

# Final Project Demonstration

Vigneshwar Ravirao

# Introduction

Cars24 is a leading AutoTech company focused on the sale, purchase, and financing of pre-owned cars.

The company offers an online marketplace for buying and selling used cars, complemented by a suite of services including car financing, quality checks, warranties, and seamless documentation for transactions.

Cars24 primarily serves the automotive industry with a customer base looking for pre-owned vehicle solutions.



# Problem Statement

The main idea I had behind using this dataset was to try and find some way to predict the selling price of a used car based on brand, model, age, no of previous owners, fuel type, kilometers driven and transmission type.

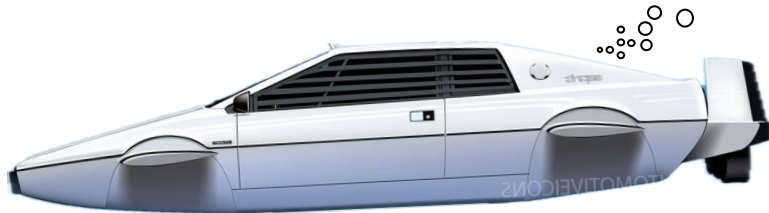
Consider a scenario like while adding a new record of a used car data, someone should make an assessment of the car and figure out what the selling price should be.

The goal of this project is to automate this task using linear regression.



# Initial EDA

- Removed rows that contains null values
- Explored each column
  - **car\_brand** : 26 different brands
  - **model** : 902 different models
  - **age**(from year) : [3, 17]
  - **fuel** : Petrol, Diesel, LPG, CNG, Electric
  - **km\_driven** : [179, 912380]
  - **gear** : Manual & Automatic
  - **ownership** : 1, 2, 3, 4
  - **price** : [91000, 6500000]
  - **monthly\_payment** : [2024, 77744]



# Hypothesis Testing

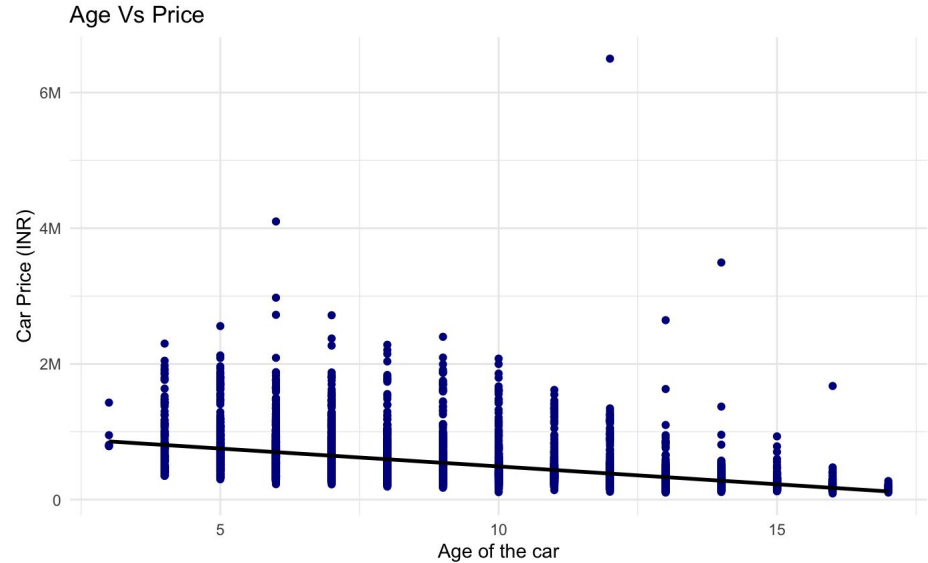
**Tested hypothesis:** Does price of the car decreases as the age increase ?

