

🚗 Car Price Prediction 🚗

```
In [15]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from scipy import stats
import seaborn as sns
import pylab
import warnings
warnings.filterwarnings("ignore")
import matplotlib inline
sns.set(style="darkgrid", font_scale=1.5)
pd.set_option("display.max.columns",None)
pd.set_option("display.max.rows",None)

from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor, AdaBoostRegressor

In [16]: import os
working_directory = os.getcwd()
print(working_directory)
/Users/ishu/Desktop

In [25]: path = working_directory + '/car.csv'
df = pd.read_csv(path)
df.head()
```

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase	carlength	carwidth	carheight	curbweight	enginetype	cylindernumber	enginesize	fuelsystem
0	1	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8	64.1	48.8	2548	dohc	four	130	mpfi
1	2	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8	64.1	48.8	2548	dohc	four	130	mpfi
2	3	1	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	171.2	65.5	52.4	2823	ohcv	six	152	mpfi
3	4	2	audi 100 ls	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	54.3	2337	ohc	four	109	mpfi
4	5	2	audi 100ls	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	54.3	2824	ohc	five	136	mpfi

```
In [26]: df.tail()
```

	car_ID	symboling	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase	carlength	carwidth	carheight	curbweight	enginetype	cylindernumber	enginesize	fuelsystem
200	201	-1	volvo 140g	gas	std	four	sedan	rwd	front	109.1	188.8	68.9	55.5	2952	ohc	four	141	mpfi
201	202	-1	volvo 140g	gas	turbo	four	sedan	rwd	front	109.1	188.8	68.8	55.5	3049	ohc	four	141	mpfi
202	203	-1	volvo 140g	gas	std	four	sedan	rwd	front	109.1	188.8	68.9	55.5	3012	ohcv	six	173	mpfi
203	204	-1	volvo 246	diesel	turbo	four	sedan	rwd	front	109.1	188.8	68.9	55.5	3217	ohc	six	145	mpi
204	205	-1	volvo 246g	gas	turbo	four	sedan	rwd	front	109.1	188.8	68.9	55.5	3062	ohc	four	141	mpfi

```
In [27]: #Checking Dimensions of the Data.
df.shape
(285, 26)
```

```
In [28]: #Checking the basic information of dataset.
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 285 entries, 0 to 284
Data columns (total 26 columns):
#   Column               Non-Null Count  Dtype
---  --
0   car_ID               285 non-null    int64
1   symboling            285 non-null    int64
2   CarName              285 non-null    object
3   fueltype            285 non-null    object
4   aspiration           285 non-null    object
5   doornumber          285 non-null    object
6   carbody             285 non-null    object
7   drivewheel          285 non-null    object
8   engineLocation       285 non-null    object
9   wheelbase           285 non-null    float64
10  carlength            285 non-null    float64
11  carwidth             285 non-null    float64
12  carheight            285 non-null    float64
13  curbweight           285 non-null    int64
14  enginetype           285 non-null    object
15  cylindernumber       285 non-null    object
16  enginesize            285 non-null    int64
17  fuelsystem           285 non-null    object
18  boreratio            285 non-null    float64
19  stroke               285 non-null    float64
20  compressionratio     285 non-null    float64
21  horsepower           285 non-null    int64
22  peakrpm              285 non-null    int64
23  citympg              285 non-null    int64
24  highwaympg           285 non-null    int64
25  price                285 non-null    float64
dtypes: float64(8), int64(8), object(18)
memory usage: 41.8+ KB

* Observation
```

