

# Classification Problem

Title: To predict value from the dataset that is categorical in nature.

## Multi-class Classification:

Multi-class Classification refers to predicting more than two classes/columns.  
To predict Type of Cars from the dataset.

## Import Libraries

```
In [1]: from warnings import filterwarnings
filterwarnings("ignore")
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
```

## Collecting Data

Reading in the Cars93.csv file into a pandas dataframe.

```
In [2]: A = pd.read_csv("C:/Users/Aksha/OneDrive/Desktop/Cars93.csv")
```

```
In [3]: A.head()
```

```
Out[3]:
```

	id	Manufacturer	Model	Type	MinPrice	Price	MaxPrice	MPG.city	MPG.highway	AirBags	Passengers	Length	Wheelbase	Width	Turncicle	Rearsearoom	Luggage room	Weight	Origin	Make
0	1	Acura	Integra	Small	12.9	15.9	18.8	25	31	None	5	177	102	68	37	26.5	11.0	2705	non-USA	Acura Integra
1	2	Acura	Legend	Midsize	29.2	33.9	38.7	18	25	Driver & Passenger	5	195	115	71	38	30.0	15.0	3560	non-USA	Acura Legend
2	3	Audi	90	Compact	25.9	29.1	32.3	20	26	Driver only	5	180	102	67	37	28.0	14.0	3375	non-USA	Audi 90
3	4	Audi	100	Midsize	30.8	37.7	44.6	19	26	None	6	193	106	70	37	21.0	17.0	3405	non-USA	Audi 100
4	5	BMW	116	Midsize	23.7	30.0	36.2	22	30	Driver only	4	186	109	69	39	27.0	13.0	3640	non-USA	BMW 116

5 rows x 28 columns

```
In [4]: Q = []
import re
for i in A.columns:
    Q.append(re.sub("[^a-zA-z]", "", i))
A.columns = Q
```

```
In [5]: Q
```

```
Out[5]: ['id',
'Manufacturer',
'model',
'type',
'minprice',
'price',
'maxprice',
'mpgcity',
'mpghighway',
'airbags',
'passengers',
'length',
'wheelbase',
'width',
'turncicle',
'rearsearoom',
'luggage room',
'weight',
'origin',
'make']
```

## Missing Data Treatment

Cleaning the data by removing NaN values and unnecessary columns from the dataset.

```
In [11]: A.isnull() #True == if it is null
#False == if it is not null
```

```
Out[11]:
```

	id	Manufacturer	Model	Type	MinPrice	Price	MaxPrice	MPG.city	MPG.highway	AirBag	Passengers	Length	Wheelbase	Width	Turncicle	Rearsearoom	Luggage room	Weight	Origin	Make
0	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	True	...	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
88	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	True	False	False	False
89	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False
90	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False
91	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False
92	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False

93 rows x 28 columns

```
In [17]: A.isnull().sum() #airbags,rear_seat,room,luggage room have missing values in them.
```

```
Out[17]:
```

id	0
Manufacturer	0
Model	0
Type	0
MinPrice	0
Price	0
MaxPrice	0
MPG.city	0
MPG.highway	0
AirBags	4
DriveTrain	0
Cylinders	0
Cylinders	0
EngineSize	0
EnginePower	0
HP	0
RevsPerMin	0
Mantransavall	0
Footcandcapacity	0
Passengers	0
Length	0
Wheelbase	0
Width	0
Turncicle	0
Rearsearoom	0
Luggage room	11
Weight	0
Origin	0
Make	0
dtype: int64	

```
In [19]: #To analyze them visually, use seaborn to create a simple heatmap to see where is the data missing.
sb.heatmap(A.isnull(), yticklabels=False, cbar=False)
```

```
Out[19]: <AxesSubplot>
```

```
In [19]: #To remove these columns
from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = "all"
A = A.dropna(inplace=True)
```

```
In [100]: A.isna().sum()
```

```
Out[100]:
```

id	0
Manufacturer	0
Model	0
Type	0
MinPrice	0
Price	0
MaxPrice	0
MPG.city	0
MPG.highway	0
AirBags	0
DriveTrain	0
Cylinders	0
Cylinders	0
EngineSize	0
EnginePower	0
HP	0
RevsPerMin	0
Mantransavall	0
Footcandcapacity	0
Passengers	0
Length	0
Wheelbase	0
Width	0
Turncicle	0
Rearsearoom	0
Luggage room	0
Weight	0
Origin	0
Make	0
dtype: int64	

```
In [111]: #to check if the values have been removed or not.
sb.heatmap(A.isnull(), yticklabels=False, cbar=False)
```

```
Out[111]: <AxesSubplot>
```

## Dropping Unnecessary Cols

```
In [112]: X = A[['Type','id','Make','Model','Manufacturer'],axis=1]
Y = A.drop(labels=['Type','id','Make','Model','Manufacturer'],axis=1)
```

## Categorical & Continuous

```
In [113]: cat = []
con = []
for i in X.columns:
    if X[i].dtype == 'object':
        cat.append(i)
    else:
        con.append(i)
```

```
In [114]: cat
```

```
Out[114]: ['AirBags', 'DriveTrain', 'Cylinders', 'Mantransavall', 'Origin']
```

```
In [115]: con
```

```
Out[115]: ['MinPrice',
'Price',
'MaxPrice',
'MPG.city',
'MPG.highway',
'EngineSize',
'EnginePower',
'HP',
'RevsPerMin',
'FuelTankCapacity',
'Passengers',
'Length',
'Wheelbase',
'Width',
'Turncicle',
'Rearsearoom',
'Luggage room',
'Weight',
'Origin',
'Make']
```