**P Value:**  $2.2 \times 10^{-16}$

$\alpha$  : 0.05

$\alpha >> \text{p value}$

**Result:** Yes it does !



# Hypothesis Testing

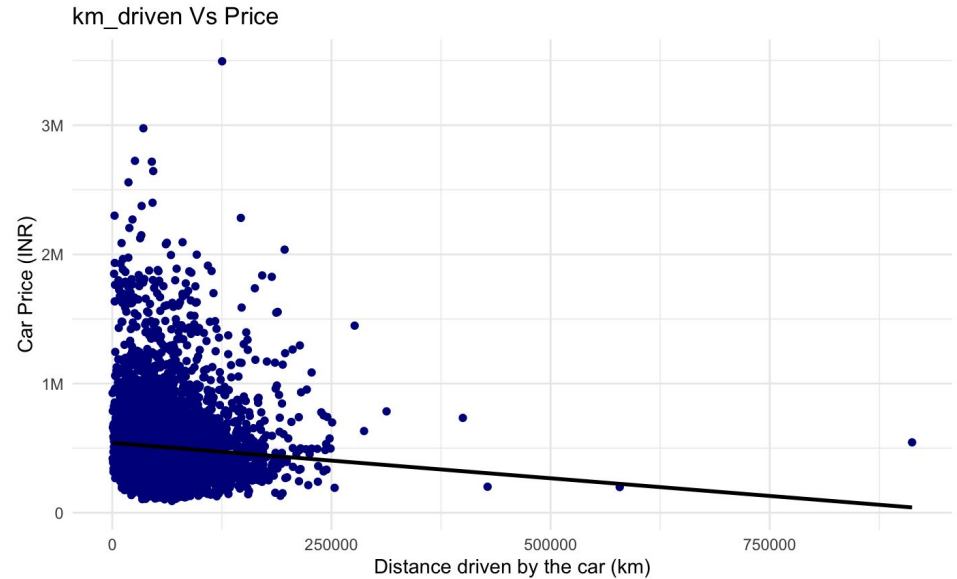
**Tested hypothesis:** Does price of the car decreases as the km\_driven increase ?

**P Value:**  $8.03 \times 10^{-8}$

$\alpha$  : 0.05

$\alpha \gg p$  value

**Result:** Yes it does !



# Preprocessing the data

## **Converting categorical variables into numerical**

- One hot encoding
- Target encoding

## **Scaling the data**

This is essential for maintaining consistent relationship between the features and improving model performance. Consider the below example

- age ranges between 3 to 17
- km\_driven ranges between 179 to 912380

Since the features have a massive difference on their range, it is better to have these values in scale.

# Before preprocessing

car_brand	model	price	year	location	fuel	km_driven	gear	ownership	emi
Hyundai	EonERA PLUS	330399	2016	Hyderabad	Petrol	10674	Manual	2	7350
Maruti	Wagon R 1.0LXI	350199	2011	Hyderabad	Petrol	20979	Manual	1	7790
Maruti	Alto K10LXI	229199	2011	Hyderabad	Petrol	47330	Manual	2	5098
Maruti	RitzVXI BS IV	306399	2011	Hyderabad	Petrol	19662	Manual	1	6816
Tata	NanoTWIST XTA	208699	2015	Hyderabad	Petrol	11256	Automatic	1	4642
Maruti	AltoLXI	249699	2012	Hyderabad	Petrol	28434	Manual	1	5554
Maruti	AltoLXI	240599	2011	Hyderabad	Petrol	31119	Manual	1	5352
Maruti	Alto K10LXI	191999	2010	Hyderabad	Petrol	10910	Manual	1	4271
Honda	Brio1.2 S MT I VTEC	362299	2013	Hyderabad	Petrol	40362	Manual	2	8059
Maruti	Wagon R 1.0VXI	385799	2013	Hyderabad	Petrol	15673	Manual	2	8582