1.)From above output we can observe that 10 categorical & 16 numerical Attributes. 2.)All the features are having correct data-types. So we don't have to do any changes.

```
In [29]: #Descriptive Statistics Analysis.
df.describe()
```

	car_ID	symboling	wheelbase	carlength	carwidth	carheight	curbweight	enginesize	boreratio	stroke	compressionratio	horsepower	peakrpm	citympg	highwaympg	price
count	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000	285.000000
mean	103.000000	0.834146	98.756685	174.049068	65.907805	53.724878	2555.565854	126.907317	3.329756	3.255415	10.142537	104.117073	5125.121951	25.219512	30.751220	13276.71057
std	59.322565	1.245307	6.021776	12.337789	2.145204	2.443522	520.680204	41.542693	0.270844	0.313597	3.973040	39.544167	476.985643	6.542142	6.886443	7988.85233
min	1.000000	-2.000000	86.600000	141.100000	60.300000	47.800000	1488.000000	61.000000	2.540000	2.070000	7.000000	48.000000	4150.000000	13.000000	16.000000	5116.000000
25%	52.000000	0.000000	94.500000	166.300000	64.100000	52.000000	2145.000000	67.000000	3.150000	3.110000	8.000000	70.000000	4800.000000	19.000000	25.000000	7788.000000
50%	103.000000	1.000000	97.000000	175.200000	65.500000	54.100000	2414.000000	120.000000	3.310000	3.290000	9.000000	95.000000	5200.000000	24.000000	30.000000	10295.000000
75%	155.000000	2.000000	102.400000	183.100000	66.800000	55.500000	2935.000000	141.000000	3.580000	3.410000	9.400000	116.000000	5500.000000	30.000000	34.000000	16503.000000
max	284.000000	3.000000	120.900000	208.100000	72.300000	59.800000	4065.000000	326.000000	3.940000	4.170000	23.000000	288.000000	6600.000000	49.000000	54.000000	45450.000000

```
In [30]: #Checking NaN(NULL) values in our Dataset.
df.isnull().sum().to_frame().rename(columns={0:"Total No. of Missing Values"})
```

	Total No. of Missing Values
car_ID	0
symboling	0
CarName	0
fueltype	0
aspiration	0
doornumber	0
carbody	0
drivewheel	0
engineLocation	0
wheelbase	0
carlength	0
carwidth	0
carheight	0
curbweight	0
enginetype	0
cylindernumber	0
enginesize	0
fuelsystem	0
boreratio	0
stroke	0
compressionratio	0
horsepower	0
peakrpm	0
citympg	0
highwaympg	0
price	0

Observation

We can observe that none of the features is having Missing values.

```
In [31]: #Showing Only Categorical Features.
df.select_dtypes(include="object").head()
```

	CarName	fueltype	aspiration	doornumber	carbody	drivewheel	engineLocation	enginetype	cylindernumber	fuelsystem
0	alfa-romero giulia	gas	std	two	convertible	rwd	front	dohc	four	mpfi
1	alfa-romero stelvio	gas	std	two	convertible	rwd	front	dohc	four	mpfi
2	alfa-romero Quadrifoglio	gas	std	two	hatchback	rwd	front	ohcv	six	mpfi
3	audi 100 ls	gas	std	four	sedan	fwd	front	ohc	four	mpfi
4	audi 100ls	gas	std	four	sedan	4wd	front	ohc	five	mpfi

```
In [32]: #Showing only the Numerical Features.
df.select_dtypes(include=["int","float"]).head()
```

	car_ID	symboling	wheelbase	carlength	carwidth	carheight	curbweight	enginesize	boreratio	stroke	compressionratio	horsepower	peakrpm	citympg	highwaympg	price
0	1	3	88.6	168.8	64.1	48.8	2548	130	3.47	2.68	9.0	111	5000	21	27	13495.0
1	2	3	88.6	168.8	64.1	48.8	2548	130	3.47	2.68	9.0	111	5000	21	27	16500.0
2	3	1	94.5	171.2	65.5	52.4	2823	152	2.68	3.47	9.0	154	5000	19	26	16500.0
3	4	2	99.8	176.6	66.2	54.3	2337	109	3.19	3.40	10.0	102	5500	24	30	13950.0
4	5	2	99.4	176.6	66.4	54.3	2824	136	3.19	3.40	8.0	115	5500	18	22	17450.0

