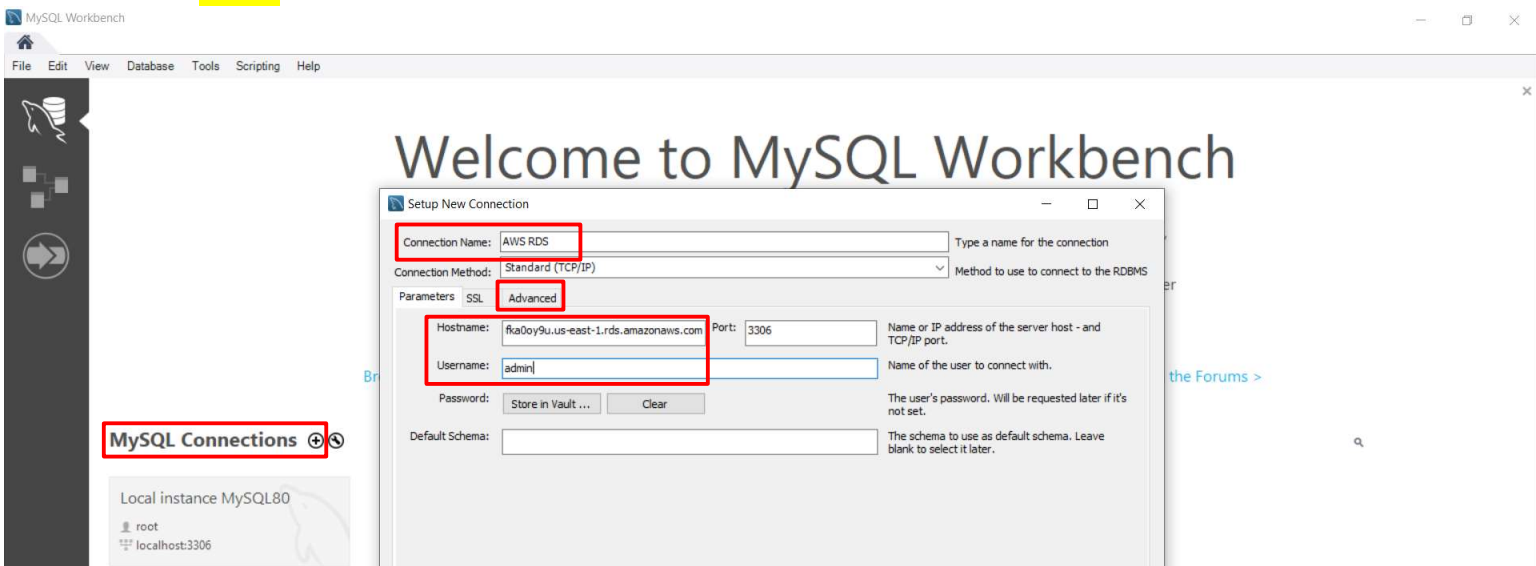


xii. Select the inbound rules tab then Edit Inbound Rules

xiii.    Click **Add Rule**, select **MYSQL/Aurora** & **My IP**, then **Save rules**
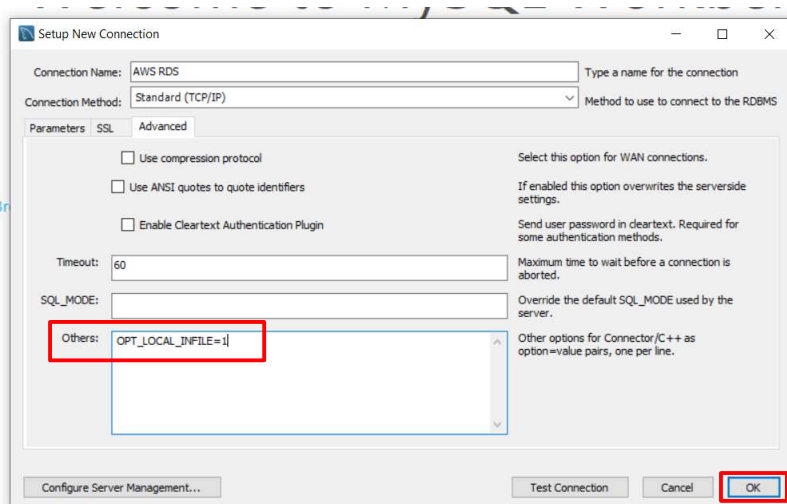


