```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from statsmodels.api import add_constant,OLS
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler,MinMaxScaler,RobustScaler
from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.stats.diagnostic import het_breuschpagan
from statsmodels.graphics.gofplots import qqplot
from sklearn.linear_model import LinearRegression,Ridge,Lasso,ElasticNet,SGDRegressor
from mlxtend.feature_selection import SequentialFeatureSelector as sfs
from sklearn.feature_selection import SequentialFeatureSelector as sfs,RFE
from sklearn.metrics import r2_score, mean_absolute_error,mean_squared_error,mean_abso
from sklearn.model_selection import train_test_split, GridSearchCV,RandomizedSearchCV
import warnings
warnings.filterwarnings('ignore')
from statsmodels.stats.diagnostic import linear_rainbow,het_breuschpagan
import statsmodels.api as sma
```

```python
! pip install mlxtend
```

# Case Study

# Data Exploration

```python
df=pd.read_csv(r"C:\Users\M.JOTHI\OneDrive\Desktop\bridge course\ML\week 2\casestudy\B
df.head()
```

```python
print('No. of rows:',df.shape[0])
print('No. of columns:',df.shape[1])
```

```python
df.dtypes
```

```python
df.info()
```

```python
df.describe().T
```

```python
df.describe(include='object').T
```

```python
# Lets treat the total_sqft first
df['total_sqft'].unique()
```

```python
pd.set_option('display.max_rows',None)
df['total_sqft'].value_counts()
```

```python
df['total_sqft']=df['total_sqft'].str.split(expand=True).loc[:,0]
```

```python
df['total_sqft'].value_counts()
```

```python
df['total_sqft'].str.isdigit().sum() # for numeric
```

```python
df=df[df['total_sqft'].str.isnumeric()] # for numeric
```

```python
df.shape
```

```python
df['total_sqft']=df['total_sqft'].astype(float)
```

```python
# lets treat the size column as well
df['size'].value_counts()
```

```python
df['size']=df['size'].str.split(expand=True).iloc[:,0]
```

```python
df['size']=df['size'].astype(float)
```

```python
df.dtypes
```

```python
# Null values
df['society'].isnull().sum()/df.shape[0]
```

- There is 41 % null values are there so we drop the null values.

```python
df.drop(columns='society',inplace=True)
# del df['society']
```

```python
df.head()
```

```python
# lets now treat the column availability
df['availability'].unique()
```

```python
def availability(s):
    if s in ['Ready To Move','Immediate Possession']:
        return 'Ready To Move'
    else:
        return 'Under Construction'
```

```python
df['availability']=df['availability'].apply(availability)
```

```python
df['availability'].value_counts()
```

```python
df.head()
```

```python
# lets check the area type as well
df['area_type'].value_counts()
```

# Univarite Analysis

```python
df.columns
```

```python
num_col=['size', 'total_sqft', 'bath',
        'balcony', 'price']
cat_col=['area_type','availability','location']
```

```python
t=1
plt.figure(figsize=(10,10))
for i in num_col:
    plt.subplot(3,3,t)
    sns.distplot(x=df[i])
    t+=1
plt.tight_layout()
plt.show()
```

```python
for i in df.select_dtypes(include='number'):
    sns.boxplot(x=df[i])
    plt.show()
```

```python
t=1
plt.figure(figsize=(10,10))
for i in cat_col:
    plt.subplot(2,2,t)
    sns.countplot(x=df[i])
    plt.xticks(rotation=90)
    t+=1
plt.tight_layout()
plt.show()
```

```python
df['location']=df['location'].str.lstrip(' ')
df.groupby('location')['price'].mean().sort_values(ascending=False
                                                   )
```

```python
top_10_location= df.groupby('location')['price'].mean().sort_values(ascending=False).h
```

```python
top_10_location.plot(kind='bar')
plt.title('Top_10_location as per property Price')
plt.show()
```

# Bivariate Analysis

```python
t=1
plt.figure(figsize=(10,10))
for i in num_col:
    if i!='price':
        plt.subplot(2,2,t)
        sns.scatterplot(x=df[i],y=df['price'])
        t+=1
plt.tight_layout()
plt.show()
```

```python
sns.heatmap(df.select_dtypes(include=np.number).corr(),vmax=1,vmin=-1,annot=True,cmap=
plt.show()
```

```python
t=1
plt.figure(figsize=(10,5))
for i in cat_col:
```

```
    if i!='location':
        plt.subplot(1,2,t)
        sns.boxplot(x=df[i],y=df['price'])
        t+=1
plt.tight_layout()
plt.show()
```

# Data Preprocessing

```
In [ ]:  ((df.isnull().sum()/df.shape[0])*100).sort_values(ascending=False)
```

```
In [ ]:  df[df['location'].isnull()]
```

```
In [ ]:  df.drop(index=[568],inplace=True)
```

```
In [ ]:  df[df['size'].isnull()]
```

```
In [ ]:  # lets impute this with mode
         df['size'].mode()
```

```
In [ ]:  df['size'].fillna(df['size'].mode()[0],inplace=True)
```

```
In [ ]:  df[df['bath'].isnull()]
```

```
In [ ]:  df.groupby('size')['bath'].median()
```

```
In [ ]:  df['size'][df['bath'].isnull()].unique()
         #So only for size of 1 to 5 are null so we are replacing
```

```
In [ ]:  df['bath']=np.where(((df['bath'].isnull())&(df['size']==1)),1,df['bath'])
         df['bath']=np.where(((df['bath'].isnull())&(df['size']==2)),2,df['bath'])
         df['bath']=np.where(((df['bath'].isnull())&(df['size']==3)),3,df['bath'])
         df['bath']=np.where(((df['bath'].isnull())&(df['size']==4)),4,df['bath'])
         df['bath']=np.where(((df['bath'].isnull())&(df['size']==5)),5,df['bath'])
```

```
In [ ]:  df[df['balcony'].isnull()]
```

```
In [ ]:  df.groupby('size')['balcony'].median()
```

```
In [ ]:  df['balcony'].fillna(df['balcony'].median(),inplace=True)
```

```
In [ ]:  df.isnull().sum()
```

```
In [ ]:  # Treat outliers
```

```
In [ ]:  t=1
         plt.figure(figsize=(10,10))
         for i in num_col:
             plt.subplot(3,2,t)
             sns.boxplot(x=df[i])
             t+=1
```

```
plt.tight_layout()
plt.show()
```

```
# lets drop the  extreme values (values > 95 percentile)
outlier_col = ['size','total_sqft','bath','price']
```

```
per_99 = df[outlier_col].quantile(0.99)
```

```
per_99
```

```
df=df[~(df[outlier_col]>per_99).any(axis=1)]
```

```
df.shape
```

```

```

```
t=1
plt.figure(figsize=(10,10))
for i in num_col:
    plt.subplot(3,2,t)
    sns.boxplot(x=df[i])
    t+=1
plt.tight_layout()
plt.show()
```

# Encoding

```
df['availability'] =np.where (df['availability']== 'Ready To Move',1,0)
```

```
df.sample(20)
```

```
(df['location'].value_counts(normalize=True)*100).cumsum()
```

```
top_10_properties=df.groupby('location')['price'].mean().sort_values(ascending=False).
```

```
top_10_properties.index
```

```
def location(s):
    if s in top_10_properties.index:
        return 1
    else:
        return 0
```

```
df['location']=df['location'].apply(location)
```

```
df.head(5)
```

```
df=pd.get_dummies(df,drop_first=True)
```

```
df
```

# predictive Modelling

```python
def model_validation(xtrain,ytrain,xtest,ytest,model):
    global m
    m=model
    m.fit(xtrain,ytrain)
    print('Training Scores')
    pred=m.predict(xtrain)
    print('R2:',r2_score(ytrain,pred))
    print('MSE:',mean_squared_error(ytrain,pred))
    print('RMSE:',mean_squared_error(ytrain,pred)**0.5)
    print('MAPE:',mean_absolute_percentage_error(ytrain,pred))


    print('Testing Scores')
    pred=m.predict(xtest)
    print('R2:',r2_score(ytest,pred))
    print('MSE:',mean_squared_error(ytest,pred))
    print('RMSE:',mean_squared_error(ytest,pred)**0.5)
    print('MAPE:',mean_absolute_percentage_error(ytest,pred))
```

# Train test split

```python
x=df.drop(columns=['price'],axis=1)
y=df['price']

x_train,x_test,y_train,y_test =train_test_split(x,y,train_size=0.8,random_state=1)
```

# Linear Regression

```python
model_validation(x_train,y_train,x_test,y_test,LinearRegression())
```

```python
# coefficients of linear regression
pd.DataFrame({'Features':x.columns,'Coef':m.coef_})
```

# Lasso

```python
# Lasso
model_validation(x_train,y_train,x_test,y_test,Lasso(alpha=0.005))
```

```python
# coefficients of Lasso
pd.DataFrame({'Features':x.columns,'Coef':m.coef_})
```

# Ridge

```python
model_validation(x_train,y_train,x_test,y_test,Ridge(alpha=0.5))# Ridge
```

```python
# coefficients of Ridge
pd.DataFrame({'Features':x.columns,'Coef':m.coef_})
```

# Elastic Net

```python
# Elastic Net
model_validation(x_train,y_train,x_test,y_test,ElasticNet(alpha=0.01,l1_ratio=0.1))
```

```python
# coefficients of Elastic Net
pd.DataFrame({'Features':x.columns,'Coef':m.coef_})
```

# Grid Search Cv

```python
## Using Grid Search CV
model_validation(x_train,y_train,x_test,y_test,ElasticNet(alpha=0.001,l1_ratio=0.25))
```

# Feature Selection

```python
from mlxtend.feature_selection import SequentialFeatureSelector
```

```python
sfs = SequentialFeatureSelector(estimator=LinearRegression(),
                                k_features='best',
                                scoring='r2')
```

```python
sfs.fit(x_train,y_train)
```

```python
k = list(sfs.k_feature_names_)
```

```python
sfs.k_score_
```

```python
model_validation(x_train[k],y_train,x_test[k],y_test,LinearRegression())
```

# Improve the prediction of the model

```python
#different aooroach in Encoding
#different techinque of missing values Traetment
#Use some advance models for prediction
```

```python
# Lets take the Linear Regression Model as final model
```

```python
model_validation(x_train,y_train,x_test,y_test,LinearRegression())
```

```python
m.feature_names_in_
```

```python
m.predict([[1,0,3,1500,2,2,1,0,0]])
```

```python
df.head()
```

# Save this model as trained model

```python
import pickle
```

```python
with open('model_lr.pkl','wb') as file:
    pickle.dump(m,file)

print('The file has been saved')
```