

```
In [6]: #import libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

#ignore/ disable warnings
import warnings
warnings.filterwarnings("ignore")

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

```
In [7]: #import data
df=pd.read_csv(r"C:\Users\Sanjay Lohar\Downloads\House_Prices.csv")
```

```
In [8]: df.head()
```

Out[8]:	Home	Price	SqFt	Bedrooms	Bathrooms	Offers	Brick	Neighborhood
0	1	114300	1780.0	2	2	2	No	East
1	2	114200	2030.0	4	2	3	No	East
2	3	114800	1740.0	3	2	1	No	East
3	4	94700	1980.0	3	2	3	No	East
4	5	119800	2130.0	3	3	3	No	East

```
In [9]: #check shape, dtypes
df.shape #128 observations, 8 variables
```

```
Out[9]: (128, 8)
```

```
In [10]: df.dtypes #all variables are stored with appropriate data types
```

Out[10]:	Home	Price	SqFt	Bedrooms	Bathrooms	Offers	Brick	Neighborhood
	int64	int64	float64	int64	int64	int64	object	object

```
In [11]: # check missing values
df.isnull().sum() #there are 2 missing values in the data
```

Out[11]:	Home	Price	SqFt	Bedrooms	Bathrooms	Offers	Brick	Neighborhood
	0	0	1	0	0	0	0	1
								dtype: int64

```
In [12]: #missing value treatment (Remove records having missing values)
df=df.dropna(axis=0)
```

```
In [13]: #check whther missing value records are removed or not
df.isnull().sum()
```

Out[13]:	Home	Price	SqFt	Bedrooms	Bathrooms	Offers	Brick	Neighborhood
	0	0	0	0	0	0	0	0
								dtype: int64

```
In [14]: #Check outliers
plt.boxplot(df['Price']) #has outliers
plt.show()
```



```
In [15]: plt.boxplot(df['SqFt']) #has outliers
plt.show()
```



```
In [16]: #outlier treatment
#we will write a user defined function to remove outliers
```

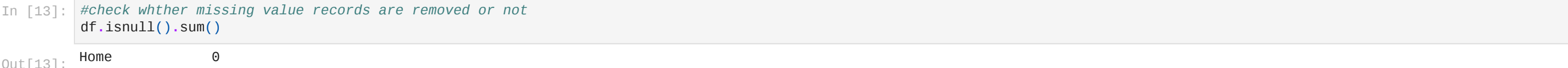
```
In [17]: def remove_outlier(d,c):
#find q1 and q3
q1=d[c].quantile(0.25)
q3=d[c].quantile(0.75)

#find iqr (inter quartile range)
iqr=q3-q1

#find upper and lower bound
ub=q3+1.5*iqr
lb=q1-1.5*iqr

final_data = d[(d[c]>lb) & (d[c]<ub)]
return final_data
```

```
In [26]: #Remove outliers from price
df=remove_outlier(df,"Price")
plt.boxplot(df["Price"])
plt.show()
```



```
In [33]: #Remove outliers from SqFt
df=remove_outlier(df,"SqFt")
plt.boxplot(df["SqFt"])
plt.show()
```

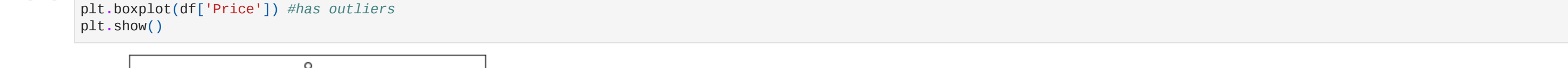


```
In [34]: #-----Start EDA (Exploratory Data Analysis)-----
```

```
In [35]: #Price distribution
sns.distplot(df["Price"])
plt.show()
```



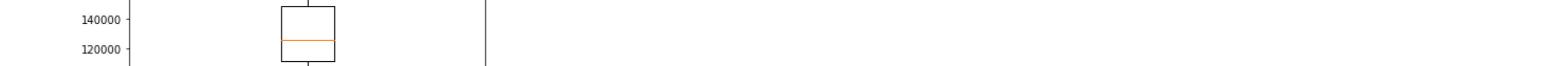
```
In [36]: #SqFt Distribution
sns.distplot(df["SqFt"])
plt.show()
```