xiv. Open MySQL Workbench
   i.    Add a new MySQL connection
   ii.   Enter a Connection Name of your choice
   iii.  Retrieve the Endpoint from Notepad and paste as the Hostname
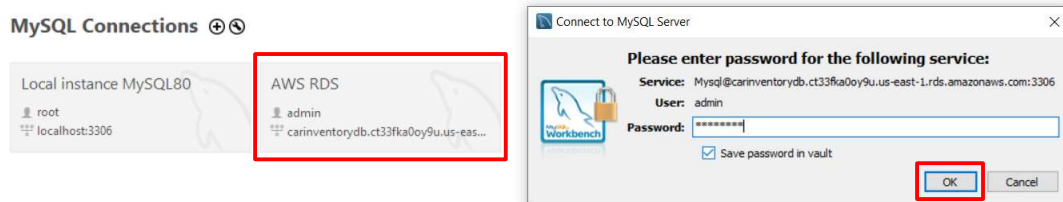   iv.   Enter admin for Username



   v.    Select the Advanced tab
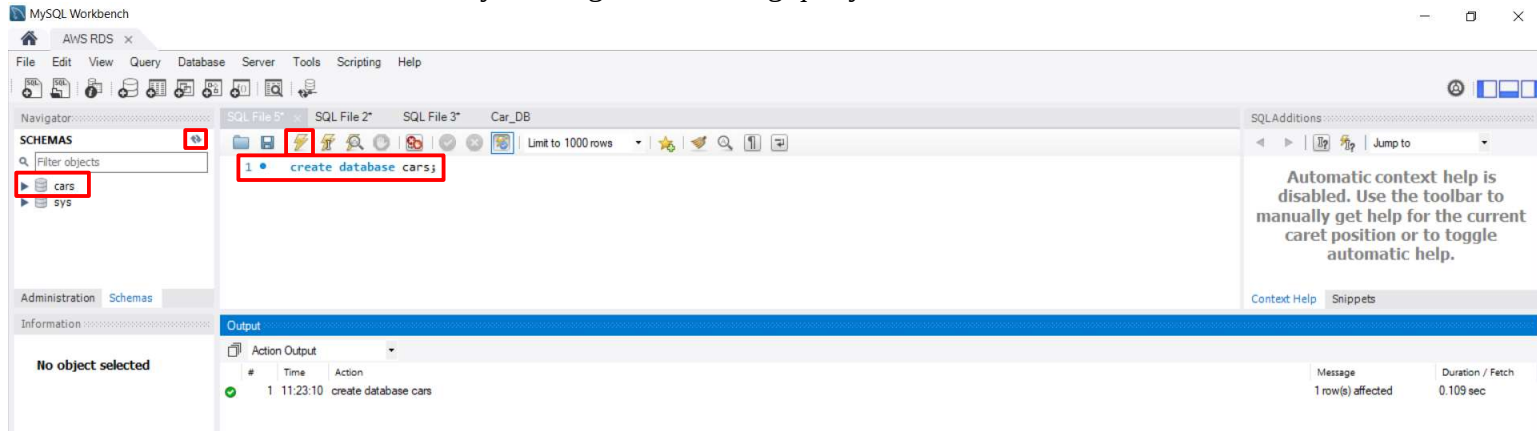   vi.   enter OPT_LOCAL_INFILE=1 into the Others box
   vii.  Click OK
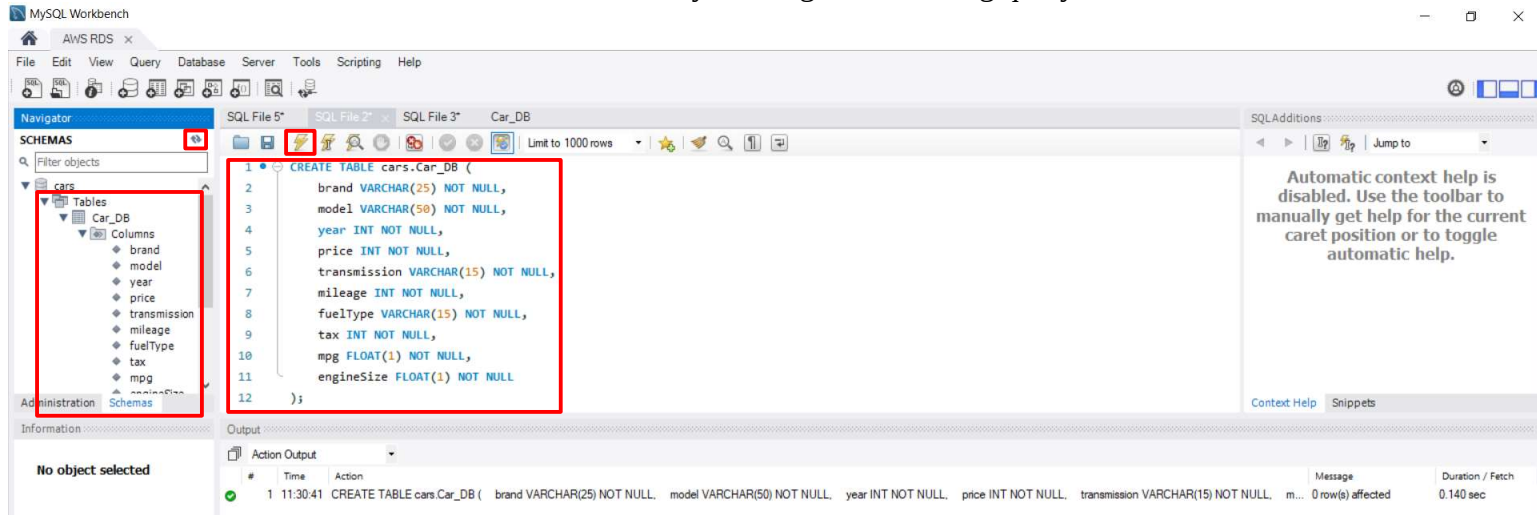