## Analyzing data through Boxplot

```
In [148]: from IPython import InteractiveShell
from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = "all"
A = A.dropna(inplace=True)
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

```
Out[148]:
```

## Cross Tabulation

To describe the relationship between two categorical variables, we use a special type of table called cross-tabulation

```
In [171]: from IPython import InteractiveShell
from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = "all"
A = A.dropna(inplace=True)
```

```
Out[171]: ['AirBags', 'DriveTrain', 'Cylinders', 'Mantransavall', 'Origin']
```

```
In [182]: pd.crosstab([Type,A,AirBags])
```

```
Out[182]:
```

AirBags	Driver & Passenger	AirBags:Driver only	AirBags:None	Cylinders_3	Cylinders_4	Cylinders_5	Cylinders_6	Cylinders:rotary	Mantransavall	Yes	Mantransavall	Yes	Origin_USA	Origin	non-USA
Compact	2	2	9	5											
Midsize	6	12	4												
Small	0	5	16												
Sporty	3	9	2												
Van	0	3	6												

```
In [200]: pd.crosstab([Type,A,DriveTrain])
```

```
Out[200]:
```

```
In [211]: pd.crosstab([Type,A,Cylinders])
```

```
Out[211]:
```

```
In [222]: pd.crosstab([Type,A,Mantransavall])
```

```
Out[222]:
```

```
In [239]: pd.crosstab([Type,A,Origin])
```

```
Out[239]:
```

```
In [246]: #convert to dummy values
cat = pd.get_dummies(cat)
Xcat = pd.get_dummies(X[cat])
```

```
In [257]: Xcat
```

```
Out[257]:
```

AirBags	Driver & Passenger	AirBags:Driver only	AirBags:None	Cylinders_3	Cylinders_4	Cylinders_5	Cylinders_6	Cylinders:rotary	Mantransavall	No	Mantransavall	Yes	Origin_USA	Origin	non-USA
0	1	0	0	1	0	0	0	0	0	0	0	0	1	0	1
1	1	0	0	0	0	0	0	1	0	0	0	0	1	0	1
2	1	0	0	0	0	0	0	0	0	0	0	0	1	0	1
3	1	0	0	0	0	0	0	1	0	0	0	0	1	0	1
4	0	1	0	0	0	1	0	0	0	0	0	0	1	0	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
88	0	0	0	1	0	0	1	0	0	0	0	0	1	0	1
89	0	0	1	0	0	1	0	0	0	0	0	0	1	0	1
90	0	0	1	0	0	0	1	0	0	0	0	0	1	0	1
91	0	1	0	0	0	1	0	0	0	0	0	0	1	0	1
92	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1

93 rows x 13 columns

## Standardization of Continuous Variables

```
In [268]: from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X = pd.DataFrame(Xcat.values.tolist(), columns=Xcat.columns)
X = sc.fit_transform(X)
```

```
In [269]: X.head()
```

```
Out[269]:
```

	MinPrice	Price	MaxPrice	MPG.city	MPG.highway	EngineSize	EnginePower	HP	RevsPerMin	FuelTankCapacity	...	Cylinders_3	Cylinders_4	Cylinders_5	Cylinders_6	Cylinders:rotary	Mantransavall	Yes	Origin_USA	Origin	non-USA
0	14805787	1697270	1932405	0.71312	0.36925	-0.84102	-0.073484	1.317489	1.129310	-1.002184	...	0	1	0	0	0	0	0	1	0	1
1	1388017	1692844	1531409	-0.781032	-0.770514	0.515869	1.078222	0.369586	0.005661	0.409445	...	0	0	0	1	0	0	0	1	0	1
2	1008558	0.986227	0.948052	0.423219	-0.581941	0.128186	0.540813	0.369586	-0.107113	0.072197	...	0	0	0	1	0	0	0	1	0	1
3	1571949	1.693374	2.693191	0.602126	-0.581941	0.128186	0.540813	0.369586	0.410559	1.359872	...	0	0	0	1	0	0	0	1	0	1
4	0.755532	1.091905	1.393535	-0.056407	0.172352	0.808601	1.231897	0.705562	0.430909	1.359872	...	0	1	0	0	0	0	0	1	0	1

5 rows x 31 columns

## Building a Logistic Regression Model

### Train Test Split

Splitting the data into training set and testing set.

```
In [269]: from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(X,y,split_size=0.3,random_state=1)
```

## Logistic Regression

Logistic Regression is used as it solves classification problems & is used to predict categorical variable.

```
In [281]: from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()
model = lr.fit(X_train,y_train)
pred = model.predict(X_test)
```

## Accuracy Score

```
In [302]: #confusion matrix
from sklearn.metrics import confusion_matrix
confusion_matrix(y_test,pred)
```

```
Out[302]:
```

array([[1, 0, 0, 0, 0, 0],					
[0, 1, 0, 0, 0, 0],					
[0, 0, 1, 0, 0, 0],					
[0, 0, 0, 1, 0, 0],					
[0, 0, 0, 0, 1, 0],					