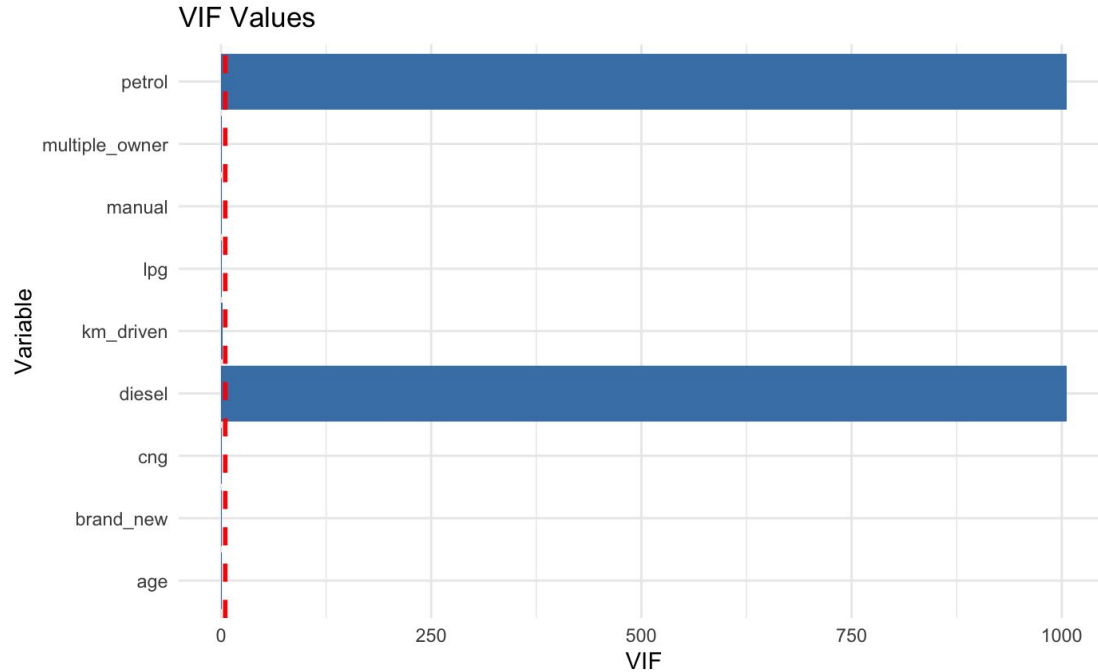
# After preprocessing

	car_brand	brand_new	petrol	diesel	lpg	cng	age	km_driven	manual	multiple_owner	price
1	Hyundai	0.1526935	1	0	0	0	0.3571429	0.011505140	1	1	0.03735357
2	Maruti	0.1257020	1	0	0	0	0.7142857	0.022801992	1	0	0.04044297
3	Maruti	0.1257020	1	0	0	0	0.7142857	0.051689266	1	1	0.02156327
4	Maruti	0.1257020	1	0	0	0	0.7142857	0.021358231	1	0	0.03360883
5	Tata	0.2437464	1	0	0	0	0.4285714	0.012143157	0	0	0.01836464
6	Maruti	0.1257020	1	0	0	0	0.6428571	0.030974533	1	0	0.02476190
7	Maruti	0.1257020	1	0	0	0	0.7142857	0.033917963	1	0	0.02334202
8	Maruti	0.1257020	1	0	0	0	0.7857143	0.011763855	1	0	0.01575893
9	Honda	0.1647287	1	0	0	0	0.5714286	0.044050598	1	1	0.04233094
10	Maruti	0.1257020	1	0	0	0	0.5000000	0.016985292	1	1	0.04599766

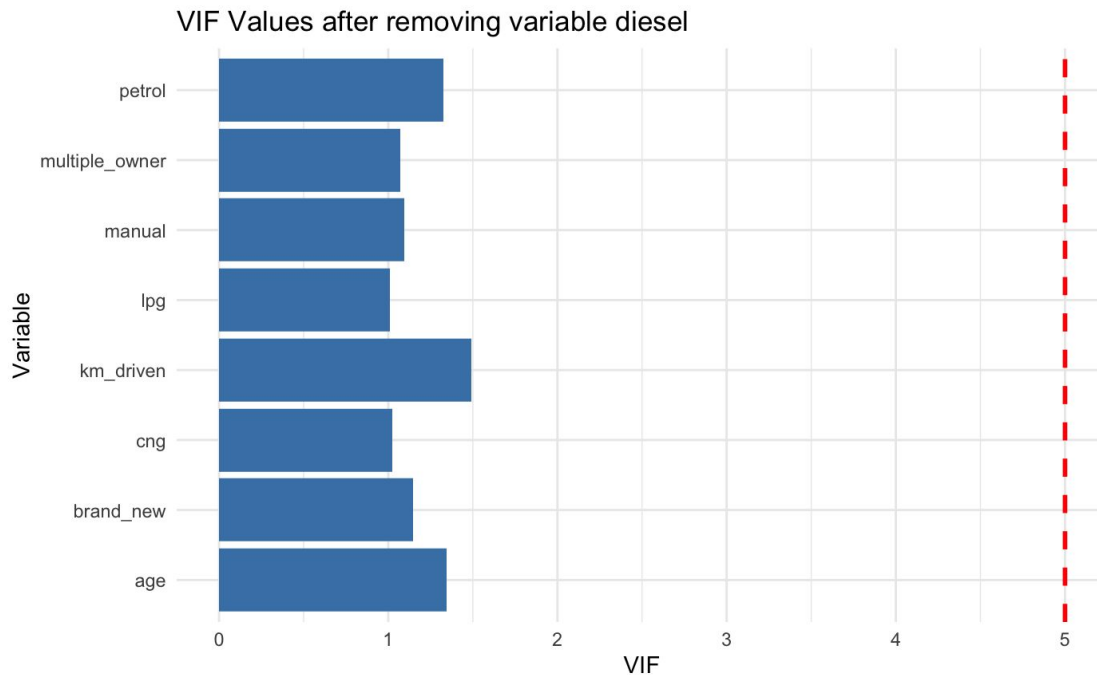
# Linear Regression Model

$$\begin{aligned}\text{Price} = & 0.062 + 0.443 \times \text{brand} - 0.184 \times \text{age} + 0.121 \times \text{petrol} \\ & + 0.163 \times \text{diesel} + 0.015 \times \text{lpg} - 0.012 \times \text{cng} \\ & - 0.067 \times \text{km\_driven} - 0.061 \times \text{manual} - 0.001 \times \text{multiple\_owner}\end{aligned}$$

# Multicollinearity check



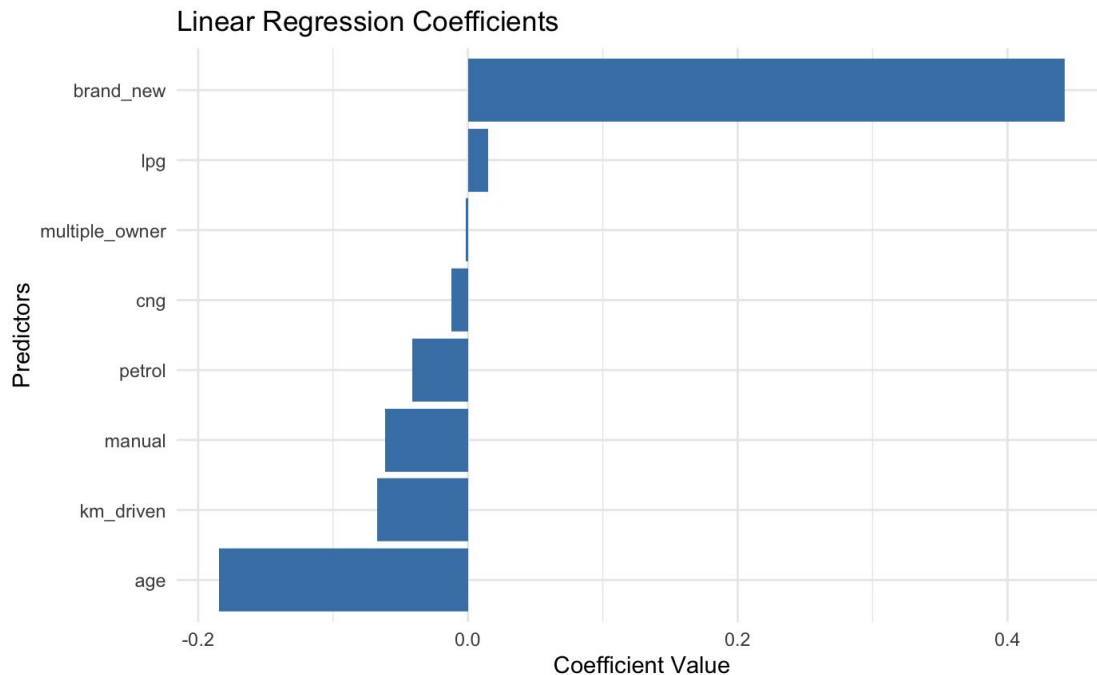
# Multicollinearity check



# New Linear Regression Model

$$\begin{aligned} \textit{Price} = & 0.225 + 0.442 \times \textit{brand} - 0.184 \times \textit{age} - 0.041 \times \textit{petrol} \\ & + 0.015 \times \textit{lpg} - 0.012 \times \textit{cng} - 0.067 \times \textit{km\_driven} \\ & - 0.061 \times \textit{manual} - 0.001 \times \textit{multiple\_owner} \end{aligned}$$

# Coefficients of new LR model1



# R - Squared value

65% of the variance in car price can be explained by the independent variables

This value makes sense, because we haven't considered variables like model, location and monthly\_payment.

```
{r}  
summary(model1_new)$r.squared  
  
[1] 0.6572958
```

# Model Evaluation



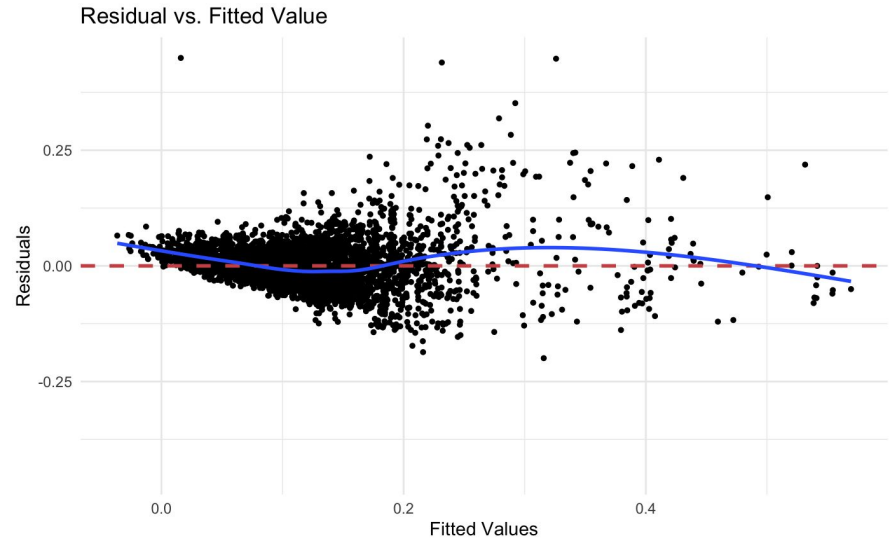


# Residual vs Fitted value

Curvature suggests possible nonlinearity issues.

Spread indicates heteroscedasticity in residuals.

Outliers may influence model results

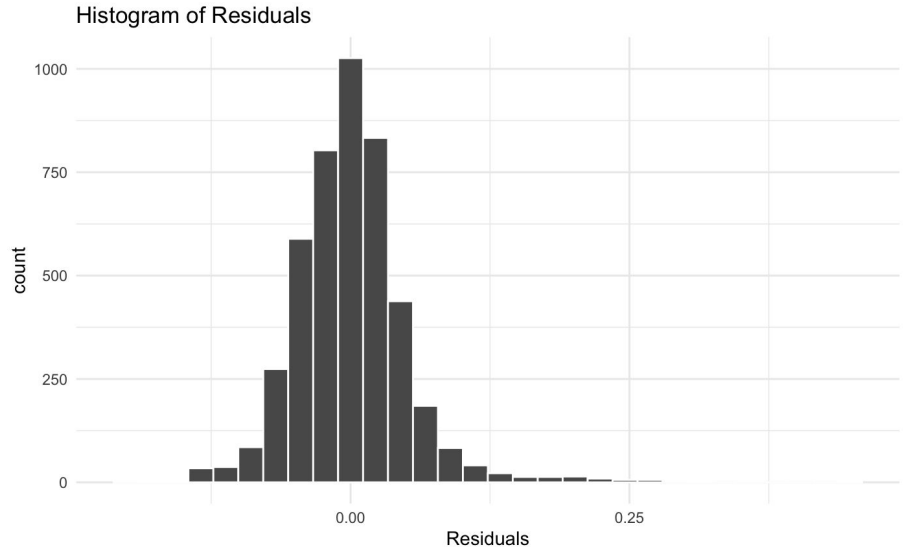


# Residual Histogram

Residuals show approximate normal distribution shape.

Outliers appear at the distribution's tails.

Minor skewness suggests slight model deviation.

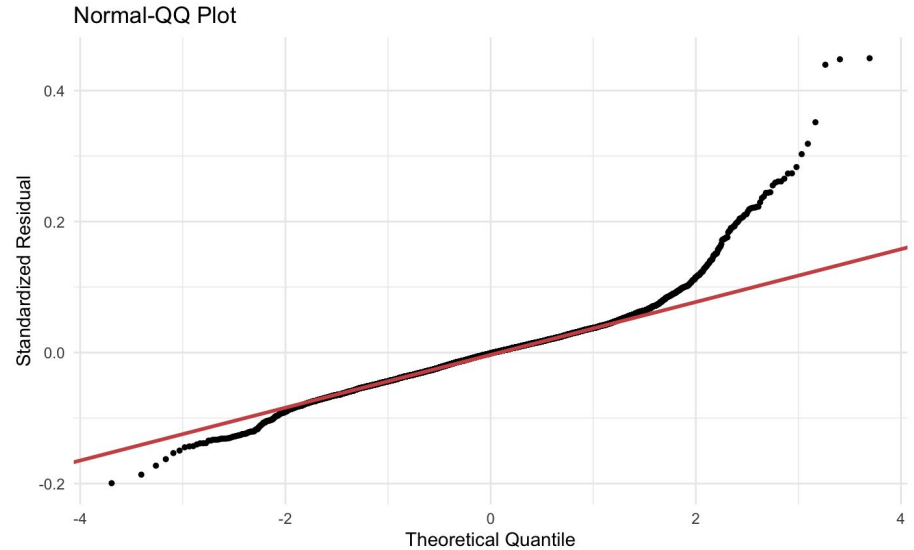


# QQ-Plots

Deviations occur at extreme quantiles.

Skewness observed in tails.

Potential outliers at upper quantile extremes.



