



# A Regression Study of Price Determinants of used cars in US Automobile markets

MENTORED BY:
MS. VIBHA SANTHANAM

BY: GROUP 5				
ANKUR KUMAR	MUSKAAN PASSI			
LAKSHYA SHARMA	KRITAGYA KASHYAP			
ADITYA KAPUR	ABHINAV TYAGI			
JASWANT SINGH				

#### PROBLEM STATEMENT

A Study of Price Determinants of used cars in US Automobile markets

#### BUSINESS IMPACT

A used car is a pre-owned or second hand vehicle, which is repaired and refurbished to regain working conditions before the sale. There is a need for a used car price prediction system to effectively determine the worthiness of the car using a variety of features.



#### **DATA OVERVIEW**

# **Used Cars Dataset Vehicles Listing From Craigslist.org**

Cat <mark>ego</mark> rio	cal Columns FIR	Numerical Columns	Unique	Date and Time
<ul> <li>Region</li> <li>Manufacturer</li> <li>Model</li> <li>Condition</li> <li>Cylinders</li> <li>Fuel</li> <li>Title Status</li> <li>Region URL</li> </ul>	<ul> <li>Transmission</li> <li>Drive</li> <li>Size</li> <li>Type</li> <li>Paint Color</li> <li>State</li> </ul>	<ul><li>Price</li><li>Lat</li><li>Long</li><li>Odometer</li><li>County</li><li>Year</li></ul>	<ul> <li>I.D.</li> <li>URL</li> <li>Image URL</li> <li>VIN</li> <li>Description</li> </ul>	Posting Date



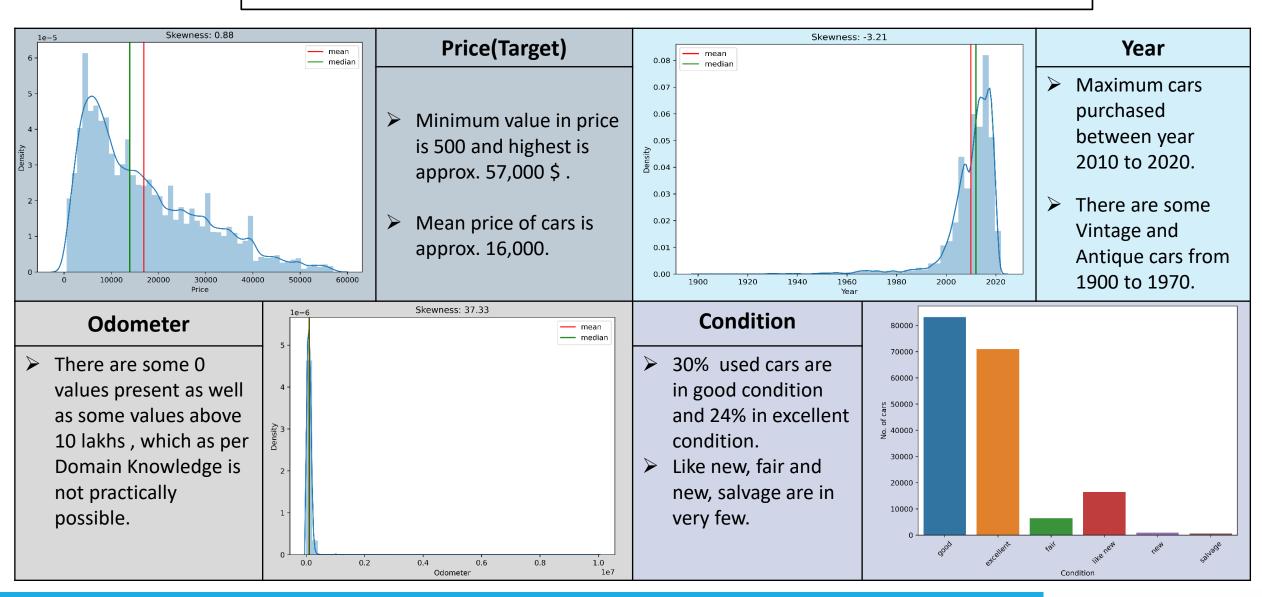
# Data Pre-Processing and Feature Engineering

