**Name : Ch Mubashir**

**SAP : 56892.**

**Lab Task 03 :**

**Data Preprocessing with heart Dataset.**
**Lab Tasks**
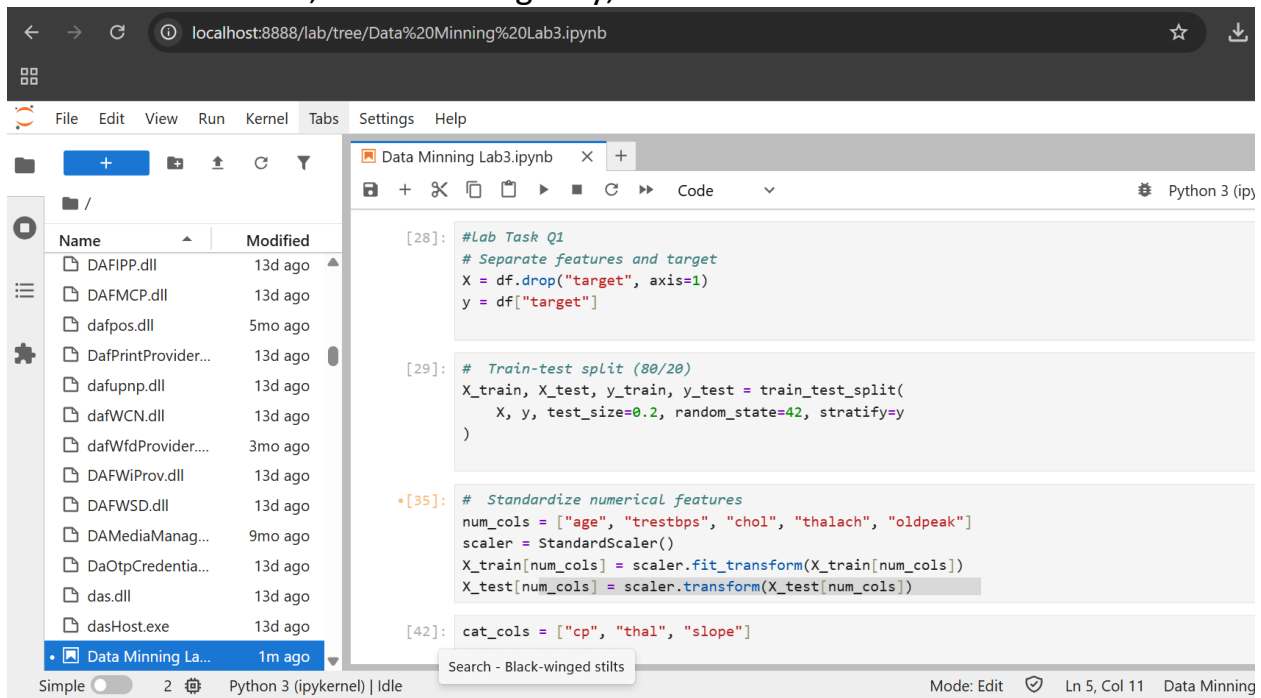1. **Separate features and target :**
 Target variable: target (1 = disease, 0 = no disease).

2. **Split the data.**
3. **Standardize numerical features :**
 Columns: age, trestbps, chol, thalach, oldpeak.
 Use StandardScaler, fit on training only, transform both train and test



4. **One-Hot Encode categorical features**
  Columns: cp, thal, slope.
  Use pd.get_dummies(..., drop_first=True).
5. **Check processed training data :**

## Exercises

1. **What are the mean and std of the standardized chol column in the training set? Why aren't they exactly 0 and 1 when checked with pandas.Series.std()?**

**Answer :** After standardizing the **chol** column, its mean becomes ~0 and std ~1 because StandardScaler uses the population formula (ddof=0). When checked with pandas.Series.std(), it's not exactly 1 since pandas by default uses the sample formula (ddof=1).

2. **Why do we use drop_first=True in one-hot encoding? What problem does it solve?**

**Answer :** The drop_first=True option in one-hot encoding avoids the **dummy variable trap** (perfect multicollinearity) by dropping one category and treating it as the baseline.

3. **If a new category of thal appeared in the test set but not in the training set, what would happen with pd.get_dummies? OneHotEncoder(handle_unknown='ignore') in scikit-learn help here?**

**Answer :** If a new category of **thal** appears in the test set but not in the training set, pd.get_dummies won't create a column for it, leading to mismatched features. Using OneHotEncoder(handle_unknown="ignore") ensures unseen categories are safely ignored, keeping columns consistent.

**Data Visualization :**

**Visualize** the heart dataset before doing any cleaning.



Here are some useful plots we can make:

1. **Target distribution** (0 = no heart disease, 1 = heart disease)

```python
[19]:  # Remove duplicates
       df = df.drop_duplicates()

[20]:  # Handle missing values (if any)
       df = df.dropna()

*[22]: # Target distribution (fixed warning)
       plt.figure(figsize=(6,4))
       sns.countplot(x="target", hue="target", data=df, palette="Set2", legend=False)
       plt.title("Target Distribution (Heart Disease: 1 = Yes, 0 = No)")
       plt.show()
```

## 2. Age distribution



```python
[23]:  # Age distribution
       plt.figure(figsize=(8,5))
       sns.histplot(df["age"], bins=20, kde=True, color="skyblue")
       plt.title("Age Distribution")
       plt.xlabel("Age")
       plt.ylabel("Count")
       plt.show()
```

## 3. Correlation heatmap of features

File   Edit   View   Run   Kernel   Tabs   Settings   Help

Data Minning Lab3.ipynb

Code

Python 3 (ipy

```
[24]: #Correlation heatmap
      plt.figure(figsize=(12,8))
      sns.heatmap(df.corr(), annot=True, cmap="coolwarm", fmt=".2f")
      plt.title("Correlation Heatmap")
      plt.show()
```

Correlation Heatmap



| | age | sex | cp | trestbps | chol | fbs | restecg | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| age | 1.00 | -0.09 | -0.06 | 0.28 | 0.21 | 0.12 | -0.11 | -0.40 | 0.09 | 0.21 | -0.16 | 0.30 | 0.07 | -0.22 |
| sex | -0.09 | 1.00 | -0.05 | -0.06 | -0.20 | 0.05 | -0.06 | -0.05 | 0.14 | 0.10 | -0.03 | 0.11 | 0.21 | -0.28 |
| cp | -0.06 | -0.05 | 1.00 | 0.05 | -0.07 | 0.10 | 0.04 | 0.29 | -0.39 | -0.15 | 0.12 | -0.20 | -0.16 | 0.43 |
| trestbps | 0.28 | -0.06 | 0.05 | 1.00 | 0.13 | 0.18 | -0.12 | -0.05 | 0.07 | 0.19 | -0.12 | 0.10 | 0.06 | -0.15 |
| chol | 0.21 | -0.20 | -0.07 | 0.13 | 1.00 | 0.01 | -0.15 | -0.01 | 0.06 | 0.05 | 0.00 | 0.09 | 0.10 | -0.08 |
| fbs | 0.12 | 0.05 | 0.10 | 0.18 | 0.01 | 1.00 | -0.08 | -0.01 | 0.02 | 0.00 | -0.06 | 0.14 | -0.03 | -0.03 |
| restecg | -0.11 | -0.06 | 0.04 | -0.12 | -0.15 | -0.08 | 1.00 | 0.04 | -0.07 | -0.06 | 0.09 | -0.08 | -0.01 | 0.13 |

Simple    2    Python 3 (ipykernel) | Idle                    Mode: Edit    Ln 5, Col 11    Data Minning