

```

1 #Importing the library
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler
7 from scipy.stats import norm
8 from sklearn.linear_model import LinearRegression
9 from sklearn.model_selection import train_test_split
10 from sklearn.feature_selection import VarianceThreshold
11 from sklearn.feature_selection import SelectKBest
12 from sklearn.feature_selection import f_regression as f_regression_for_r
13 from sklearn.metrics import mean_squared_error
14 from sklearn.preprocessing import PolynomialFeatures
15 from sklearn.linear_model import Ridge, Lasso
16 from sklearn.metrics import mean_squared_error
17 from sklearn.model_selection import train_test_split
18 from sklearn.pipeline import
19
20 #Load the dataset
21 df=pd.read_csv('dataset/Lah20_prepared.csv')
22
23 #peek of dataset
24 df.head()
25
26 #about the dataset
27 df.info()
28
29 #checking for the null values
30 df.isnull().sum()

```

### b. Correlation & heatmap

```

M 1 #correlation of the features
   df.corr()

M 1
   2 #plotting the heatmap
   3 plt.figure(figsize=(10,10))
   4 sns.heatmap(df.corr(),
   5             cbar=True,
   6             square=True,
   7             annot=True,
   8             fmt='.1f',
   9             annot_kws={'size':12},
10             yticklabels=df.columns,
11             xticklabels=df.columns)

```

- Observations: This shows that 'log\_inst\_review' has high correlation with 'log\_enrollment' and 'log\_number\_ratings'. 'log\_inst\_review' has high correlation with 'log\_enrollment' and 'log\_number\_ratings'. 'log\_inst\_student' and 'log\_inst\_student' has 1.0 correlation. This shows that 'log\_inst\_review' and 'log\_inst\_student' have same data.

```

11 #Multicollinearity analysis
12 df_multi=df
13 df_independent = df_multi.select_dtypes(include=np.number).drop('log_inst_review', axis=1)
14 # create a dataframe for VIF data
15 vif_data = pd.DataFrame()
16 vif_data["feature"] = df_independent.columns
17
18
19 vif_data["vif"] = [variance_inflation_factor(df_independent.values, i)
20                    for i in range(len(df_independent.columns))]
21
22 ]
23
24 print(vif_data)

```

## 6. Linear Regression Model with Ridge

Observations: It shows the high multicollinearity in the Log\_enrollment and log\_number\_ratings. we should handle it for analysis.

```

1 scaled_df1.drop('log_enrollment', axis=1, inplace=True)
2 create a dataframe for VIF data
3 vif_data = pd.DataFrame()
4 vif_data['feature'] = scaled_df1.columns
5
6 vif_data['vif'] = [variance_inflation_factor(scaled_df1.values, i)
7                   for i in range(len(scaled_df1.columns))]
8
9
10 print(vif_data)

```

### c. Univariate Analysis

```
1 #Univariate analysis
2 #Distribution plot
3 sns.displot(df.log inst student)
```

- Observation: this graph is left skewed, it has a long left tail.

```
1 #find the index and drop
2 x=df[df.log_instudent<7].index
3 print(x)
4 df.drop(x,inplace=True)
5 # plotting the graph
6 sns.distplot(df.log_instudent,fit=norm)
```

#### d. Multivariate Analysis

```
1 #MultiVariate Analysis
```

- Observation: They are growing in the same direction and not disperse. Means they are

- Observation: They are growing in same direction but they are a lot disperse

#### 4.Feature Selection

```

# 1. Created a column of the dataframe to a new variable called target
H 1 # create a copy of data frame
H 2 # df.copy()
H 3 # df['target']
H 4 # target=df.copy(log_log_review)
H 5 # features=df.copy(df.drop('log_log_review',axis=1))

H 1 # df.copy.info()

# b. Correlation Based Selection(Threshold)
H 1 # imported the function
H 2 # from Functions/Functions_W4146384.ipynb
H 3 #
H 4 # using threshold value to select the best values
H 5 # called the method and assigned it to df
H 6 # df=correlation_based_selection(features,df.copy)
H 7 # df.correlation.info()

# c. Variance Threshold Selection
H 1 # df.copy.info()

H 1 #
H 2 # # run the function
H 3 # from Functions/Functions_W4146384.ipynb
H 4 #
H 5 # calling the function and using the variance threshold
H 6 # df=select_variance(features,df)
H 7 #
H 8 # printing the df
H 9 # df.variance.info()

# d. Select K-Best method
H 1 #
H 2 # # run the function
H 3 # from Functions/Functions_W4146384.ipynb
H 4 #
H 5 # calling the function
H 6 # df=select_kbest(features,k)
H 7 #
H 8 # output the dataset
H 9 # df.selectbest.head()

H 1 # df.selectbest.info()