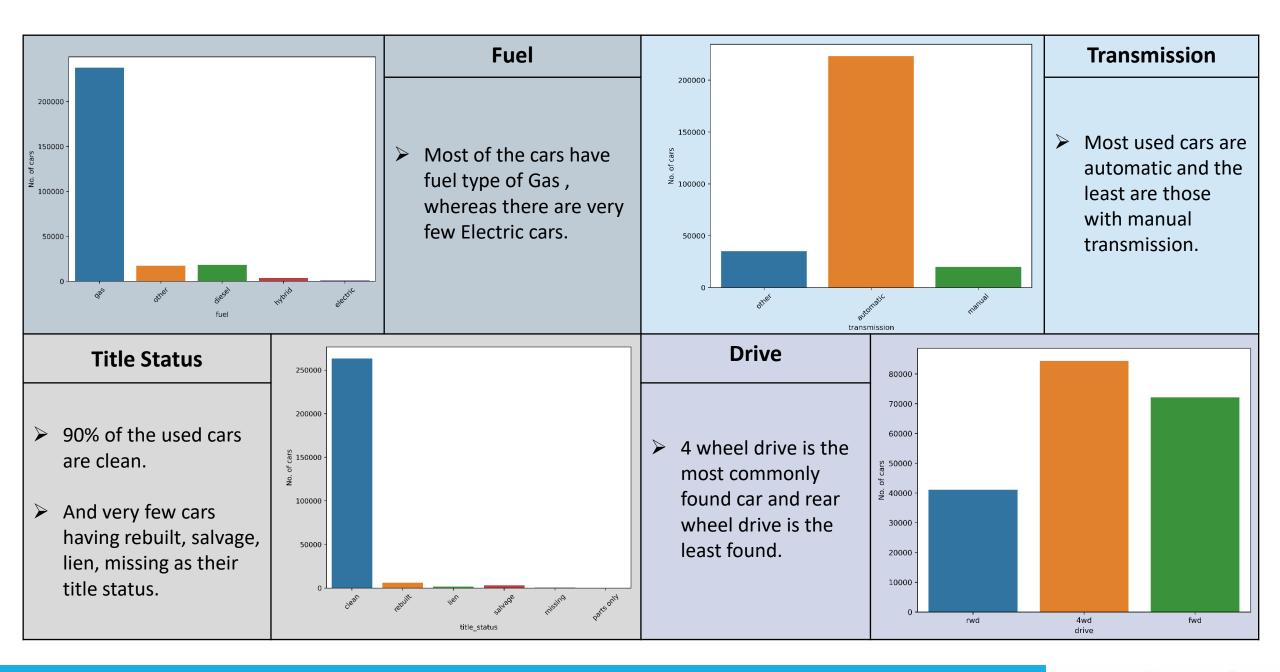
Removing Redundant Columns		Columns with too-many Sub categories		
<ul> <li>I.D.</li> <li>Region Url</li> <li>Image Url</li> <li>Lat</li> </ul>	FOR SALE	<ul> <li>Region (404)</li> <li>State (51)</li> <li>Model (29,668)</li> <li>Region URL (413)</li> </ul>		
<ul><li>Long</li><li>County</li></ul>		➤ Manufacturer (43)		

Columns that have Outliers		Columns which are converted into bins
> Odometer	>	State (into State Zone)
> Price	>	Paint colour
> Year	>	Year (into Vint car)
	> Odometer > Price	> Odometer > Price

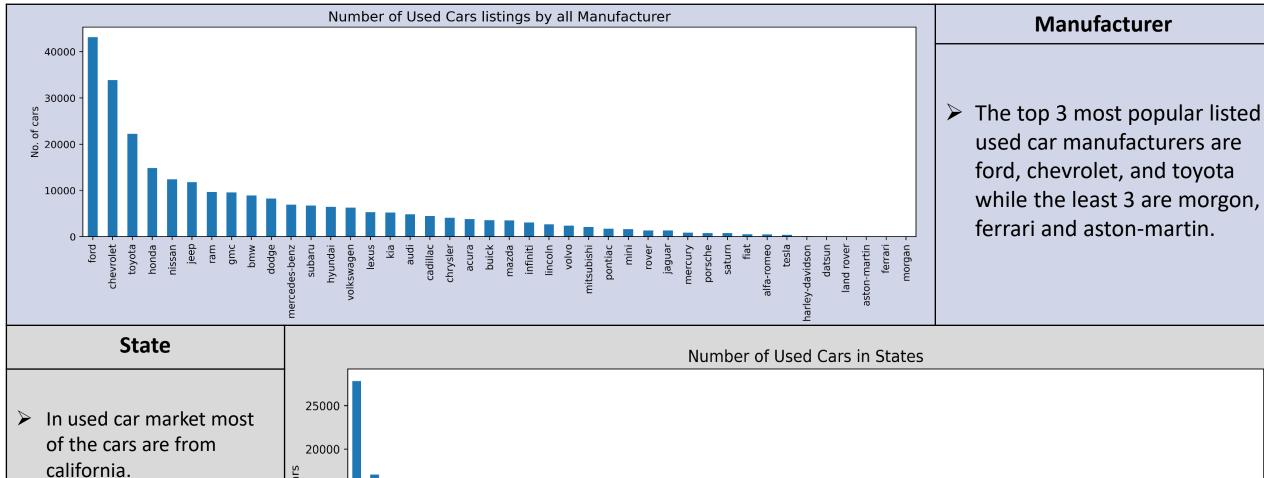


#### **EXPLORATORY DATA ANALYSIS** – Univariate Analysis









ca fl tx ny oh mi pa or co nc il wi wa tn nj az ma ia mn id va ga ks ok sc in ct mdnm mt al ky mo nv dc ar nh me hi ak la ri vt sd ut ne de ms wv wy nd

15000

10000

5000

Top 5 states are California,

Florida, Texas, New York, Ohio in used car market.

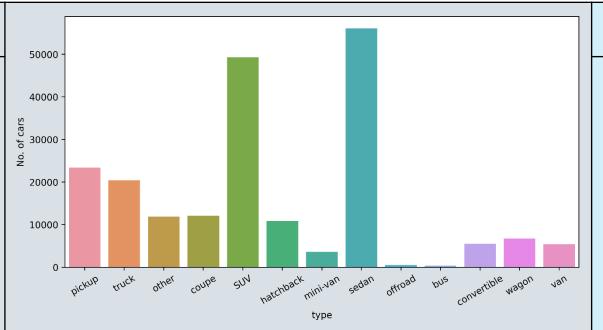
greatlearning

Learning for Life

#### Top 20 regions based on listing of Used Cars Region **Size** 2500 2000 No. of Cars > Full size cars have the most 1000 Central Nj have highest count followed by mid size 500 number of listed used cars. and sub-compact size are found the least. **Paint Colour** 50000 40000 White and Black are 40000 most popular No. of Cars colours in used cars. ∮ 20000 J 20000 ➤ Where Purple, 10000 10000 Orange & yellow are least popular. white blue red black silver full-size mid-size compact sub-compact size

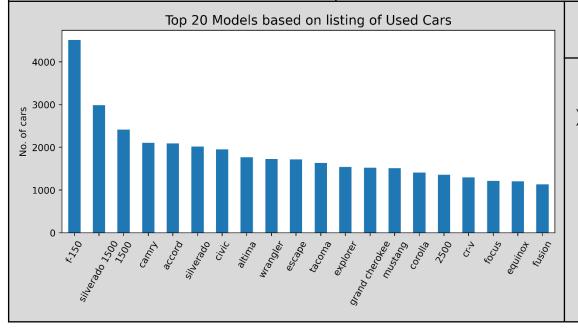
#### Type

- Sedan and SUVs are the most found cars in Used Car market.
- ➤ Where offroad and bus are the least in used car market.



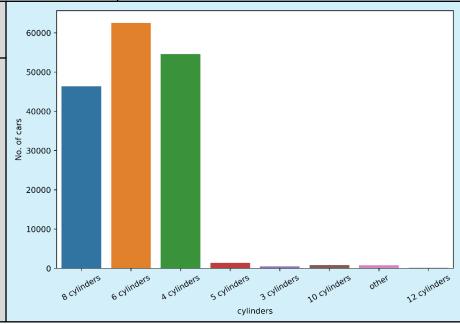
#### **Cylinders**

- Cars having 6, 4 and 8 cylinders occupy the most used car market shares.
- Cars having 3, 5, 10, 12 are very less car market shares.



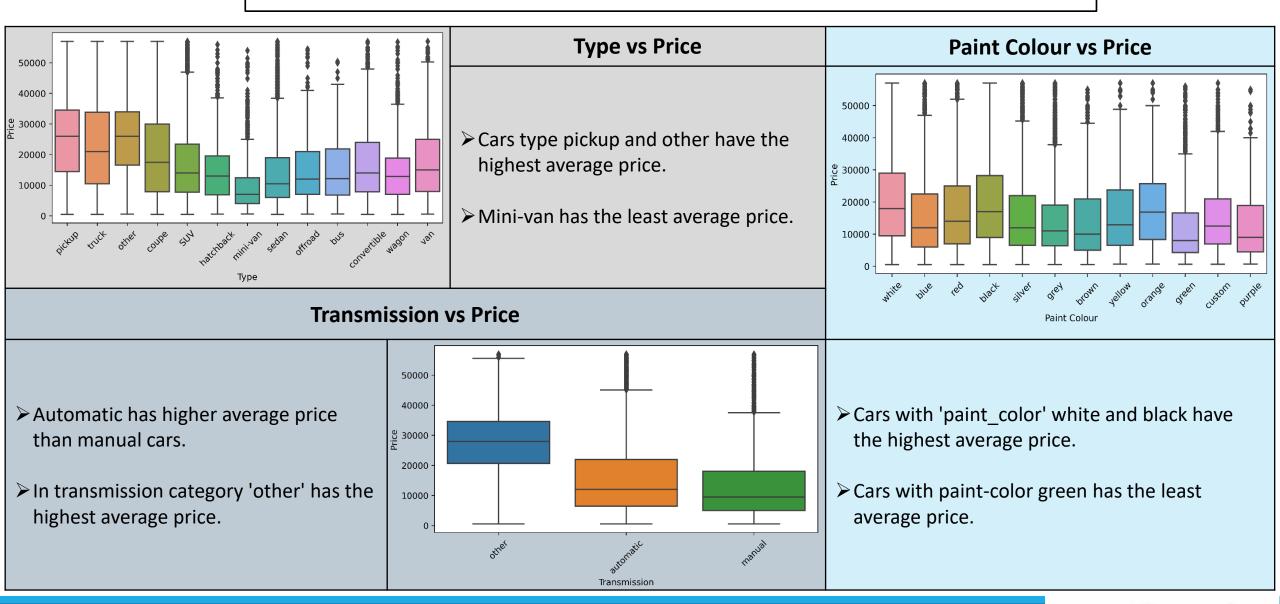
#### **Models**

➤ F150 car model is the most found model in used car market and it belongs to Ford.

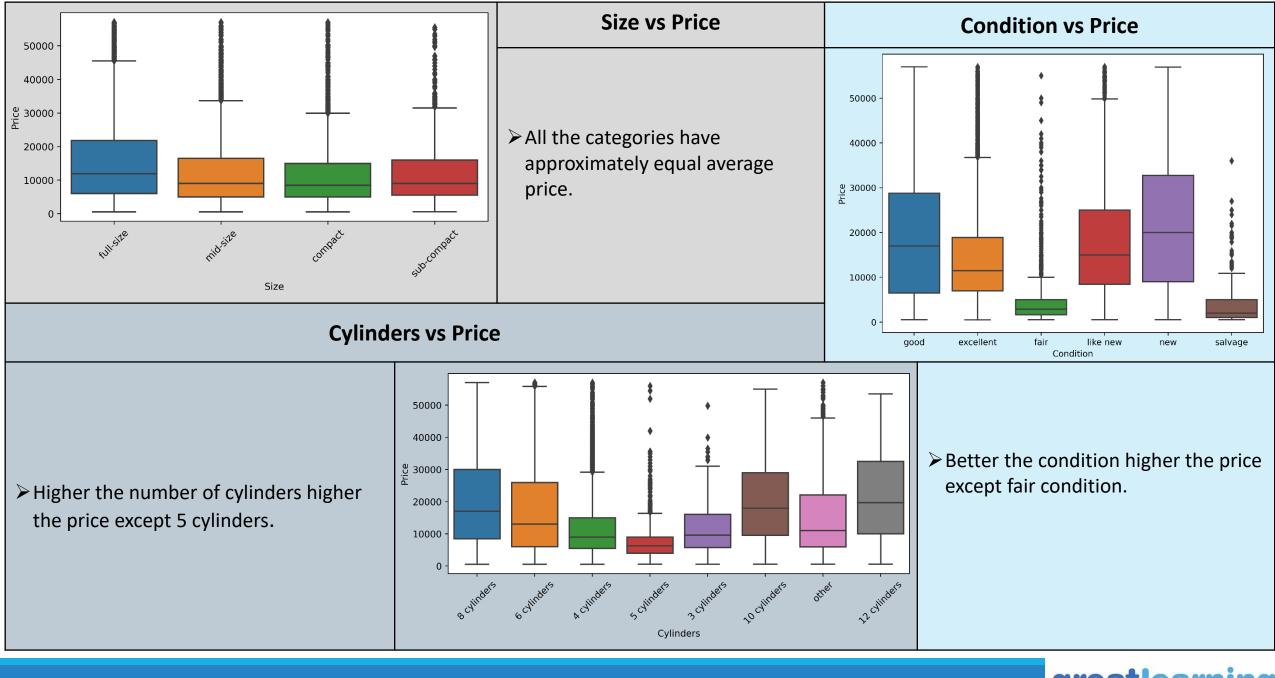


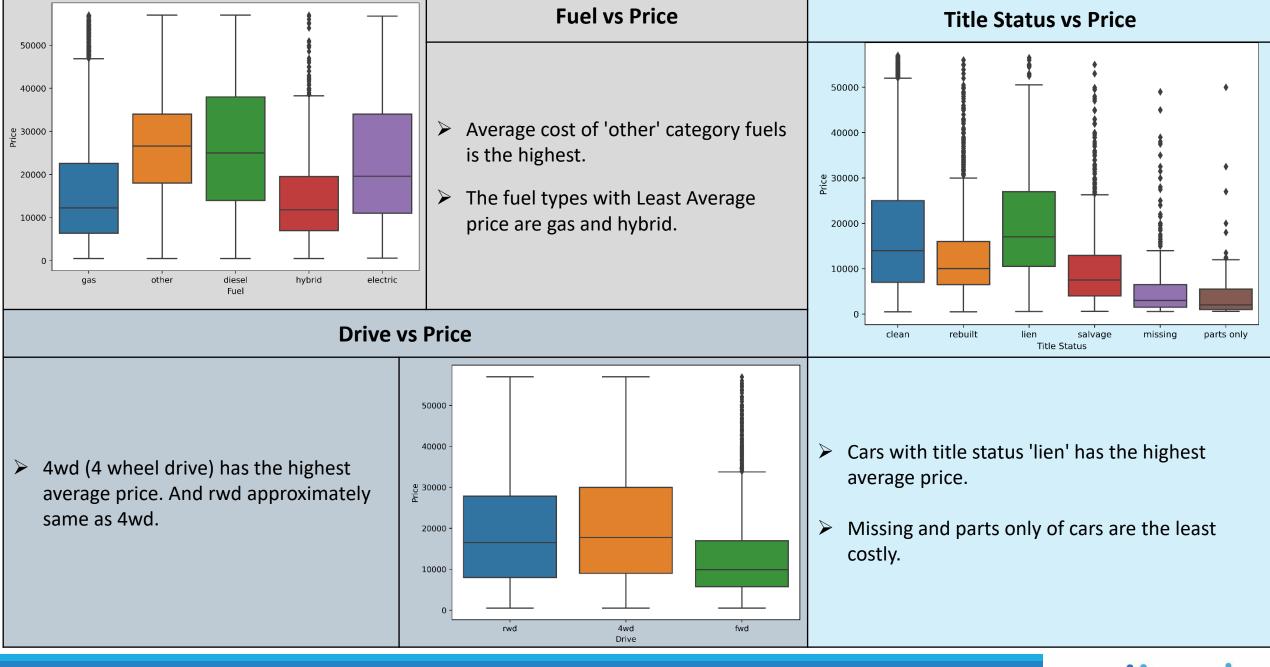


#### **Bivariate Analysis - With Respect to Target (Price)**

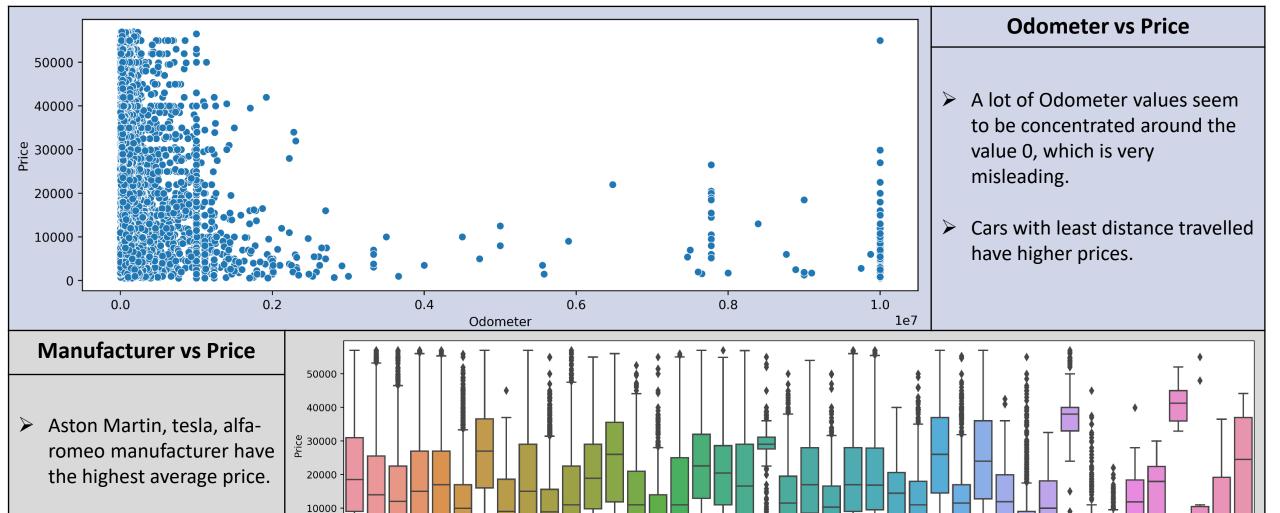












Manufacturer

Saturn, Mercury and pontiac manufacturer

have the least average

price.

toyota

honda dodge lexus jaguar buick chrysler



# **Missing Values**

Total		Percentage of Missing Values	Total Percentage of Missi		Percentage of Missing Values
county	<b>county</b> 426880 100.000000		odometer	4400	1.030735
size	306361	71.767476	fuel	3013	0.705819
cylinders	177678	41.622470	transmission	2556	0.598763
condition	174104	40.785232			
VIN	161042	37.725356	year	1205	0.282281
drive	130567	30.586347	30.586347 image_url	68	0.015930
paint_color	130203	30.501078	posting_date	68	0.015930
type	92858	21.752717	url	0	0.000000
manufacturer	17646	4.133714	33714 price	0	0.000000
title_status	8242	1.930753	state	0	0.0000000
lat	6549	1.534155	region_url	0	0.0000000
long	6549	1.534155	region	0	0.000000
model	5277	1.236179	id	0	0.000000



#### MISSING VALUES TREATMENT

- ➤ **Year** The values were filled using →
  - 1) URL
  - 2) VIN (using vin function in python)
  - 3) Remaining values imputed with mode year values.
- ➤ Manufacturer Filled using →
  - 1) VIN
  - 2) Model
  - 3) URL
  - 4) Remaining Values Binned as Unknown
- Condition Filled using >
  - 1) Title Status
  - 2) Rest using python Forward Fill function.
- ➤ Odometer Filled using →