xv. Double-click on the new connection then enter password as the Password and click OK
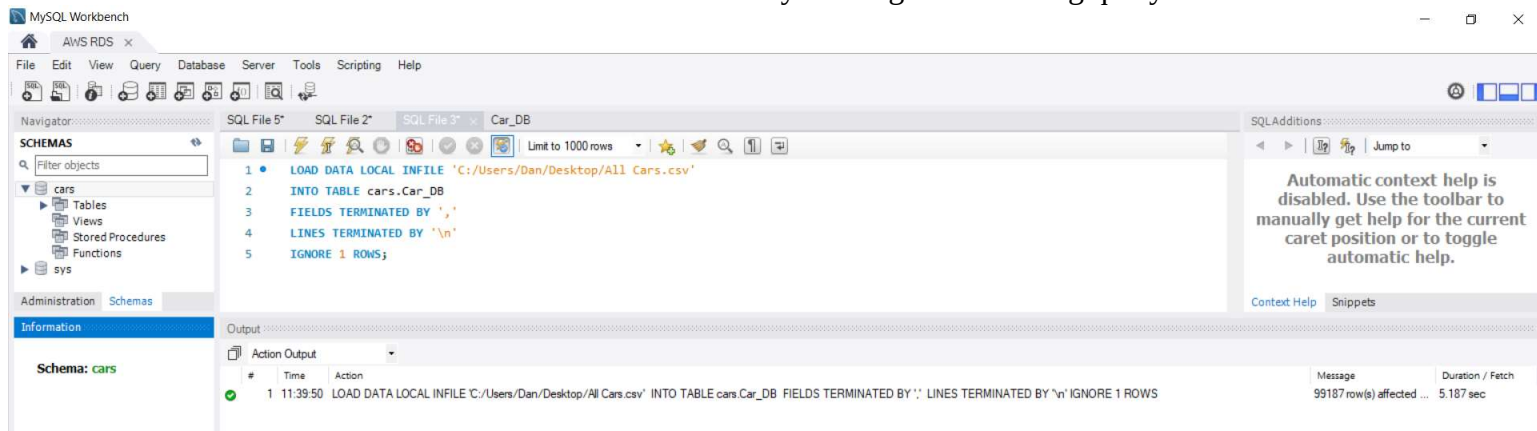


xvi.    Create a new database by running the following query



xvii.    Create a new table in the new database by running the following query



xviii.    Load the new table with the data from the csv by running the following query

xix.      Run a couple of checks:

  i.  SELECT COUNT(*) FROM cars.Car_DB;



  ii.  SELECT * FROM cars.Car_DB;