```

```

def method = ['correlation', 'variance', 'silhouette']

if method=='correlation':
    features=df.correlation
elif method=='variance':
    features=df.variance
elif method=='silhouette':
    features=df.silhouette

for transform in ['No Transformation', 'Polynomial Transform', 'Log Transformation']:
    if transform=='Polynomial Transform':
        from Functions import PolynomialTransform
        transformed_df=PolynomialTransform(features)
    elif transform=='Log Transformation':
        from Functions import LogTransform
        transformed_df=LogTransform(features)
    X_train, X_test, Y_train, Y_test = train_test_split(transformed_df, y, test_size=0.2, random_state=0)

    # Create a scalar instance
    scalar = StandardScaler()

    # Fit the scalar on the training data
    scaled_train=scalar.fit_transform(X_train)
    scaled_test=scalar.transform(X_test)

    # reshape the dataframe
    X_train_scaled=pd.DataFrame(scaled_train)
    X_test_scaled=pd.DataFrame(scaled_test)

    # build the model
    model = LinearRegression()
    model.fit(X_train, Y_train)

    # get the predictions for the test set
    Y_pred = model.predict(X_test)

    # R-squared
    r2 = model.score(X_test, Y_test)
    print("R^2 = ", r2)

    from sklearn.metrics import mean_squared_error
    mse = mean_squared_error(Y_test, Y_pred)
    print("MSE is ", mse)

    #appending the values to Lists
    feature_list.append(transform)
    transform_list.append(transforms)
    r2_list.append(r2)
    mse_list.append(mse)

```

## 7. Plot and summary analysis

```

1 # Making a table using the values of list
2 result=np.vstack((fselection_list,tftransform_list,r2_list,rmse_list)).T
3 result_df=pd.DataFrame(result,columns=['Feature Selection', 'Feature Transform', 'R2', 'RMSE'])
4 result_df
5
6 # Selecting the values based on the lowest value for RMSE
7 result_df.sort_values(by=['RMSE', 'R2'],ascending=[True,False]).head(1)
8
9 # Running the query to make the best regression model and make predictions
10 features_df.select
11
12 # Run functions/Functions_K0146384.ipynb
13 transformed_df=make_poly(features)
14 X_train, X_test, Y_train, Y_test = train_test_split(transformed_df,target, test_size=0.25, random_state=0)
15 # Create a scaler instance
16 scaler = StandardScaler()
17
18 # Fit the scaler on the training data
19 scaled_train_X=scaler.fit_transform(X_train)
20 scaled_test_X =scaler.fit_transform(X_test)
21
22 #Rename the dataframe
23 X_train_scaled=pd.DataFrame(scaled_train_X,columns=X_train.columns)
24 X_test_scaled=pd.DataFrame(scaled_test_X,columns=X_test.columns)
25
26 # build the model
27 model = LinearRegression()
28 model.fit(X_train, Y_train)
29 # get the predictions for the test dataset
30 Y_pred = model.predict(X_test)
31
32 # printing the prediction
33 print(Y_pred)
34
35 transformed_df.head()
36
37 # Print the coefficients of the best linear model
38 print(pd.Series(model.coef_, index=transformed_df.columns))
39
40 # Scatter plot using Y_pred and Y_test
41 plt.scatter(Y_train, Y_pred)

```

## 8. Out of Sample Prediction

```

M 1 df_selBest.head()
M 1 df_selBest.describe()
M 1 synthetic_data = df_selBest.mean().to_frame().T
   synthetic_data = df_selBest.quantile(0.05).to_frame().T
3 synthetic_data = df_selBest.quantile(0.75).to_frame().T
4 # Example selecting rows 0 and 2
5 synthetic_data = df_selBest.iloc[[0, 2]]
6
M 1 #creating a DataFrame
2 synthetic_data=pd.DataFrame(synthetic_data, columns=df_
M 1 synthetic_data.head()
M 1 # Apply the same feature transformation as the best lin
2 poly_transform = make_poly(synthetic_data)
3
M 1 # Apply the scaling as the best linear regression model
2 scaler = StandardScaler()
3 # Fit the scaler on the training data
4 scaled_train_X=scaler.fit_transform(poly_transform)
5
M 1 # Make predictions on the synthetic dataset
2 y_synthetic_pred = model.predict(scaled_train_X)
3
4 # Print the predictions
5 print("Predictions on Synthetic Dataset:")

```

#### 5.4 Feature Scaling

```

# 1 from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler

DO NOT SCALE the target variable

5.4.1 Standard Scaler

The standard scaler tries to standardize each value in a feature by removing the mean and scaling the values

# 1 ss = StandardScaler()
# 2 scaled_features = ss.fit_transform(df_manual)

# you need to rescale the dataframe
# 3 scaled_df = pd.DataFrame(scaled_features, columns=df_manual.columns)
# 4 scaled_df

5.4.2 MinMax Scaler

In MinMax Scaler, features will be scaled into range [0, 1] giving it by the range between the Max() and
lower bound using feature_range variable.

# 1 ms = MinMaxScaler()
# 2 scaled_features = ms.fit_transform(df_manual)
# 3 you need to rescale the dataframe
# 4 scaled_df = pd.DataFrame(scaled_features, columns=df_manual.columns)
# 5 scaled_df

5.4.3 Robust Scaler

The robust scaler tries to remove the disadvantage of the MinMax scaler such that it will not be affected by
measures known as inter-quartile range(IQR) that is the range between the 75% percentile and 25% percentile

# 1 rs = RobustScaler()
# 2 scaled_features = rs.fit_transform(df_manual)

```

### Multicollinearity analysis

```

# We need to make sure that none of the features have a high multicollinearity among themselves to have a
# good fit.

# 1. df_independent = df.select_types(include=op.number).drop('price', axis=1)
# 2. # create a dataframe for VIF data
# 3. vif_data = pd.DataFrame()
# 4. vif_data['feature'] = df_independent.columns
# 5. vif_data['vif'] = [variance_inflation_factor(df_independent.values, i)
# 6.                     for i in range(len(df_independent.columns))]
# 7. ]
# 8. print(vif_data)

# 9. scaled_features = StandardScaler().fit_transform(df_independent)
# 10. scaled_dfi = pd.DataFrame(scaled_features, columns=df_independent.columns)
# 11. # create a dataframe for VIF data
# 12. vif_data = pd.DataFrame()
# 13. vif_data['feature'] = scaled_dfi.columns
# 14. vif_data['vif'] = [variance_inflation_factor(scaled_dfi.values, i)
# 15.                     for i in range(len(scaled_dfi.columns))]
# 16. ]
# 17. print(vif_data)

# 18. scaled_dfi_drop['tax', axis=1, inplace=True)
# 19. # create a dataframe for VIF data
# 20. vif_data = pd.DataFrame()
# 21. vif_data['feature'] = scaled_dfi.columns
# 22. vif_data['vif'] = [variance_inflation_factor(scaled_dfi.values, i)
# 23.                     for i in range(len(scaled_dfi.columns))]
# 24. ]
# 25. print(vif_data)

# 26. df_drop['tax', axis=1, inplace=True)

```

### Feature Selection, Feature Transformation and Feature Scaling

Univariate Analysis

### Training and Test

```
1 # build the model
2 model = LinearRegression()
3 model.fit(X_train, y_train)

1 # get the predictions for the test dataset
2 y_pred = model.predict(X_test)
```