# Root Mean Square Error

If RMSE is 0, the model's predictions are perfect.

A lower RMSE indicates better model performance, but it should be compared to other models or benchmarks for context.

```
print(rmse)
```

```
[1] 0.05099379
```

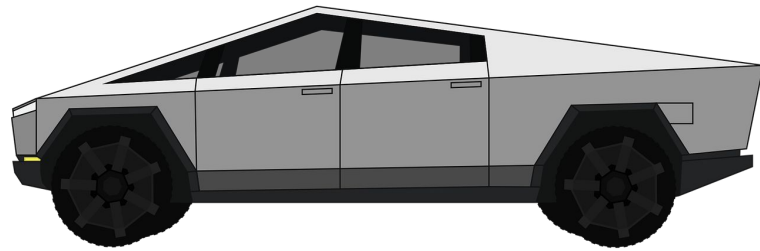
# Conclusion

- Older cars significantly lower the price.
- More distance driven by a car, the lower the price.
- Petrol and manual reduce car price substantially.
- Model explains 65% variation in car price.
- This method can be used by any companies who want to automate finding the price of a car based on different explanatory variables discussed above.

# Applying the model to test set

car_brand <chr>	price <dbl>	predicted_price_1 <dbl>
1-10 of 10 rows	249699	285770.4
Maruti	240599	240231.3
Hyundai	401599	441284.8
Maruti	241599	241316.6
Maruti	356099	505644.5
Hyundai	401699	402981.1
Maruti	238099	284154.1
Maruti	613499	551834.6
Hyundai	222699	114447.2
Hyundai	346399	258846.8

Out of Curiosity



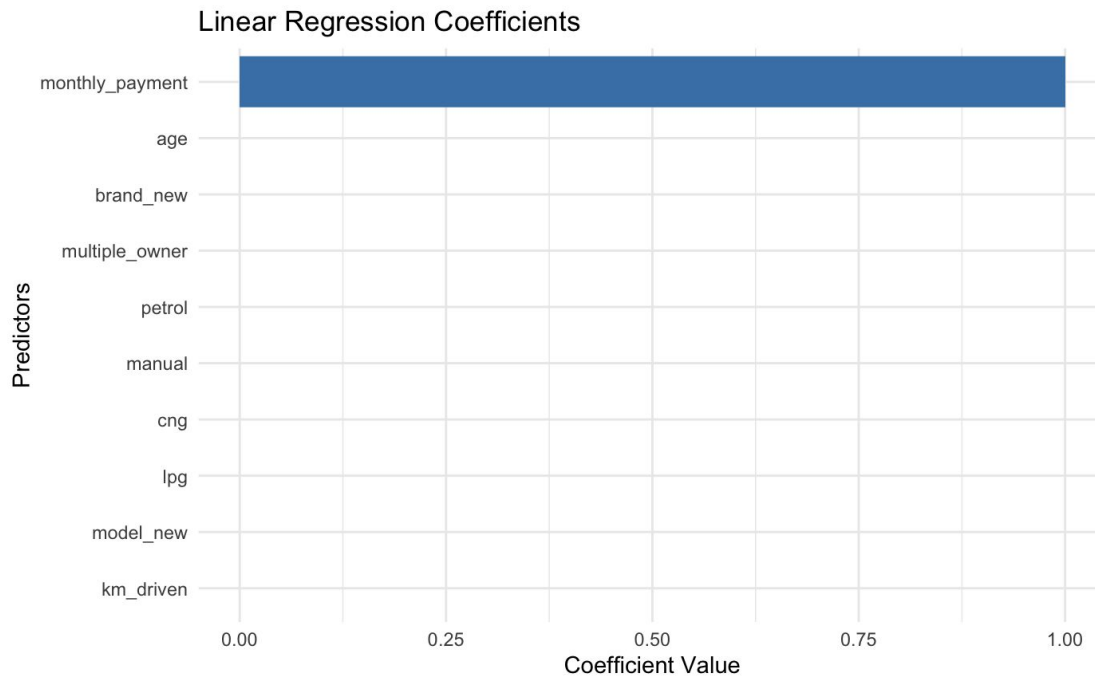
# Linear Regression Model 2

I have created a another linear regression model incorporating car's model and monthly\_payment

