```python
"""
Version: 1.0
Created on: Sat Mar 20 18:55:27 2021
Author: Balaji Kannan
Description: Lower Back Pain Dataset from Kaggle
    1. https://www.kaggle.com/sammy123/lower-back-pain-symptoms-dataset
    2. This pipeline is custom built for a single target variable with Datatype
    Object.
    3. Binary classification problem.

"""

#%% Library

import pandas as pd
pd.set_option('display.max_columns', 20)
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

#%% IO Path & Dataframe Definitions

PATH = r"C:\DSML_Case_Studies\02_Logistic_Regression\Input"
FNAME = r"\Dataset_Lower_Back_Pain.csv" # Prefix LrBkPn

OUTPATH = r"C:\DSML_Case_Studies\02_Logistic_Regression\Output"
PREFIX = r"\LrBkPn_"

df = pd.read_csv(f"{PATH}{FNAME}")
df = df.round(decimals=3) # rounding the decimals

targetvar = 0
while targetvar <= 0:
    targetvar = int(input("Enter # of Target Variables: ")) # user specifies # of
    targets

collst = []
for columns in df.columns:
    collst.append(columns)

featlst = collst[0:len(collst)-targetvar]
targlst = collst[-targetvar:]

# Sanity Checks

print(df.head(), sep='\n')
print("List of Features:", featlst, sep='\n')
print("List of Targets:", targlst, sep='\n')

#%% Exploratory Data Analysis

OFNAME1 = r"01_Descriptive_Stats.txt"

desc_stat = df.describe() # Univariate analyses
print(desc_stat)
FOUT1 = open(f"{OUTPATH}{PREFIX}{OFNAME1}", 'w+')
desc_stat.to_string(FOUT1)
FOUT1.close()

# Generating n*4 matrix of box plots

n_rows = len(collst)//4
fig, axes = plt.subplots(n_rows, 4, figsize = (15,15))
axes = axes.flatten()

FIG1 = r"_FIG01_Boxplot"
for i in range(0,len(df.columns)-targetvar):
    sns.boxplot(x=targlst[0], y=df.iloc[:,i], data=df, orient='v', ax=axes[i])
    plt.tight_layout()
    plt.savefig(f"{OUTPATH}{PREFIX}{FIG1}")

# Linear Correlation Heatmap
```

```python
71
72     cormethod = {0:'pearson', 1:'kendall', 2:'spearman'}
73     for i in range(0, 3, 1):
74         temp = 'linear_cor' + str(i)
75         temp = df.corr(method=cormethod[i])
76         ftemp = cormethod[i].title()
77         FIG2 = r"_FIG02_Corr_"
78         mask = np.zeros(temp.shape, dtype=bool)
79         mask[np.tril_indices(len(mask))] = False
80         plt.subplots(figsize=(20,15))
81         plt.title(f"{ftemp} Corrlelation")
82         sns.heatmap(temp, annot=True, vmin=-1, vmax=1, center=0,
83                     cmap='coolwarm', square=True, mask=mask)
84         plt.savefig(f"{OUTPATH}{PREFIX}{FIG2}{ftemp}")
85
86     # Non-Linear Correlation Predictive Power Score - Heatmap
87
88     import ppscore as pps
89
90     FIG3 = r"_FIG03_Predictive_Power_Score"
91
92     ppscorr = pps.matrix(df) # Predictive Power Score - PPS
93     matrix_df = pps.matrix(df)[['x', 'y', 'ppscore']].pivot(columns='x', index='y',
94                                                             values='ppscore')
95     plt.subplots(figsize=(20,15))
96     sns.heatmap(matrix_df, cmap="Greens", annot=True, linewidth=0,
97     annot_kws={"size":12}, fmt='.2g')
97     plt.savefig(f"{OUTPATH}{PREFIX}{FIG3}")
98
99     # Scatter Plot with Hue
100
101    FIG4 = r"_FIG04_Scatter_Plot"
102
103    grid1 = sns.PairGrid(df, hue=targlst[0])
104    grid1.map(plt.scatter)
105    grid1.map_diag(sns.kdeplot)
106    grid1.add_legend()
107    grid1.fig.suptitle("Scatter Plot", y=1.01)
108    grid1.savefig(f"{OUTPATH}{PREFIX}{FIG4}")
109
110    #%% Pandas Profile Report
111
112    from pandas_profiling import ProfileReport
113
114    OFNAME2 = r"02_Descriptive_Stats.html"
115    report = ProfileReport(df) # Descriptive statistics report
116    report.to_file(f"{OUTPATH}{PREFIX}{OFNAME2}") # Rendering to HTML
117
118    #%% High Dimensional Interactive Plot - HD Plot
119
120    import hiplot as hip
121
122    OFNAME3 = r"03_Parallel_Plot.html"
123    parplot = hip.Experiment.from_dataframe(df)
124    parplot.to_html(f"{OUTPATH}{PREFIX}{OFNAME3}")
125
126    #%% Machine Learning Model - Pre-Processing
127
128    from sklearn import preprocessing
129    from sklearn.preprocessing import StandardScaler
130    from sklearn.model_selection import train_test_split
131    from sklearn.linear_model import LogisticRegression
132    from sklearn.metrics import classification_report
133
134    import statsmodels.api as sm
135
136    df.loc[df.Class_att=='Abnormal',targlst[0]] = 1
137    df.loc[df.Class_att=='Normal', targlst[0]] = 0
138    X = df.drop(columns=targlst)
139    y = df.filter(targlst, axis=1)
140
141
```

```python
142    def data_split(X,y):
143        """
144        Parameters
145        ----------
146        X : Feature Variables
147        y : Target Variables
148        Returns
149        -------
150        Split dataframe into train and test
151        """
152        X_train, X_test, y_train, y_test = train_test_split(X, y.values.ravel(),
           test_size=0.3,
153                                                             random_state=39)
154        scaler = StandardScaler(copy=True, with_mean=True, with_std=True)
155        scaler.fit(X_train)
156        train_scaled = scaler.transform(X_train) # Only feature variables are scaled
157        test_scaled = scaler.transform(X_test) # only features varaibels are scaled
158        y_train=y_train.astype('int')
159        y_test=y_test.astype('int')
160        return(train_scaled, test_scaled, y_train, y_test)
161
162    #%% Logistic Regression Model
163
164    def logistic_regression(x,y):
165        """
166        Parameters
167        ----------
168        x : Feature Variables
169        y : Target Variables
170
171        Returns
172        -------
173        Logistic Regression Fit
174
175        """
176        logreg = LogisticRegression().fit(x, y)
177        return(logreg)
178
179    X_train_scaled, X_test_scaled, y_train, y_test = data_split(X,y)
180    logreg_result = logistic_regression(X_train_scaled, y_train)
181
182    print("Training set score:
       {:.3f}".format(logreg_result.score(X_train_scaled,y_train)))
183    print("Test set score: {:.3f}".format(logreg_result.score(X_test_scaled,y_test)))
184
185    #%% Logit with Stats Model
186    """
187    1. Scikit LogisticRegression is good in predicting target variable on a test set.
188    2. It did not interpret anything about the individual features.
189    3. Which variable(s) influence the Target variable more?
190    4. Use logit from stats model to answer questions 2 & 3
191    """
192
193    logit_model = sm.Logit(y_train, X_train_scaled)
194    result = logit_model.fit()
195    print(result.summary2())
196
197    """
198    1. Output such as:
199        Maximum Likelihood optimization failed to converge -
200        indicates that there is multicoliniarity.
201        Removing those variable would improve the convergence.
202    """
203
204    #%% Feature Reduction Tryouts
205
206    features_to_drop = ['pelvic_incidence', 'lumbar_lordosis_angle'] # From Correlation
       Heatmap
207
208    new_lst = targlst + features_to_drop
209    X = df.drop(columns=new_lst)
210    y = df.filter(targlst, axis=1)
```

```python
211
212    X_train_scaled, X_test_scaled, y_train, y_test = data_split(X,y)
213    logreg_result = logistic_regression(X_train_scaled, y_train)
214
215    print("Training set score:
       {:.3f}".format(logreg_result.score(X_train_scaled,y_train)))
216    print("Test set score: {:.3f}".format(logreg_result.score(X_test_scaled,y_test)))
217
218    logit_model = sm.Logit(y_train, X_train_scaled)
219    result = logit_model.fit()
220    print(result.summary2())
221
222    #%% Consider only Variables for which P-value is < 0.05
223
224    features_to_drop1 = ['sacral_slope', 'pelvic_slope', 'Direct_tilt',
225                          'thoracic_slope', 'cervical_tilt', 'sacrum_angle',
226                          'scoliosis_slope', 'pelvic_tilt'] #
227    new_lst1 = targlst + features_to_drop + features_to_drop1
228    X = df.drop(columns=new_lst1)
229    y = df.filter(targlst, axis=1)
230    print(X.head())
231
232    X_train_scaled, X_test_scaled, y_train, y_test = data_split(X,y)
233    logreg_result = logistic_regression(X_train_scaled, y_train)
234
235    print("Training set score:
       {:.3f}".format(logreg_result.score(X_train_scaled,y_train)))
236    print("Test set score: {:.3f}".format(logreg_result.score(X_test_scaled,y_test)))
237
238    logit_model = sm.Logit(y_train, X_train_scaled)
239    result = logit_model.fit()
240    print(result.summary2())
241
242    #%% Model Metrics
243
244    # assigning the model predicted values to y_pred
245    y_pred = logreg_result.predict(X_test_scaled)
246
247
248    # assigning the string Normal and Abnormal to the 0 and 1 values respectively. This
       is useful in plotting
249    # the confusion matrix
250    y_pred_string = y_pred.astype(str)
251    y_pred_string[np.where(y_pred_string == '0')] = 'Normal'
252    y_pred_string[np.where(y_pred_string == '1')] = 'Abnormal'
253
254    y_test_string = y_test.astype(str)
255    y_test_string[np.where(y_test_string == '0')] = 'Normal'
256    y_test_string[np.where(y_test_string == '1')] = 'Abnormal'
257
258    from sklearn.metrics import confusion_matrix
259
260    FIG5 = r"_FIG05_Confusion_Matrix"
261    ax= plt.subplot()
262    labels = ['Abnormal','Normal']
263    cm = confusion_matrix(y_test_string, y_pred_string, labels)
264    sns.heatmap(cm, annot=True, ax = ax, cmap='coolwarm'); #annot=True to annotate cells
265    plt.savefig(f"{OUTPATH}{PREFIX}{FIG5}")
266
267    # labels, title and ticks
268    ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
269    ax.set_title('Confusion Matrix');
270    ax.xaxis.set_ticklabels(['Abnormal', 'Normal']);
       ax.yaxis.set_ticklabels(['Abnormal', 'Normal']);
271    plt.show()
272
273    print(classification_report(y_test, y_pred, target_names=labels))
```