Calculate the performance metrics

```

1 # R-squared
2 r2 = model.score(test, 'r2')
3 print("R score is: ", r2)
4
5 from sklearn.metrics import mean_squared_error
6 # MSE
7 mse = mean_squared_error(test, y_pred, squared=False)
8 print("MSE is: ", mse)
9
10 # plot the predictions vs actual
11 plt.scatter(y_pred, test, alpha=0.7)
12 plt.title('Linear Regression Model')
13 plt.xlabel('Predicted Value')
14 plt.ylabel('Actual Value')
15 plt.show()
16
17 # the coefficient list
18 feat = list(X_train.columns)
19 coef = model.coef_.transpose()
20 pd.DataFrame(list(feat,coef), columns=['feature', 'coefficients'])

```

```

1 # Let's ridge them!
2 rmse_list = []
3 r2_list = []
4 model_list = []
5
6 for a in alphas:
7     ridge = Ridge(alpha=a, max_iter=1000)
8     ridge.fit(X_train, Y_train)
9     pred = ridge.predict(X_test)
10
11     r2_list.append(ridge.score(X_train,Y_train))
12     rmse = mean_squared_error(Y_test, pred)**0.5;
13     rmse_list.append(rmse)
14     model_list.append(ridge)
15     print("Alpha",a,"RMSE",rmse)
16
17 ridge_result = np.vstack((alphas, rmse_list, r2_list)).T
18 ridge_df = pd.DataFrame(ridge_result, columns=['Alpha', 'RMSE', 'R2'])
19 ridge_df

```

```
1 ridge_df.sort_values(by=['RMSE', 'R2'], ascending=[True, False]).head(1)
```

```
1 # find out the coefficient for the best result
```

```
2 best_model = model_list[14]
```

```
1 print(pd.Series(best_model.coef_, index=features.columns))
```

## Ridge Regression

```
1 # import the libraries
2 from sklearn.linear_model import Ridge, Lasso
3 from sklearn.metrics import mean_squared_error
4 from sklearn.model_selection import train_test_split
5 import numpy as np
6 import pandas as pd
7 import matplotlib.pyplot as plt

1 # Load the dataset
2 df = pd.read_csv("boston_modified.csv")
3 df.head()
4
5 # just drop all NaN
6 df.dropna(inplace=True)
7
8 # just transform all categorical
9 # we need to be sure that all the NaN are removed and all other columns are numerical
10 df = pd.get_dummies(df, drop_first=True)
```

The `Ridge()` function has an `alpha` argument that is used to tune the model. We'll generate an array of `alpha` values `alpha_vals` essentially covering the full range of scenarios from the null model containing only the intercept, to the least squares

## Univariate Analysis

In univariate analysis, we try to see the distribution of some interesting feature. We can use histogram, distplot, bar chart, boxplot or violin plot for this purpose.

## Lasso Regression

The `Lasso()` function has an `alpha` argument that is used to tune the model. We'll generate an array of `alpha` values ranging 1 to a small value:

```
1 alphas = np.linspace(1, 0.001, 15)
2 alphas
```

We will use the same dataset

```

1 xtrain, xtest, ytrain, ytest = train_test_split(features, target, test_size=0.25)

2 # Let's put them in the Lasso yee haw
3 rmse_list = []
4 r2_list = []
5 model_list = []
6
7 for a in alphas:
8     lasso = Lasso(alpha=a, max_iter=1000)
9     lasso.fit(X_train, y_train)
10    pred = lasso.predict(X_test)
11
12    r2_list.append(lasso.score(X_train, y_train))
13    rmse = mean_squared_error(y_test, pred)**0.5;
14    rmse_list.append(rmse)
15    model_list.append(lasso)
16    print("Alpha", a, "RMSE", rmse)
17
18 lasso_df = pd.DataFrame(zip(alphas, rmse_list, r2_list), columns=['Alpha', 'RMSE', 'R2'])

```

```
1 lasso_df.sort_values(by=['RMSE', 'R2'], ascending=[True, False]).head(1)
```

  

```
1 # find out the coefficient of the best result obtained from our lasso
2 best_model = model_list[14]
3 print(pd.Series(best_model.coef_, index=features.columns))
```