🌟 Checking the Unique Car Company Names.

```
In [42]: df[["CompanyName"]].unique()
```

```
Out[42]: array(['alfa-romero', 'audi', 'bmw', 'chevrolet', 'dodge', 'honda',
       'isuzu', 'jaguar', 'mazda', 'mercury', 'buick', 'mercury',
       'mitsubishi', 'nissan', 'nissan', 'peugeot', 'plymouth', 'porsche',
       'porsche', 'renault', 'saab', 'subaru', 'toyota', 'toyota',
       'volksagen', 'volksagen', 'vw', 'volvo'], dtype=object)
```

```
In [43]: #Creating a Function to Replace the Values.
def replace(a,b):
    df["CompanyName"].replace(a,b,inplace=True)

replace('mazda','mazda')
replace('porsche','porsche')
replace('toyota','toyota')
replace('volksagen','volksagen')
replace('vw','volksagen')
df["CompanyName"].unique()
```

```
Out[43]: array(['alfa-romero', 'audi', 'bmw', 'chevrolet', 'dodge', 'honda',
       'isuzu', 'jaguar', 'mazda', 'buick', 'mercury', 'mitsubishi',
       'Nissan', 'Nissan', 'peugeot', 'plymouth', 'porsche', 'renault',
       'saab', 'subaru', 'toyota', 'volksagen', 'volvo'], dtype=object)
```

Observation


Now all the car company name seems correct. So we don't need to do any more cleaning. Now we can go to next step which is exploratory data analysis.

🌟 Exploratory data Analysis

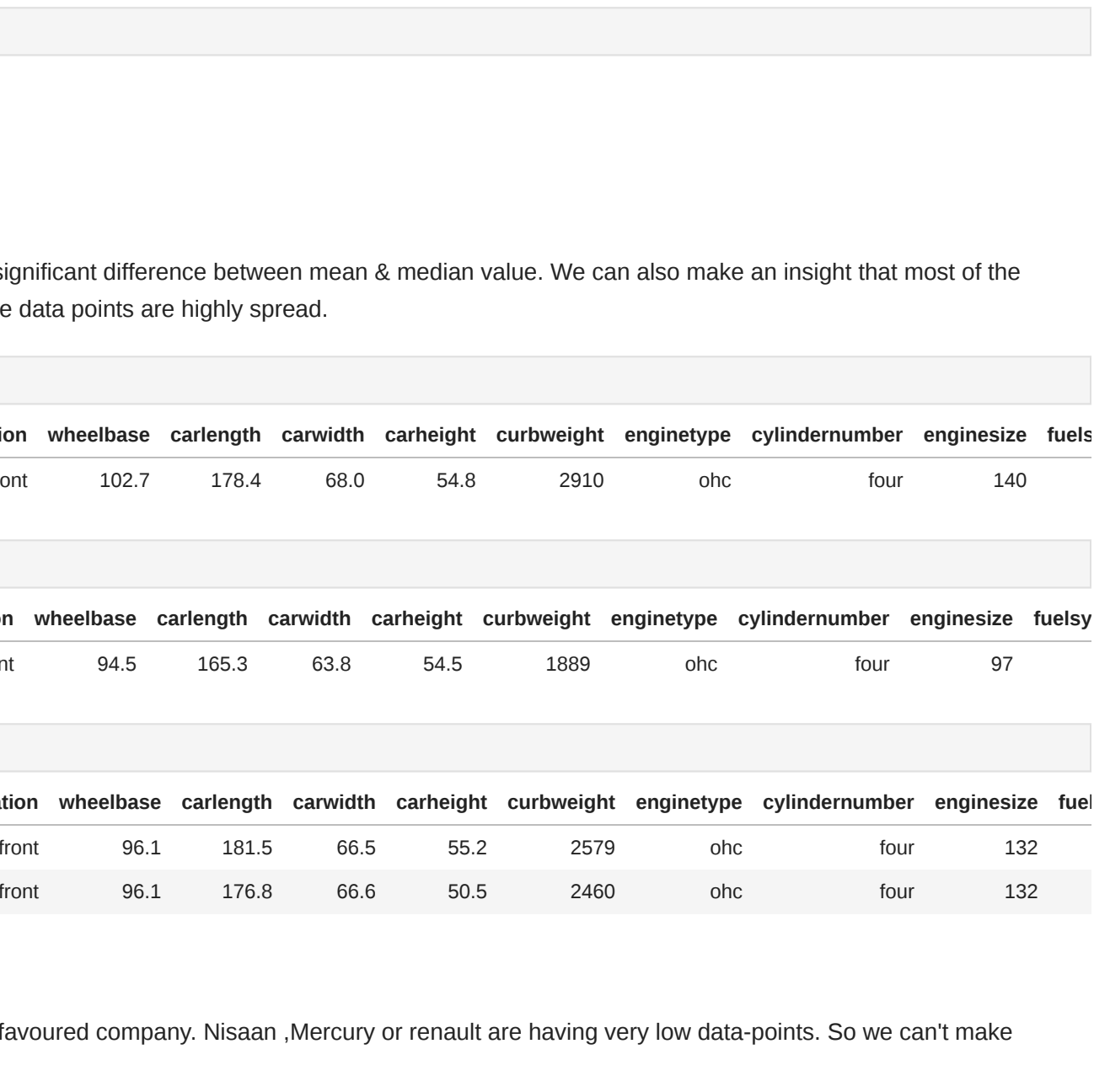
```
In [44]: #Visualizing our Target Feature.
plt.figure(figsize=(15,6))
plt.subplot(1,2,1)
sns.distplot(df["price"],color="red",kde=True)
plt.title("Car Price Distribution",fontweight="black",pad=20,fontsize=20)

plt.subplot(1,2,2)
sns.boxplot(y=df["price"],palette="Set2")
plt.title("Car Price Spread",fontweight="black",pad=20,fontsize=20)
plt.tight_layout()
plt.show()
```