```
In [37]: #Correlation between Price and Sqft
df.plot(kind='scatter', x='SqFt', y='Price')
plt.show()
```



```
In [66]: #count of house by bedrooms
df.groupby('Bedrooms')['Bedrooms'].count().plot(kind='bar')
plt.show()
```



```
In [39]: #count of house by neighborhood
df.groupby('Neighborhood')['Neighborhood'].count().plot(kind='bar')
plt.show()
```



```
In [40]: #Check the correlation of numeric variables
df_numeric = df.select_dtypes(include=['float64', 'int64'])
df_numeric.head()
```

Out[40]:	Home	Price	SqFt	Bedrooms	Bathrooms	Offers
0	1	114300	1780.0	2	2	2
1	2	114200	2030.0	4	2	3
2	3	114800	1740.0	3	2	1
3	4	94700	1980.0	3	2	3
4	5	119800	2130.0	3	3	3

```
In [41]: # drop Home variable because it is just a serial number
df_numeric = df_numeric.drop('Home', axis=1)
df_numeric.head()
```

Out[41]:	Price	SqFt	Bedrooms	Bathrooms	Offers
0	114300	1780.0	2	2	2
1	114200	2030.0	4	2	3
2	114800	1740.0	3	2	1
3	94700	1980.0	3	2	3
4	119800	2130.0	3	3	3

```
In [42]: # correlation matrix
cor_mat = df_numeric.corr()
cor_mat
```

Out[42]:		Price	SqFt	Bedrooms	Bathrooms	Offers
	Price	1.000000	0.512462	0.499146	0.501489	-0.386868
	SqFt	0.512462	1.000000	0.436659	0.490578	0.293352
	Bedrooms	0.499146	0.436659	1.000000	0.382643	0.064872
	Bathrooms	0.501489	0.490578	0.382643	1.000000	0.114522
	Offers	-0.386868	0.293352	0.064872	0.114522	1.000000

```
In [43]: # plot correlations on a heatmap
```

```
# figure size
plt.figure(figsize=(10,5))
```

```
# heatmap
sns.heatmap(cor_mat, cmap="YlGnBu", annot=True) #YlGnBu
plt.show()
```



```
In [44]: #----- END OF EDA -----
```

```
In [45]: #print column names
df.columns
```

```
Out[45]: Index(['Home', 'Price', 'SqFt', 'Bedrooms', 'Bathrooms', 'Offers', 'Brick', 'Neighborhood'], dtype='object')
```

```
In [46]: df.dtypes
```

Out[46]:	Home	Price	SqFt	Bedrooms	Bathrooms	Offers	Brick	Neighborhood
	int64	int64	float64	int64	int64	int64	object	object

```
In [47]: #transform Bedrooms and bathrooms to object
df['Bedrooms']=df['Bedrooms'].replace([1,2,3,4,5],[1'1BR',2'2BR',3'3BR',4'4BR',5'5BR'])
df['Bathrooms']=df['Bathrooms'].replace([2,3,4],[2'Bath',3'3Bath',4'4Bath'])
```

```
In [48]: df.head()
```

Out[48]:	Home	Price	SqFt	Bedrooms	Bathrooms	Offers	Brick	Neighborhood
0	1	114300	1780.0	2BR	2Bath	2	No	East
1	2	114200	2030.0	4BR	2Bath	3	No	East
2	3	114800	1740.0	3BR	2Bath	1	No	East
3	4	94700	1980.0	3BR	2Bath	3	No	East
4	5	119800	2130.0	3BR	3Bath	3	No	East

```
In [49]: df.Bathrooms.unique()
```

```
Out[49]: array(['2Bath', '3Bath', '4Bath'], dtype=object)
```

```
In [50]: # create dummy variables for categorical variables
```

```
# subset all categorical variables
df_categorical = df.select_dtypes(include=['object'])
df_categorical.head()
```

Out[50]:	Bedrooms	Bathrooms	Brick	Neighborhood
0	2BR	2Bath	No	East
1	4BR	2Bath	No	East
2	3BR	2Bath	No	East
3	3BR	2Bath	No	East
4	3BR	3Bath	No	East

