# Data Collection and Preprocessing Phase

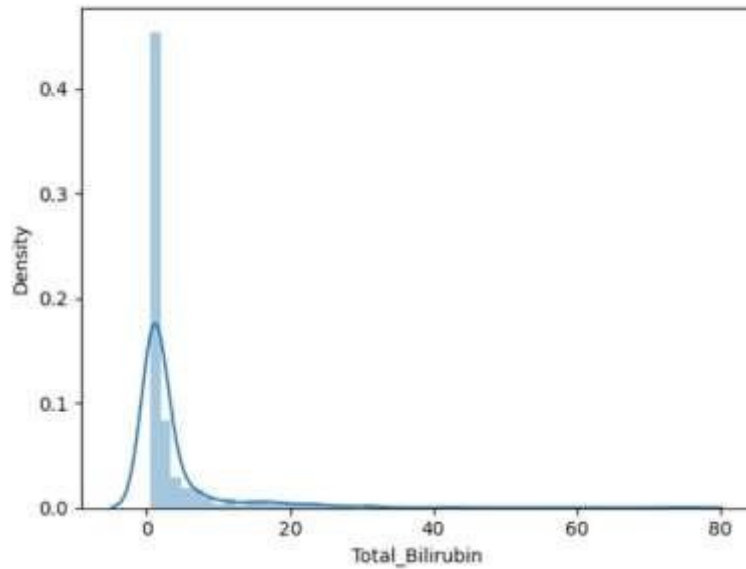| | |
|---|---|
| Date | 12 july 2024 |
| Team ID | 739651 |
| Project Title | Prediction and Analysis of Liver Patient Data Using Machine Learning |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

Identifies data sources, assesses quality issues like missing values and duplicates,
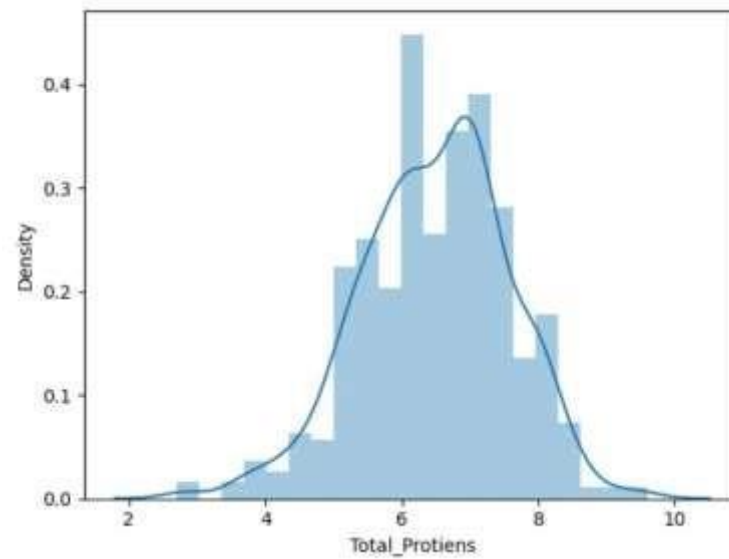and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
|---|---|
| Data Overview | 583 rows × 11 columns  |

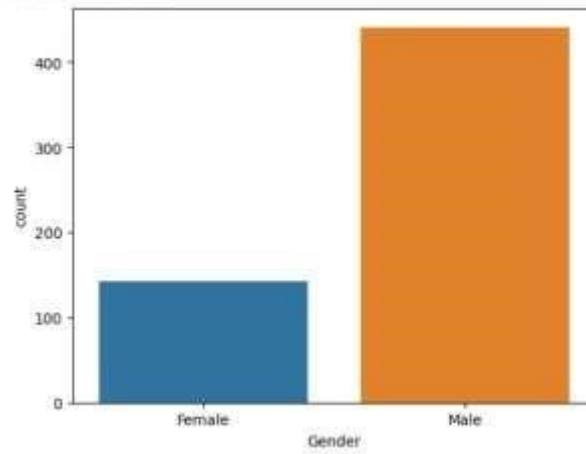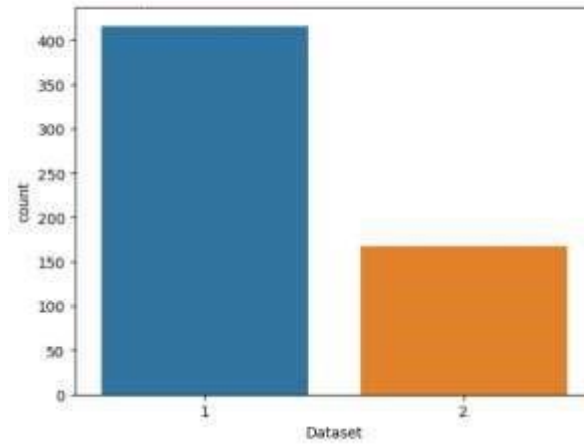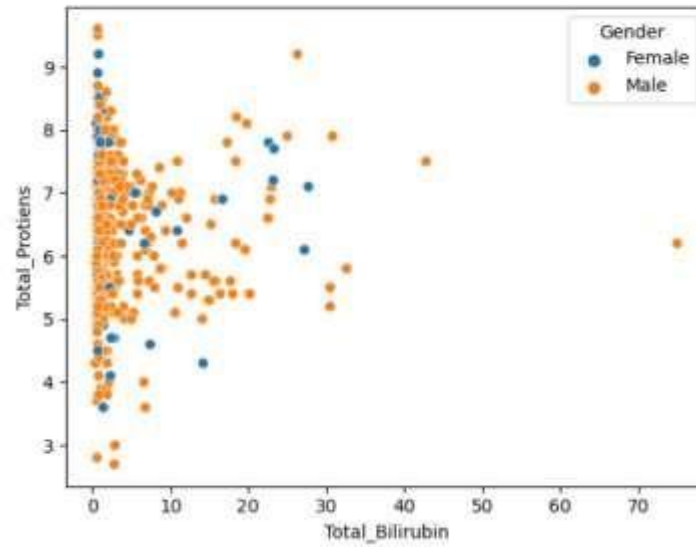| Univariate Analysis | <Axes: xlabel='Total_Bilirubin', ylabel='Density'><br><br><Axes: xlabel='Total_Protiens', ylabel='Density'><br> |

| | |
|---|---|
| Bivariate Analysis | <br>No. of Males: 441<br>No. of Females: 142<br><br><br>Liver disease patients: 416<br>Non-Liver disease patients: 167 |

| Multivariate Analysis |  |

<Axes: >



**Outliers and Anomalies**

```
sns.boxplot(data.Albumin_and_Globulin_Ratio,orient='h')
```

<Axes: >



## Data Preprocessing Code Screenshots

**Loading Data**

```
# loading the dataset
data = pd.read_csv("indian_liver_patient.csv")
```

| | |
|---|---|
| |  |
| Handling Missing Data |  |
| Data Transformation |  |
| Feature Engineering |  |

| Save Processed Data | ```python
import pickle
pickle.dump(svm , open('model.pkl','wb'))
pickle.dump(sc , open('sc.pkl','wb'))
``` |
| --- | --- |