  - 1) All values greater than 999999 imputed as 999999.
  - 2) Null and 0 values filled with mean values basis the model .
  - 3) Remaining values filled using median.

- Imputed mode values on basis Model & Manufacturer are: Transmission, Type, Fuel.
- Imputed mode values on the basis Model, Manufacturer & Type are: Drive, Size, Paint colour.
- Cylinders The values were filled using >
  - 1) Electric cars with other.
  - 2) Description.
  - 3) Model.
  - 4) Manufacturer.



# **Statistical Tests** (Predictor vs Target)

S	Significant Variable	es	Non significant variables			
Features	Test	P value	Features	Test	P value	
Year	Mann-Whitney U	0.0	Region	Anova	1.0	
Odometer	Mann-Whitney U 0.0		Manufacturer	Anova	1.0	
Condition	Anova 0.000001		Cylinder	Anova	0.390706	
State	Anova	0.032514	Fuel	Anova	0.964539	
Is tax	T-test independent	0.009575	Transmission	Anova	0.270099	
Age	T-test independent 0.0005157		Drive	Anova	0.057264	
Vint_car	Anova 0.000009		Size	Anova	0.208265	
N 4 o d o l	Due to high categories we can't run statistical test and from domain knowledge we keep this		Туре	Anova	0.388538	
Model	variable for modelling.	iowicage we keep tilis	Paint colour	Anova	0.294331	
			State Zone	Anova	0.588757	



# **Correlation** (Predictor vs Target)

Feature	Correlation	
Region	-0.072851	
Price	1.000000	
Year	0.344641	
Manufacturer	-0.078222	
Model	-0.004239	
Condition	0.028708	
Cylinders	0.244160	
Odometer	-0.412041	
Size	0.184555	
Туре	0.035677	

