Course: IST 615 Cloud Management
Title: Used Car Retail Website

## Project Description:

Built a used car retail website using Microsoft Azure Cloud services where buyers can buy and enlist cars based on manufacturers, mileage, distance traveled, etc.

## Dataset Description:

The dataset was taken from Kaggle in the form of a csv file. It is originally from Craigslist. The dataset consisted of 426881 rows and 26 columns. Columns consisted of details of used car information like price, year, manufacturer, model, make, transmission, cylinders, mileage, distance travelled, condition, etc.
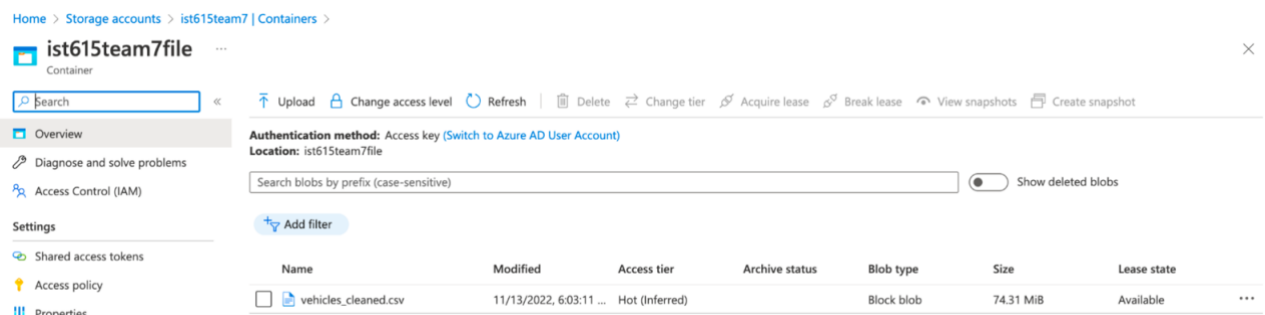
## Data Exploration and Cleansing:

Dataset was cleaned by removing rows with null values in columns like condition, manufacturer, model, fuel, odometer, and transmission. Column with the County name was also dropped as it had more than 50% null values. Similarly, the size and region URL were also removed due to a lot of missing values. The description column was also dropped as the column contained identical descriptions for more than 50% of the records.
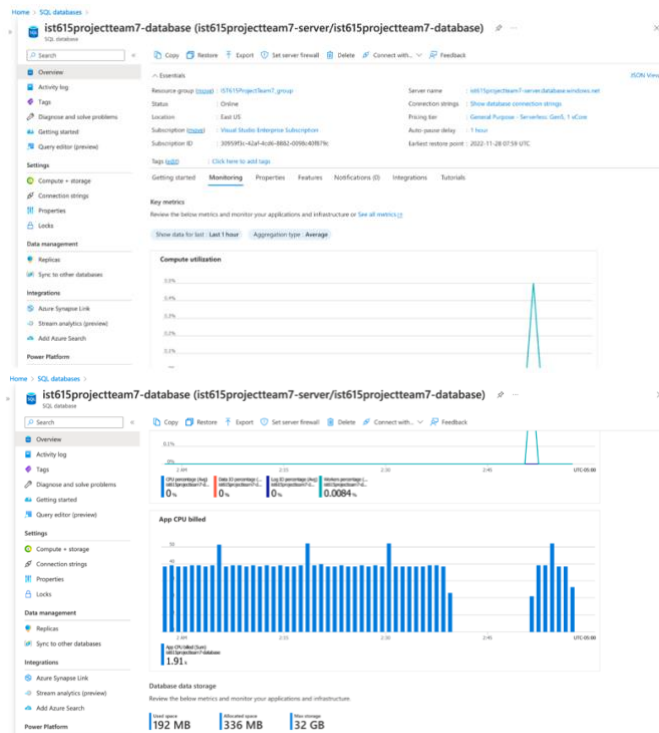
## Cloud Services Implemented:

Azure Blob Storage:

Our cleaned dataset file 'vehicles_cleaned.csv' was uploaded to Azure blob service and stored in a container.
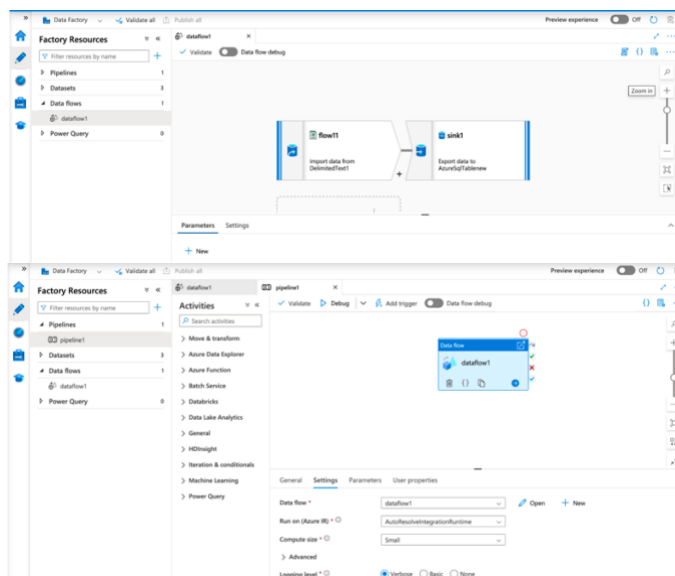


