

MTMC-512 Programming Lab IV(Machine Learning)

Lab Assignment 2: Data Preprocessing Pipeline

1 Task 1: Load and Explore the Dataset

1.1 Loading the Dataset

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 import warnings
7 warnings.filterwarnings('ignore')
8
9 df = pd.read_csv('heart_disease_uci.csv')
10 df.head()
```

Output:

	id	age	sex	dataset	cp	trestbps	chol	fbs \
0	1	63	Male	Cleveland	typical angina	145.0	233.0	True
1	2	67	Male	Cleveland	asymptomatic	160.0	286.0	False
2	3	67	Male	Cleveland	asymptomatic	120.0	229.0	False
3	4	37	Male	Cleveland	non-anginal	130.0	250.0	False
4	5	41	Female	Cleveland	atypical angina	130.0	204.0	False

1.2 Dataset Characteristics

```
1 df.shape
2 df.dtypes
3 df.describe()
4 df.isnull().sum()
5 df.duplicated().sum()
```

Output:

Shape: (920, 16)

Data Types:

id	int64
age	int64
sex	object
dataset	object
cp	object
trestbps	float64
chol	float64
fbs	object
restecg	object
thalch	float64
exang	object
oldpeak	float64
slope	object

```
ca          float64
thal        object
num         int64
```

Missing Values:

```
trestbps    59
chol        30
fbs         90
restecg      2
thalch      55
exang       55
oldpeak     62
slope       309
ca          611
thal        486
num         0
```

2 Task 2: Data Cleaning

2.1 Handling Missing Values

```
1 numeric = df.select_dtypes(include=['float64', 'int64']).columns
2 df[numeric] = df[numeric].fillna(df[numeric].median())
3
4 categoric = df.select_dtypes(include=['object']).columns
5 df[categoric] = df[categoric].fillna(df[categoric].mode())
```

Output:

Missing values after imputation:

```
id          0
age         0
sex         0
dataset     0
cp          0
trestbps    0
chol        0
fbs         0
restecg     0
thalch      0
exang       0
oldpeak     0
slope       0
ca          0
thal        0
num         0
```

2.2 Boxplot for Visualizing Outliers

```
1 num_columns = df.select_dtypes(include=['float64', 'int64']).columns
2 plt.figure(figsize = (20,12))
3 sns.boxplot(data=df[num_columns])
4 plt.show()
```

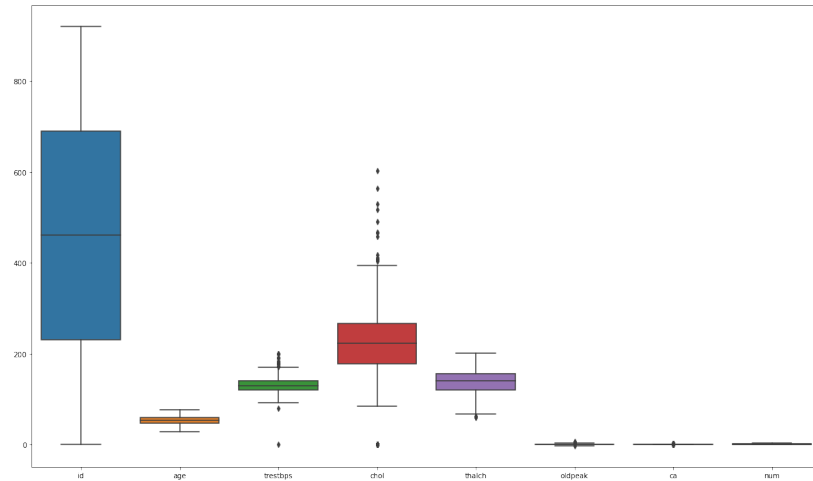


Figure 1: Box Plot For Outliers

2.3 Handling Outliers Using Z-scores