Feature	Correlation		
State	-0.000448		
State_Zone	0.068220		
Age	-0.344641		
is_tax	0.001162		
fuel_diesel	0.184503		
fuel_electric	0.027906		
fuel_gas	-0.252713		
fuel_hybrid	-0.021164		
fuel_other	0.184598		
transmission_automatic	-0.227737		

Feature	Correlation
transmission_manual	-0.080006
transmission_other	0.333662
drive_4wd	0.212183
drive_fwd	-0.287314
drive_rwd	0.083461
vint_car_antique	-0.002434
vint_car_classic	-0.207867
vint_car_modern	0.186869
vint_car_vintage	-0.000414



# **DATA MODELLING**

- BASE MODEL: Decision Tree Regressor
- Let's see final comparison for the performance of all the models that we have tried to make the final conclusion through Data Frame:

	Model Name	Train R2 Score	Test R2 Score	Train RMSE	Test RMSE	MAPE
0	Decision Tree	0.999604	0.696382	247.414093	6816.527200	40.538615
1	Ridge	0.551710	0.550717	8323.137291	8291.999291	65.752941
2	Lasso	0.551695	0.550725	8323.279008	8291.927331	65.720329
3	Elastic Net	0.420750	0.422053	9461.076248	9404.668752	62.097979
4	Random Forest	0.978557	0.851884	1820.349500	4761.028578	43.168481
5	AdaBoost	0.635130	0.635568	7508.906508	7468.054537	60.193155
6	Gradient Boost	0.781405	0.783581	5812.027572	5755.016932	50.773069
7	XGBoost	0.856421	0.837678	4710.345448	4984.118219	42.128192
8	CatBoost	0.851534	0.841018	4789.844196	4932.564387	42.989694
9	Tuned CatBoost	0.882055	0.851753	4269.212741	4763.129450	41.563495

#### **Final Insights**

- ✓ From this table, we can observe that CatBoost model is our best fit model.
- ✓ After parameter tuning of CatBoost our Tuned CatBoost Model is our Final Model.
- ✓ Tuned CatBoost has given best R-Squared i.e. 85.17%.

Regression

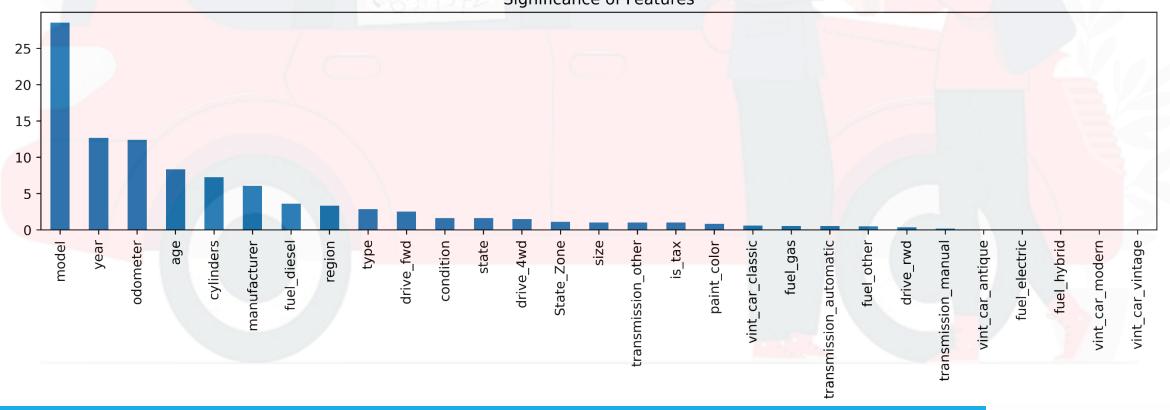


#### **Feature Importance**

## Bar plot to check the importance of the features.

1st most significant feature : Model2nd most significant feature : Year

Significance of Features



### **Conclusion**

- 1. Used cars with high Mileage are cheaper.
- 2. Used cars with better Appearance are expensive.
- 3. Used cars which come from big manufacturer are cost more.
- 4. Used cars with pickup or truck or coupe or convertible type are cost more.
- 5. Used cars with white or black paint-color are expensive.
- 6. Used cars with more cylinders are expensive.
- 7. Used cars with higher vehicle age are cost less.