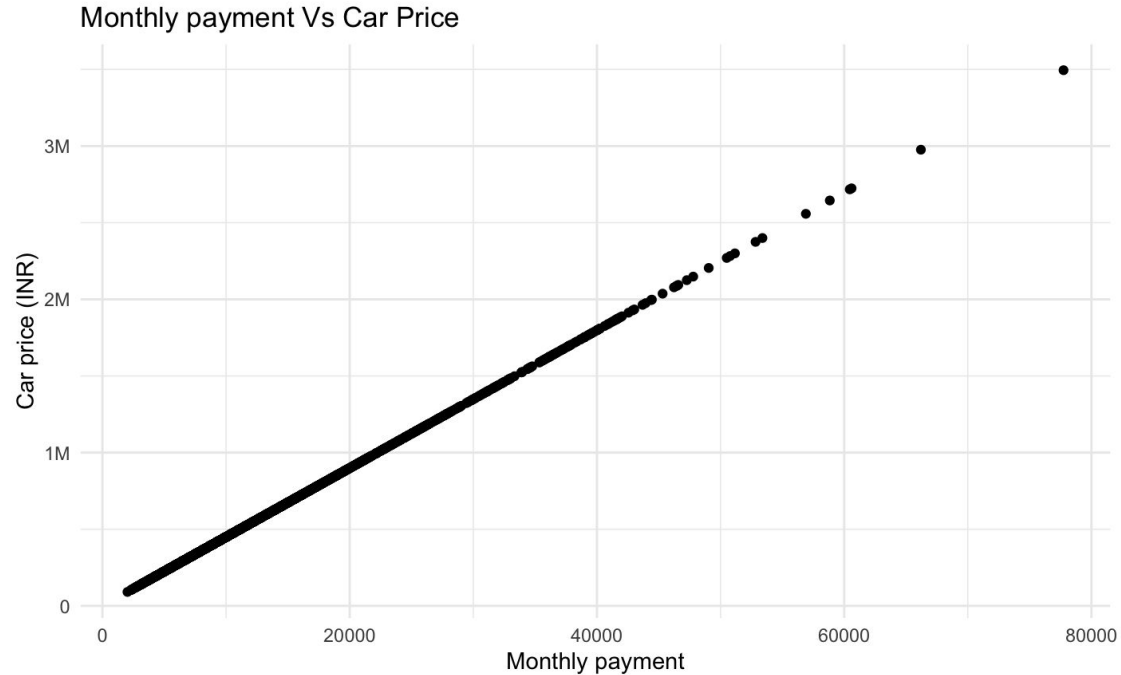
$$\begin{aligned} \textit{Price} = & 0.000003 + 0.0000002 \times \textit{brand} - 0.0000007 \times \textit{model} - 0.0000004 \times \textit{age} \\ & - 0.000000008 \times \textit{petrol} - 0.0000004 \times \textit{lpg} - 0.0000002 \times \textit{cng} \\ & - 0.000001 \times \textit{km\_driven} - 0.0000001 \times \textit{manual} \\ & + 0.00000002 \times \textit{multiple\_owner} + 0.999 \times \textit{monthly\_payment} \end{aligned}$$