Azure SQL Database:

The SQL Database was set up so that the data from the blob storage can be integrated with the Database which was communicating with the website.
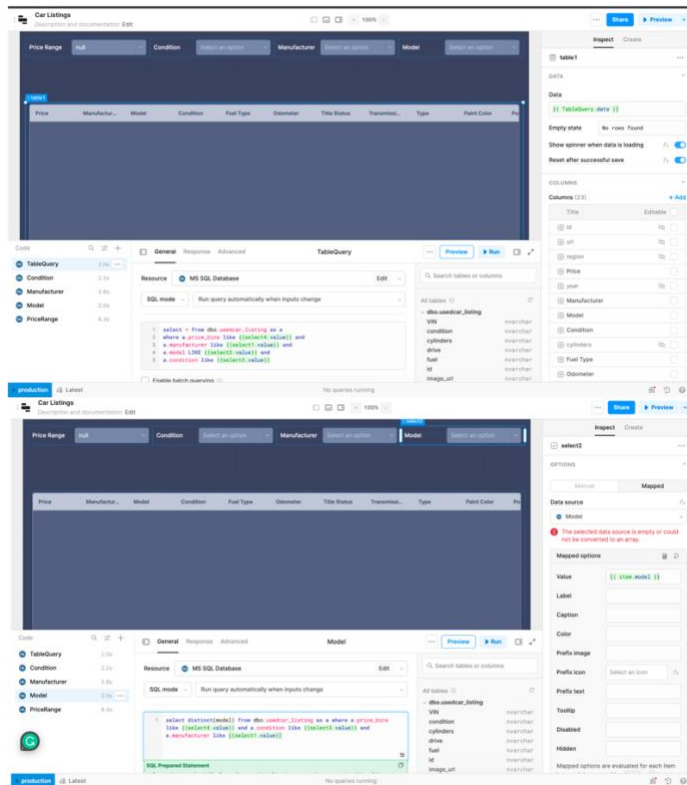
## Azure Data Factory:

Azure Data Factory was used to facilitate data engineering and data integration. A data flow was created to get data from the Azure Blob storage container to Azure SQL DB. This was executed using a pipeline for which we had to set up Azure Blob storage in a right way so that the column headers and the values are accurately analyzed. Data flow was used so that each field in the csv file from the Azure Blob storage was mapped to the column headers.
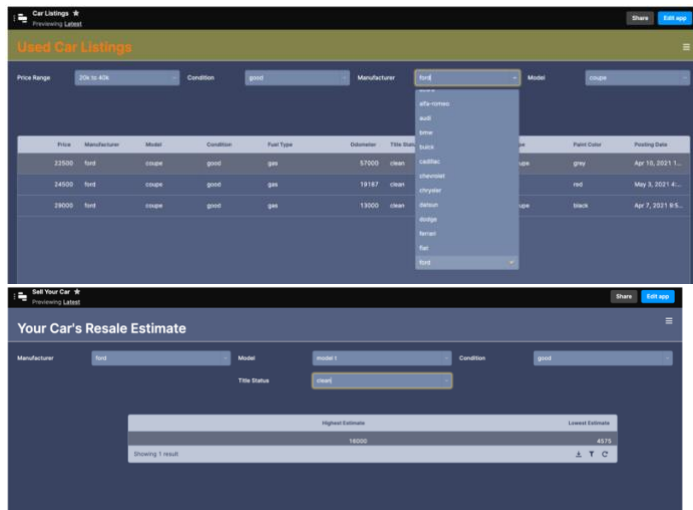




## Retool Service:

Retool is not a Microsoft Azure service but is a platform where users can connect to databases and create an application which acts like a CRUD application. To create our

