

In [1]:

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib as mpl
5 %matplotlib inline
6 # mpl.style.use('ggplot')

```

In [2]:

```
1 car = pd.read_csv('quikr_car.csv')
```

In [3]:

```
1 car.head()
```

Out[3]:

	name	company	year	Price	kms_driven	fuel_type
0	Hyundai Santro Xing XO eRLX Euro III	Hyundai	2007	80,000	45,000 kms	Petrol
1	Mahindra Jeep CL550 MDI	Mahindra	2006	4,25,000	40 kms	Diesel
2	Maruti Suzuki Alto 800 Vxi	Maruti	2018	Ask For Price	22,000 kms	Petrol
3	Hyundai Grand i10 Magna 1.2 Kappa VTVT	Hyundai	2014	3,25,000	28,000 kms	Petrol
4	Ford EcoSport Titanium 1.5L TDCi	Ford	2014	5,75,000	36,000 kms	Diesel

In [4]:

```
1 car.shape
```

Out[4]:

(892, 6)

In [5]:

```
1 car.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 892 entries, 0 to 891
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   name            892 non-null   object
1   company         892 non-null   object
2   year            892 non-null   object
3   Price           892 non-null   object
4   kms_driven      840 non-null   object
5   fuel_type       837 non-null   object
dtypes: object(6)
memory usage: 41.9+ KB

```

In [6]:

```
1 backup = car.copy()
```

## Data Quality

- names are pretty inconsistent
- names have company names attached to it
- some names are spam like 'Maruti Ertiga showroom condition with' and 'Well mentained Tata Sumo'
- company: many of the names are not of any company like 'Used', 'URJENT', and so on.
- year has many non-year values
- year is in object. Change to integer
- Price has Ask for Price
- Price has commas in its prices and is in object
- kms\_driven has object values with kms at last.
- It has nan values and two rows have 'Petrol' in them
- fuel\_type has nan values

## Cleaning Data

**Year has many non-year valus**

In [7]:

```
1 car = car[car['year'].str.isnumeric()]
```

**Year is in object. change to intger**

In [8]:

```
1 car['year'] = car['year'].astype(int)
```

**Price has "Ask for Price"**

In [9]:

```
1 car = car[car['Price'] != 'Ask For Price']
```

In [10]:

```
1 # car['Price'].unique()
```

**Price has commas in its prices and is in object**

In [11]:

```
1 car['Price']=car['Price'].str.replace(',','').astype(int)
```

**kms\_driven has object values with kms at last**

In [12]:

```
1 car['kms_driven']=car['kms_driven'].str.split().str.get(0).str.replace(',','')
```

**It has nan value values and two rows have 'Petrol' in them**

In [13]:

```
1 car=car[car['kms_driven'].str.isnumeric()]
```

In [14]:

```
1 car['kms_driven']=car['kms_driven'].astype(int)
```

**fuel\_type has nan values**

In [15]:

```
1 car=car[~car['fuel_type'].isna()]
```

In [16]:

```
1 car.shape
```

Out[16]:

(816, 6)

**name and company had spammed data...but with previous cleaning, those rows got removed.**

**company does not need any cleaning now. WChwanging car names. keeping only the first 3 words**

In [17]:

```
1 car['name']=car['name'].str.split().str.slice(start=0, stop=3).str.join(' ')
```

**Resetting the index of the final cleaned data**

In [18]:

```
1 car=car.reset_index(drop=True)
```

## Cleaned Data

In [19]:

```
1 car.to_csv('Cleaned_Car_data.csv')
```

In [20]:

```
1 car.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 816 entries, 0 to 815
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   name            816 non-null    object
1   company         816 non-null    object
2   year            816 non-null    int32
3   Price           816 non-null    int32
4   kms_driven      816 non-null    int32
5   fuel_type       816 non-null    object
dtypes: int32(3), object(3)
memory usage: 28.8+ KB
```

In [21]:

```
1 car.describe(include='all')
```

Out[21]:

	name	company	year	Price	kms_driven	fuel_type
<b>count</b>	816	816	816.000000	8.160000e+02	816.000000	816
<b>unique</b>	254	25	NaN	NaN	NaN	3
<b>top</b>	Maruti Suzuki Swift	Maruti	NaN	NaN	NaN	Petrol
<b>freq</b>	51	221	NaN	NaN	NaN	428
<b>mean</b>	NaN	NaN	2012.444853	4.117176e+05	46275.531863	NaN
<b>std</b>	NaN	NaN	4.002992	4.751844e+05	34297.428044	NaN
<b>min</b>	NaN	NaN	1995.000000	3.000000e+04	0.000000	NaN
<b>25%</b>	NaN	NaN	2010.000000	1.750000e+05	27000.000000	NaN
<b>50%</b>	NaN	NaN	2013.000000	2.999990e+05	41000.000000	NaN
<b>75%</b>	NaN	NaN	2015.000000	4.912500e+05	56818.500000	NaN
<b>max</b>	NaN	NaN	2019.000000	8.500003e+06	400000.000000	NaN

In [22]:

```
1 car=car[car['Price']<6000000]
```

## Extracting Data

In [23]:

```
1 x=car.drop(columns='Price')
2 y=car['Price']
```

In [24]:

```
1 from sklearn.model_selection import train_test_split
2 x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.2)
```

In [25]:

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score
3 from sklearn.preprocessing import OneHotEncoder
4 from sklearn.compose import make_column_transformer
5 from sklearn.pipeline import make_pipeline
```

In [26]:

```
1 ohe = OneHotEncoder()
2 ohe.fit(x[['name', 'company', 'fuel_type']])
```

Out[26]:

OneHotEncoder()

In [27]:

```
1 column_trans = make_column_transformer((OneHotEncoder(categories=ohe.categories_), ['name', 'company', 'fuel_type']),
2                                         remainder='passthrough')
```

In [28]:

```
1 lr = LinearRegression()
```

In [29]:

```
1 pipe = make_pipeline(column_trans, lr)
```

In [30]:

```
1 pipe.fit(x_train, y_train)
```

Out[30]:

```
Pipeline(steps=[('columntransformer',
                  ColumnTransformer(remainder='passthrough',
                                     transformers=[('onehotencoder',
                                                    OneHotEncoder(categories=
[array(['Audi A3 Cabriolet', 'Audi A4 1.8', 'Audi A4 2.0', 'Audi A6 2.0',
'Audi A8', 'Audi Q3 2.0', 'Audi Q5 2.0', 'Audi Q7', 'BMW 3 Series',
'BMW 5 Series', 'BMW 7 Series', 'BMW X1', 'BMW X1 sDrive20d',
'BMW X1 xDrive20d', 'Chevrolet Beat', 'Chevrolet Beat...

array(['Audi', 'BMW', 'Chevrolet', 'Datsun', 'Fiat', 'Force', 'Ford',
'Hindustan', 'Honda', 'Hyundai', 'Jaguar', 'Jeep', 'Land',
'Mahindra', 'Maruti', 'Mercedes', 'Mini', 'Mitsubishi', 'Nissan',
'Renault', 'Skoda', 'Tata', 'Toyota', 'Volkswagen', 'Volvo'],
dtype=object),

array(['Diesel', 'LPG', 'Petrol'], dtype=object)]),
                  ['name', 'company',
                  'fuel_type'])))],
        ('linearregression', LinearRegression()))]
```

In [31]:

```
1 y_pred = pipe.predict(x_test)
```

In [32]:

```
1 r2_score(y_test, y_pred)
```

Out[32]:

```
0.6954629091832226
```

**Finding the model with a random state of TrainTestSplit where the model was found to give almost 0.88**

In [33]:

```
1 scores=[]
2 for i in range(1000):
3     x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.2,random_state=i)
4     lr=LinearRegression()
5     pipe=make_pipeline(column_trans, lr)
6     pipe.fit(x_train, y_train)
7     y_pred=pipe.predict(x_test)
8     # print(r2_score(y_test, y_pred), i)
9     scores.append(r2_score(y_test, y_pred))
```

In [34]:

```
1 np.argmax(scores)
```

Out[34]:

661

In [35]:

```
1 scores[np.argmax(scores)]
```

Out[35]:

0.889770931767991

In [36]:

```
1 x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.2,random_state=np.  
2 lr=LinearRegression()  
3 pipe=make_pipeline(column_trans, lr)  
4 pipe.fit(x_train, y_train)  
5 y_pred=pipe.predict(x_test)  
6 r2_score(y_test, y_pred)
```

Out[36]:

0.889770931767991

In [37]:

```
1 import pickle
```

In [38]:

```
1 pickle.dump(pipe,open('LinearRegressionModel.pkl','wb'))
```

In [39]:

```
1 pipe.predict(pd.DataFrame([[ 'Maruti Suzuki Swift', 'Maruti', 2019,100,'Petrol']], columnr
```

Out[39]:

array([400795.54550476])

In [40]:

```
1 pipe.predict(pd.DataFrame([[ 'Maruti Suzuki Swift', 'Maruti', 2017,100,'Petrol']], columnr
```

Out[40]:

array([363365.76708183])

In [41]:

```
1 pipe.predict(pd.DataFrame([[ 'Maruti Suzuki Swift', 'Maruti', 2015,1000,'Petrol']], columns=
```

Out[41]:

```
array([325857.22134978])
```

In [42]:

```
1 pipe.predict(pd.DataFrame([[ 'Maruti Suzuki Swift', 'Maruti', 2015,500000,'Petrol']], columns=
```

Out[42]:

```
array([282185.12441023])
```

In [ ]:

```
1
```