```
In [51]: # convert into dummies
df_dummies = pd.get_dummies(df_categorical)
df_dummies
```

Out[51]:	Bedrooms_2BR	Bedrooms_3BR	Bedrooms_4BR	Bedrooms_5BR	Bathrooms_2Bath	Bathrooms_3Bath	Bathrooms_4Bath	Brick_No	Brick_Yes	Neighborhood_East	Neighborhood_North	Neighborhood_West
0	1	0	0	0	1	0	0	1	0	1	0	0
1	0	0	1	0	1	0	0	1	0	1	0	0
2	0	1	0	0	1	0	0	1	0	1	0	0
3	0	1	0	0	1	0	0	1	0	1	0	0
4	0	1	0	0	0	1	0	1	0	1	0	0

```
122 rows x 12 columns
```

```
In [52]: #Create master data by combining only numeric columns and dummies
df_numeric = df_numeric.drop(['Bedrooms', 'Bathrooms', 'Offers'], axis=1)
df_numeric.head()
```

Out[52]:	Price	SqFt
0	114300	1780.0
1	114200	2030.0
2	114800	1740.0
3	94700	1980.0
4	119800	2130.0

```
In [53]: #Combine data
master_df=pd.concat([df_numeric,df_dummies], axis=1)
```

```
In [54]: #master master df to excel to validate the data
#master_df.to_excel(r"C:\Users\Sanjay Lohar\Desktop\prepared_data.xlsx")
```

```
In [55]: #Create x and y
y=master_df['Price']
x=master_df.drop('Price', axis=1)
```

```
In [56]: #Split x and y into training and test samples
from sklearn.model_selection import train_test_split
```

```
In [57]: xtrain, xtest, ytrain,ytest= train_test_split(x,y, test_size=0.3,random_state=0)
```

```
In [58]: print(xtrain.shape,ytrain.shape, xtest.shape,ytest.shape)
(85, 13) (85,) (37, 13) (37,)
```

```
In [ ]:
```

## Linear regression

```
In [59]: #Build Linear Regression model
from sklearn.linear_model import LinearRegression
```

```
#create an instance of Linear Regression
model=LinearRegression()
```

```
#fit the model using training sample
model.fit(xtrain,ytrain)
```

```
Out[59]: LinearRegression()
```

```
In [60]: #check accuracy of training model
model.score(xtrain,ytrain)
```

```
Out[60]: 0.836396621740817
```

```
In [61]: #predict house price
y_pred=model.predict(xtest)
```

```
In [62]: #check accuracy of test model
model.score(xtest,ytest)
```

```
Out[62]: 0.7837367468033985
```

```
In [63]: # print coefficients and intercept
print(model.coef_) #B1, B2, B3,...,Bn
print(model.intercept_) #B0
```

[	38.27609761	-536.97256168	-822.6711758	8944.00851895
-	10869.78389424	10669.78389424	-8228.66613913	-10382.75582474
-	18611.4239637	67587.39272921886		

```
In [64]: y_pred
```

Out[64]:	array([157891.33865404, 110137.29967318, 134396.47711957, 115835.85436085, 128885.68889607, 112953.66040863, 124204.48893494, 154639.56578979, 159332.31572027, 117366.89826542, 110422.99828729, 124454.24311827, 110849.23731115, 115115.19236384, 161966.54108292, 145916.83681149, 129998.01047815, 133017.46276467, 121719.74717424, 128760.47828663, 141845.2843425, 114738.43138689, 95112.55924156, 111767.96179452, 105546.55817601, 108692.45888789, 130231.7052213, 182966.46038562, 106212.62754971, 106425.79789835, 160897.83767256, 111145.6682851, 148889.12288214, 166187.9045109, 118177.83602692, 111143.10443836, 109471.39593766])
----------	---

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```