Car Price Distribution



Car Price Spread



```
In [45]: df["price"].agg(["min","mean","median","max","std","skew"]).to_frame().T
```

	min	mean	median	max	std	skew
price	5118.0	13276.710571	10295.0	45400.0	7988.852332	1.777678

Insights

We can clearly observe that our Car Price Feature is Right Skewed. We can clearly observe that there is a significant difference between mean & median value. We can also make an insight that most of the car's price is below 14000. We can also see that the skewness of the car price is above 1.5 which means that the data points are highly spread.

```
In [49]: df[df["CompanyName"]=="mercury"]
```

	car_ID	symboling	CompanyName	fueltype	aspiration	doornumber	carbody	drivewheel	engineLocation	wheelbase	carlength	carwidth	carheight	curbweight	enginetype	cylindernumber	enginesize	fuelsystem
75	76	1	mercury	gas	turbo	two	hatchback	rwd	front	102.7	178.4	66.0	54.8	2910	ohc	four	140	

```
In [50]: df[df["CompanyName"]=="nissan"]
```

	car_ID	symboling	CompanyName	fueltype	aspiration	doornumber	carbody	drivewheel	engineLocation	wheelbase	carlength	carwidth	carheight	curbweight	enginetype	cylindernumber	enginesize	fuelsystem
89	90	1	Nissan	gas	std	two	sedan	fwd	front	94.5	165.3	63.8	54.5	1889	ohc	four	97	

```
In [51]: df[df["CompanyName"]=="renault"]
```

	car_ID	symboling	CompanyName	fueltype	aspiration	doornumber	carbody	drivewheel	engineLocation	wheelbase	carlength	carwidth	carheight	curbweight	enginetype	cylindernumber	enginesize	fuelsystem
130	131	0	renault	gas	std	four	wagon	fwd	front	96.1	181.5	66.5	55.2	2579	ohc	four	132	
131	132	2	renault	gas	std	two	hatchback	fwd	front	96.1	176.8	66.6	50.5	2460	ohc	four	132	

Insights

Toyota company has sold the highest number of cars. So we can say that Toyota is kind of customers most favoured company. Nissan, Mercury or renault are having very low data-points. So we can't make any inference of least sold car companies.

```
In [58]: #Visualizing Car Company w.r.t Price.
plt.figure(figsize=(15,6))
plt.subplot(1,2,1)
sns.boxplot(y=df["price"])
plt.xticks(rotation=90)
plt.title("Car Company vs Price",pad=10,fontweight="black",fontsize=20)
```

```
Out[58]: Text(0.5, 1.0, 'Car Company vs Price')
```