# Coefficients of LR model2



# Why it has coeff as 0.999?

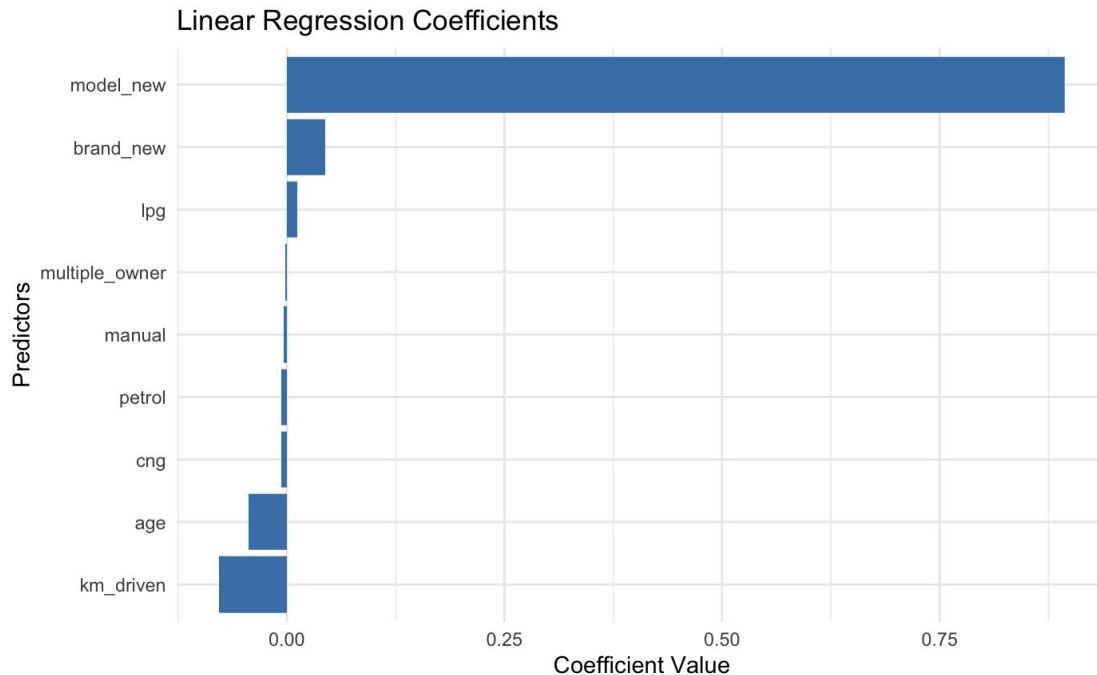


# New Linear Regression Model 2

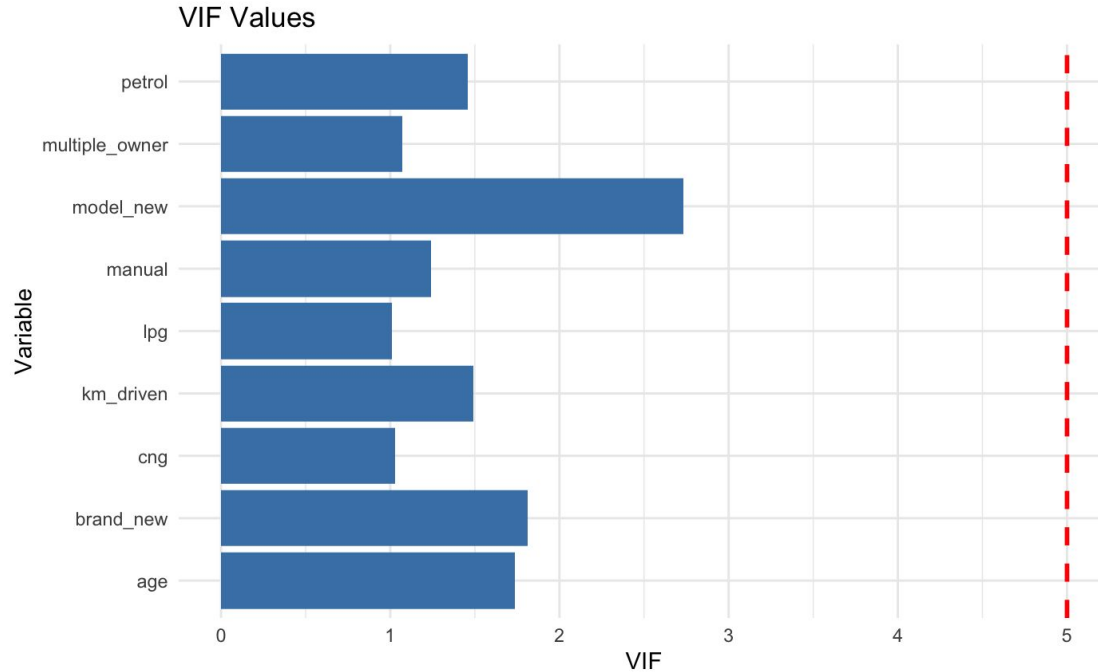
Dropping the monthly\_payment from the current linear regression model.

$$\begin{aligned} \textit{Price} = & 0.047 + 0.044 \times \textit{brand} + 0.89 \times \textit{model} - 0.044 \times \textit{age} \\ & - 0.006 \times \textit{petrol} + 0.011 \times \textit{lpg} - 0.006 \times \textit{cng} \\ & - 0.078 \times \textit{km\_driven} - 0.004 \times \textit{manual} - 0.002 \times \textit{multiple\_owner} \end{aligned}$$

# Coefficients of new LR model2



# Multicollinearity check



# R - squared comparison

model1_new	model2_new
0.657	0.932

This shows us that the second model captures 28% more variability compared to the first model.

# RMSE comparison

model1_new	model2_new
0.050	0.022

This shows us that the second model is a better compared to the first one.

# Applying the model to test set

<b>car_brand</b> <chr>	<b>price_</b> <dbl>	<b>predicted_price_2</b> <dbl>
Maruti	249699	224233.2
Maruti	240599	212745.2
Hyundai	401599	395449.9
Maruti	241599	214006.5
Maruti	356099	338990.3
Hyundai	401699	394005.8
Maruti	238099	222354.7
Maruti	613499	513521.1
Hyundai	222699	155263.5
Hyundai	346399	259325.8



# Questions ?



Thank you !