```

1 from scipy.stats import zscore
2 z_scores = df[numeric].apply(zscore)
3
4 zero_outliers = df[(z_scores < 3).all(axis=1)]
5 df = zero_outliers
6
7 df_shape = df.shape

```

Output:

Data after removing outliers using Z-scores: (887, 16)

3 Task 3: Feature Engineering

3.1 Encoding Categorical Features

```

1 df = pd.get_dummies(df, columns=categorical, drop_first=True)

```

3.2 Feature Scaling

```

1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 df[numeric] = scaler.fit_transform(df[numeric])

```

3.3 Correlation Analysis

```

1 corr = df.corr()
2 sns.heatmap(corr, annot=True, cmap="coolwarm")
3 plt.savefig("heatmap.png")
4 plt.show()

```

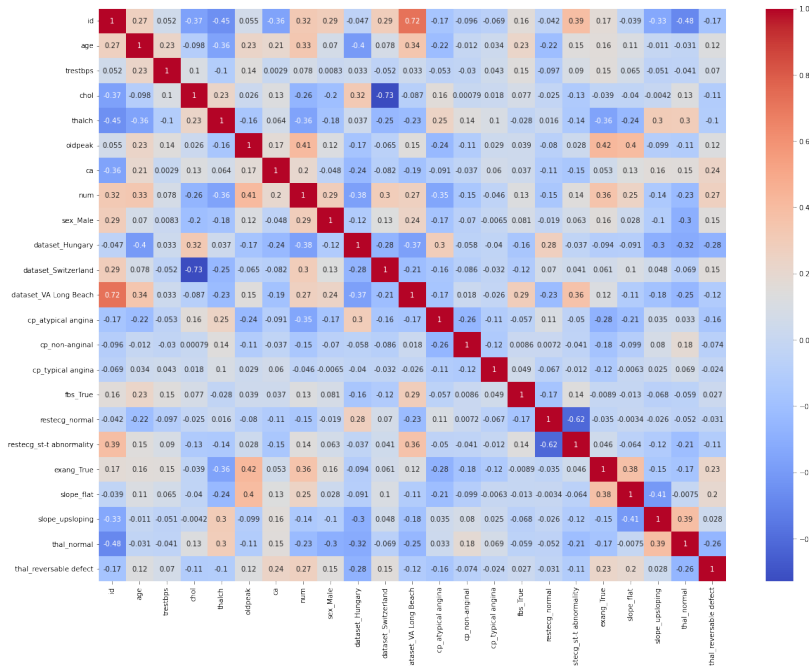


Figure 2: Feature Correlation Heatmap

4 Task 4: Classification using kNN

4.1 Splitting the Dataset

```
1 from sklearn.model_selection import train_test_split
2 X = df.drop('num',axis=1)
3 y = df['num']
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
```

4.2 Training k-Nearest Neighbors (kNN)

```
1 from sklearn.neighbors import KNeighborsClassifier
2 knn = KNeighborsClassifier(n_neighbors=5)
3 knn.fit(X_train, y_train)
```

4.3 Model Evaluation

```
1 from sklearn.metrics import accuracy_score, classification_report,
    confusion_matrix
2 y_pred = knn.predict(X_test)
3 accuracy = accuracy_score(y_test, y_pred)
4 print(accuracy)
5 print(classification_report(y_test, y_pred))
6 print(confusion_matrix(y_test, y_pred))
```

Output:

Accuracy: 0.8595505617977528

Classification Report:

	precision	recall	f1-score	support
0	0.88	0.97	0.92	155
1	0.33	0.10	0.15	20
2	0.00	0.00	0.00	3
accuracy			0.86	178
macro avg	0.40	0.36	0.36	178
weighted avg	0.80	0.86	0.82	178

Confusion Matrix:

```
[[151  4  0]
 [ 18  2  0]
 [  3  0  0]]
```

5 Conclusion

In conclusion, while the current model demonstrates a solid classification performance, there's room for improvement, particularly by experimenting with more advanced algorithms and fine-tuning the current model. Preprocessing steps, including outlier handling, missing value imputation, and feature scaling, were crucial in ensuring the dataset was ready for training and ultimately contributed to the overall model accuracy.