Car Company vs Price



```
In [61]: df[df["enginetype"]=="rotor"]
```

	car_ID	symboling	CompanyName	fueltype	aspiration	doornumber	carbody	drivewheel	engineLocation	wheelbase	carlength	carwidth	carheight	curbweight	enginetype	cylindernumber	enginesize	fuelsystem
56	57	3	mazda	gas	std	two	hatchback	rwd	front	95.3	169.0	65.7	49.6	2380	rotor	two	70	
56	57	3	mazda	gas	std	two	hatchback	rwd	front	95.3	169.0	66.7	49.6	2380	rotor	two	70	
57	58	3	mazda	gas	std	two	hatchback	rwd	front	95.3	169.0	65.7	49.6	2385	rotor	two	70	
58	59	3	mazda	gas	std	two	hatchback	rwd	front	95.3	169.0	65.7	49.6	2500	rotor	two	80	

Insights

Cars having Overhead Camshaft (OHC) engines are mostly sold. Only one car has been sold having engine type dohc. There are very few data-points of engine type dohcvt. So we can say that cars having overhead engines are mostly expensive. Cars having Overhead Camshaft (OHC) engines are least expensive cars.

🌟 Data Preprocessing

```
In [75]: #Deriving New Features From "Company Name" Feature.
#As we found an insight above that we can split the car company name into different price ranges. Like Low Range, Medium Range, High Range cars.
z = round(df.groupby(["CompanyName"])[["price"]].agg(["mean"])).T
z
```

	CompanyName	Nissan	Subaru	Alfa Romeo	Audi	Bmw	Isuzu	Jaguar	Dodge	Honda	Mercury	Mitsubishi	Nissan	Peugeot	Plymouth	Porsche	Renault	Saab	Subaru	Toyota		
	mean	5489.0	15498.33	17859.17	26118.75	33647.0	6007.0	7875.44	8184.69	8916.5	34600.0	10652.88	16503.0	9239.77	10704.88	15489.09	7963.43	31400.5	9595.0	15223.33	8541.25	8885

```
In [76]: df = df.merge(z.T,how="left",on="CompanyName")
bins = [0,10000,20000,40000]
cars_bin["Budget","Medium","Highend"]
df["CarsRange"] = pd.cut(df["mean"],bins,right=False,labels=cars_bin)
df.head()
```

	car_ID	symboling	CompanyName	fueltype	aspiration	doornumber	carbody	drivewheel	engineLocation	wheelbase	carlength	carwidth	carheight	curbweight	enginetype	cylindernumber	enginesize	fuelsystem
0	1	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8	64.1	48.8	2548	dohc	four	130	
1	2	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8	64.1	48.8	2548	dohc	four	130	
2	3	1	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	171.2	65.5	52.4	2823	ohcv	six	152	
3	4	2	audi	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	54.3	2337	ohc	four	109	
4	5	2	audi	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	54.3	2824	ohc	five	136	

```
In [77]: #Creating new DataFrame with all the useful Features.
new_df = df[["fueltype","aspiration","doornumber","carbody","drivewheel","enginetype","cylindernumber","fuelsystem","wheelbase","carlength","carwidth","curbweight","enginesize","boreratio","horsepower","citympg","highwaympg","price","CarsRange"]]
new_df.head()
```

	fueltype	aspiration	doornumber	carbody	drivewheel	enginetype	cylindernumber	fuelsystem	wheelbase	carlength	carwidth	curbweight	enginesize	boreratio	horsepower	citympg	highwaympg	price
0	gas	std	two	convertible	rwd	dohc	four	mpfi	88.6	168.8	64.1	2548	130	3.47	111	21	27	13495.0
1	gas	std	two	convertible	rwd	dohc	four	mpfi	88.6	168.8	64.1	2548	130	3.47	111	21	27	16500.0
2	gas	std	two	hatchback	rwd	ohcv	six	mpfi	94.5	171.2	65.5	2823	152	2.68	154	19	26	16500.0