website we created three different pages, one in which we can look for the cars listed on the website. To do that basic SQL + JavaScript was used. For Retool to communicate with the database some IPs were added to the firewall rules of the database. Each of the cells were attached to a query that got values from the database based on the selections made by the user.
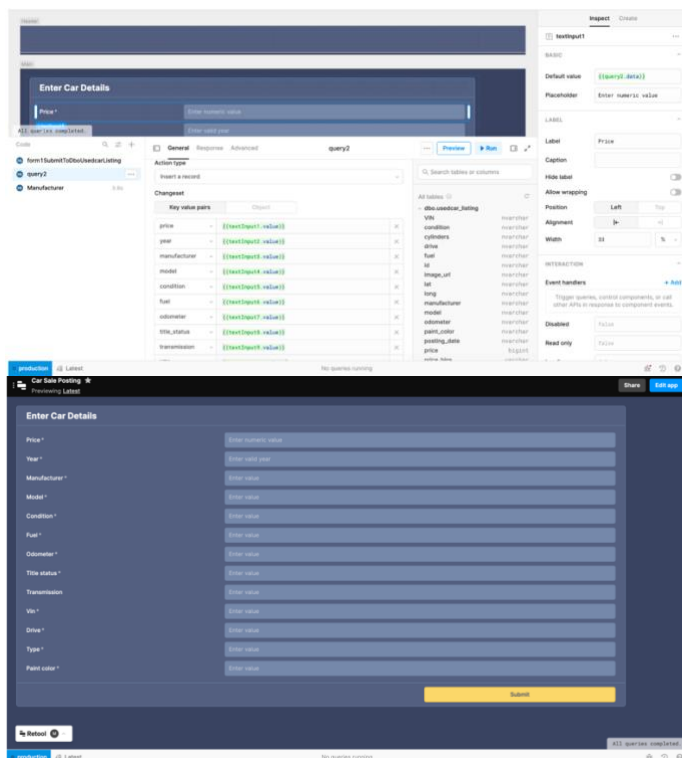


Similarly, another webpage was created for the user to check the value of their car - the highest and the lowest estimated resale value according to the data that we have in our database. The query to find the highest/lowest estimate is the highest/lowest value for that car that has been listed in our database i.e by the manufacturer, model, condition and the title status.

Once the user selects all the fields in a webpage, they will be able to the listings that are available in our database for that specific webpage. Different webpages could be accessed by clicking the navigation button on the top.
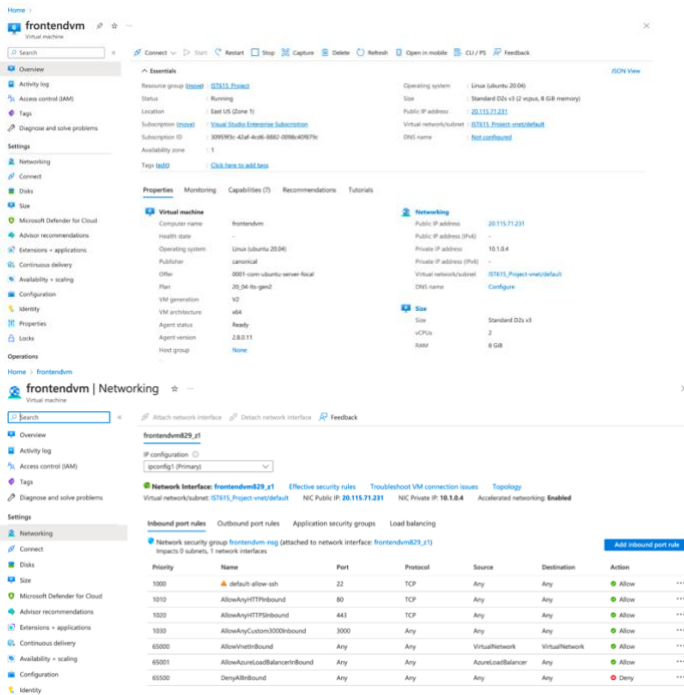
A third page was created for users to post their own car for sale. Each of the fields had to be filled to post their car for sale. To post the car, a query with Insert action was used which works exactly like the SQL insert into database command. Each field was given some constraints to enter the information.





Azure Virtual Machine:

Azure Virtual Machine was used to host out website in which we cloned a git docker environment which is responsible for hosting our website. This environment is created by Retool itself. Necessary security and network rules were added to it so that Retool can communicate with Azure. Retool service was hosted on port 3000. To access it we had to put in our public IP address {{pulic_ip}}/3000.

## Data Modelling:

Azure Databricks:

Azure Databricks was used for the purpose of python scripting, data analysis and model creation in our project. We wrote a script to integrate Azure Blob storage and Azure Databricks. Using this script, we can read data from our container in blob storage. Some additional data cleansing was done and new columns were added to add on to the features to build price estimation models for users' cars that they had to post online for sale. Three models - Linear Regression, Random Forest and XGBoost were built to predict estimated price. Our best prediction model was Random Forest and it was used to predict prices for these used cars.

## Future Scope:

- Include our model predicted prices in the price estimator webpage of our website.
- Establish connection to PowerBI and build dashboards to display our data analysis, predicted prices and our model performance over time.

## Specific contribution to the project:

- Assisted with data exploration and data cleansing.
- Assisted with Retool to build our website.
- Databricks analysis and multiple model building.

## Learning outcomes from the project:

- Even though I was aware of the theoretical aspects of cloud services before and had implemented parts of the services for IST 615 assignments, but I never had the chance to practically implement a whole project using cloud services which was very insightful.
- Integration of various cloud services to implement a project.