# Data Collection and Preprocessing Phase

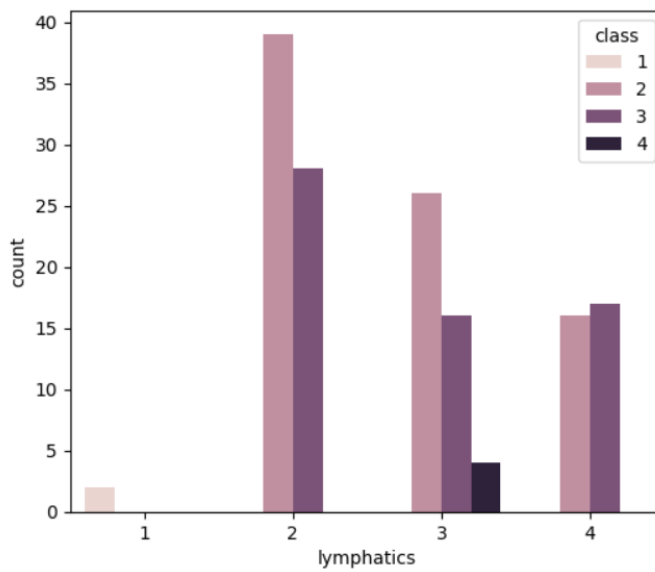| Date | 12 JULY 2024 |
|------|--------------|
| ID | 740036 |
| Project Title | Lymphography Classification using ML |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Report**

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.
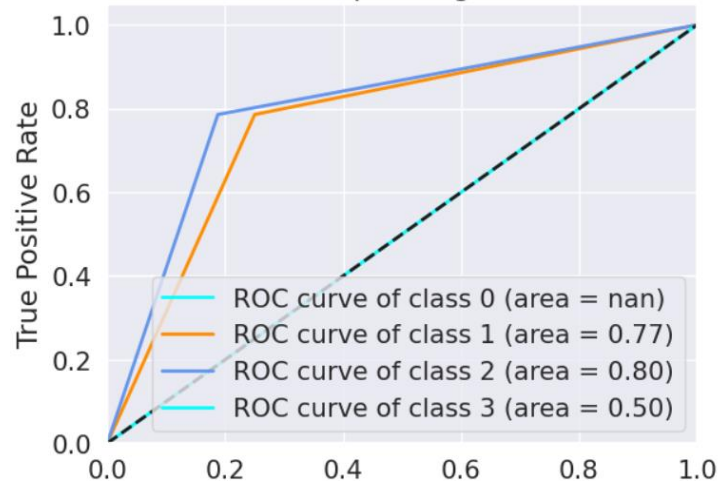
| Section | Description |
|---------|-------------|
| Data Overview | **Dimension:** 148 rows × 19 columns <br><br> **Descriptive statistics:** <br><br> *(see table below)* |
| | |

|  | class | lymphatics | block of affere | bl. of lymph. c | bl. of lymph. s | by pass | extravasates | regeneration of | early uptake in | lym.nodes dimin | lym.nodes enlar | changes in lym. | defect in node | changes in node | changes in stru |
|------|-------|-----------|----------|----------|----------|---------|--------------|-------------|-----------|-----------|-----------|-----------|---------|---------|---------|
| count | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 | 148.000000 |
| mean | 2.452703 | 2.743243 | 1.554054 | 1.175676 | 1.047297 | 1.243243 | 1.506757 | 1.067568 | 1.702703 | 1.060811 | 2.472973 | 2.398649 | 2.966216 | 2.804054 | 5.216216 |
| std | 0.575396 | 0.817509 | 0.498757 | 0.381836 | 0.212995 | 0.430498 | 0.501652 | 0.251855 | 0.458621 | 0.313557 | 0.836627 | 0.568323 | 0.868305 | 0.761834 | 2.171368 |
| min | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 25% | 2.000000 | 2.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 | 4.000000 |
| 50% | 2.000000 | 3.000000 | 2.000000 | 1.000000 | 1.000000 | 1.000000 | 2.000000 | 1.000000 | 2.000000 | 1.000000 | 2.000000 | 2.000000 | 3.000000 | 3.000000 | 5.000000 |
| 75% | 3.000000 | 3.000000 | 2.000000 | 1.000000 | 1.000000 | 1.000000 | 2.000000 | 1.000000 | 2.000000 | 1.000000 | 3.000000 | 3.000000 | 4.000000 | 3.000000 | 8.000000 |
| max | 4.000000 | 4.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 | 2.000000 | 3.000000 | 4.000000 | 3.000000 | 4.000000 | 4.000000 | 8.000000 |

| | |
|---|---|
| Univariate Analysis |  |
| Bivariate Analysis |  |

| | |
|---|---|
| Multivariate Analysis | Some extension of Receiver operating characteristic to multi-class<br><br>ROC curve of class 0 (area = nan)<br>ROC curve of class 1 (area = 0.77)<br>ROC curve of class 2 (area = 0.80)<br>ROC curve of class 3 (area = 0.50) |

| | |
|---|---|
| Outliers and Anomalies | - |

## Data Preprocessing Code Screenshots

| | |
|---|---|
| Loading Data |  |

```
1 df=pd.read_csv("https://archive.ics.uci.edu/ml/machine-learning-databases/lymphography/lymphography.data",names=col_names)
2 print("Size of dataset:",df.shape)
3 df.head()
```

Size of dataset: (148, 19)

| | class | lymphatics | block of affere | bl. of lymph. c | bl. of lymph. s | by pass | extravasates | regeneration of | early uptake in | lym.nodes dimin | lym.nodes enlar | changes in lym. | defect in node | changes in node | changes in stru | special forms | dislocation of | exclusion of no | no. of nodes in |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | 4 | 2 | 1 | 1 | 1 | 1 | 1 | 2 | 1 | 2 | 2 | 2 | 4 | 8 | 1 | 1 | 2 | 2 |
| 1 | 2 | 3 | 2 | 1 | 1 | 2 | 2 | 1 | 2 | 1 | 3 | 3 | 2 | 3 | 4 | 2 | 2 | 2 | 2 |
| 2 | 3 | 3 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 4 | 3 | 3 | 4 | 8 | 3 | 2 | 2 | 7 |
| 3 | 3 | 3 | 1 | 1 | 1 | 1 | 2 | 1 | 2 | 1 | 3 | 3 | 4 | 4 | 4 | 3 | 1 | 2 | 6 |
| 4 | 2 | 3 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 4 | 3 | 5 | 1 | 2 | 2 | 1 |

| | |
|---|---|
| Handling Outliers | ```python
for col in df.columns:
    q1 = np.quantile(df[col],0.25)
    q3 = np.quantile(df[col],0.75)
    iqr = q3-q1
    lower_bound = q1 - (1.5*iqr)
    upper_bound = q3 + (1.5*iqr)
    df[col] = np.where(df[col]> upper_bound,upper_bound,df[col])
    df[col] = np.where(df[col]< lower_bound,lower_bound,df[col])
    sns.boxplot(df[col])
    print("")
    plt.show()
``` |
| Training and Testing | ```python
# Assuming 'class' is your target variable and the rest are features
y = df['class']  # Create y to hold your target variable
x = df.drop('class', axis=1)  # Create x to hold your features
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=42)
```
```python
###check shape to make sure it is all in order
print("size of x_train: {} \t size of x_test: {} \nsize of y_train:{} \t sixe of y_test: {}".format(x_train.shape,x_test.shape,y_train.shape,
``` |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | - |