

ASAD ASHRAF KARL

PYTHON LIBRARIES

To avoid warnings:

```
In [ ]: 1 from warnings import filterwarnings
        2 filterwarnings('ignore')
```

Python Fundamental libraries:

```
In [1]: 1 import numpy as np
        2
        3 import pandas as pd
        4
        5 import math
        6
        7 import random
```

Visualizations:

```
In [10]: 1 import matplotlib.pyplot as plt
        2
        3 import seaborn as sns
```

STATISTICS:

```
In [3]: 1 import scipy.stats as st
2
3 from itertools import combinations
4 # list(combinations(data,x)) , x=1,2,3,... number of combinations
5
6 import statsmodels.api as sm
7 from statsmodels.formula.api import ols
8 # model = ols('numerical_variable'~'Categorical_variable', data).fit()
9 # sm.stats.anova_lm(model)
10
11 # If we reject null hypothesis, and wish to all possibilites in difference o
12 from statsmodels.stats.multicomp import pairwise_tukeyhsd
13 # pairwise_tukeyhsd('numerical_variable', 'categorecal_variable', alpha).sum
14
15 # Varias functions from scipy:
16 from scipy import stats
17 from scipy.stats import shapiro
```

Split

```
In [ ]: 1 from sklearn.model_selection import train_test_split
2
3 x=df.drop('target_column', axis = 1)#drop the target variable
4 y=df['target_column']
5 xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size = 0.2, random_state = 42)
```

FEATURE ENGINEERING:

STANDARDIZATION:

z_score

```
In [12]: 1 from scipy.stats import zscore
2 # df['New_column'] = zscore(df['On_which_column'])
```

NORMALIZATION:

MinMax scaling

```
In [5]: 1 from sklearn.preprocessing import MinMaxScaler
2 # mm = MinMaxScaler(feature_range=(0, 1)) To fix the range according to use
3 # df['New_column'] = mm.fit_transform(df[['On_which_column']])
```

OR

```
In [ ]: 1 df['New_Column'], lmda = st.boxcox(df['Column_name'])
```

LABEL ENCODING:

```
In [7]: 1 from sklearn.preprocessing import LabelEncoder
2 # LE = LabelEncoder()
3 # df['New_column'] = LE.fit_transform(df['On_which_column'])
```

ONE HOT ENCODING:

```
In [9]: 1 # pd.get_dummies(df['Column_name'])
```

Various functions from statsmodel to perform linear regression

```
In [ ]: 1 import statsmodels
2 import statsmodels.api as sm
3 import statsmodels.stats.api as sms
4 from statsmodels.compat import lzip
5 from statsmodels.stats.outliers_influence import variance_inflation_factor
6 from statsmodels.graphics.gofplots import qqplot
7 from statsmodels.stats.anova import anova_lm
8 from statsmodels.formula.api import ols
9 from statsmodels.tools.eval_measures import rmse
10 from statsmodels.formula.api import ols
```

Sklearn libraries:

```
In [ ]: 1 from sklearn.model_selection import train_test_split, GridSearchCV
2 from sklearn.linear_model import Linear_regression, Lasso, Ridge, LassoCV, R
3 from sklearn.feature_selection import import RFE
4 from sklearn.ensemble import RandomForestRegressor, BaggingRegressor
5 # RandomForestRegressor Learns the maximum of the data. (Suitable model) , B
```

'metrics' from sklearn is used for evaluating the model performance

```
In [ ]: 1 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_scor
```

END