

## Data Collection and Preprocessing Phase

|               |                                                        |
|---------------|--------------------------------------------------------|
| Date          | 16July 2024                                            |
| Team ID       | SWTID1720074204                                        |
| Project Title | prediction and analysis of liver patient data using ml |
| 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          | Basic statistics, dimensions, and structure of the data.                                                                                                                                                                                                                                                                                               |
| Univariate Analysis    | Median:<br>0.93<br>Mean:<br><div><div>44.7461413.2987991.486106290.57632980.713551109.9108066</div></div>                                                                                                                                                                                                                                              |
| Bivariate Analysis     | Histogram<br>countplot                                                                                                                                                                                                                                                                                                                                 |
| Multivariate Analysis  | Usually we drop that feature which has above 0.85% multicollinearity between two independent feature. Here we have only 'Total_Bilirubin' and 'Direct_Bilirubin' feature which has 0.87% mutlicollinearity. So we drop one of the feature from them and other independent feature has less multicollinearity, less than 0.80% So we keep that feature. |
| Outliers and Anomalies | There is no need of Standardization and Normalization of our dataset, as we using Ensemble Technique.                                                                                                                                                                                                                                                  |

|                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
|--------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                                            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
| <b>Data Preprocessing Code Screenshots</b> |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
| Loading Data                               | <pre> # Reading Dataset: dataset = pd.read_csv("Dataset/Liver_data.csv") # Top 5 records: dataset.head()</pre>                                                                                                                                                                                                                                                                                                                                                                         |
| Handling Missing Data                      | <pre> # Cheaking Missing (NaN) Values: dataset.isnull().sum()  [ ] # Filling NaN Values of "Albumin_and_Globulin_Ratio" feature with Median : dataset['Albumin_and_Globulin_Ratio'] = dataset['Albumin_and_Globulin_Ratio'].fillna(dataset['Albumin_and_Globulin_Ratio'].median())</pre>                                                                                                                                                                                               |
| Data Transformation                        | There is no need of Standardization and Normalization of our dataset, as we using Ensemble Technique.                                                                                                                                                                                                                                                                                                                                                                                  |
| Feature Engineering                        | <pre> # Target feature: print("Liver Disease Patients :", dataset['Dataset'].value_counts()[1]) print("Non Liver Disease Patients :", dataset['Dataset'].value_counts()[2])  # Visualization: sns.countplot(dataset['Dataset']) plt.show()  # Target feature: print("Liver Disease Patients :", dataset['Dataset'].value_counts()[1]) print("Non Liver Disease Patients :", dataset['Dataset'].value_counts()[2])  # Visualization: sns.countplot(dataset['Dataset']) plt.show()</pre> |
| Save Processed Data                        | <pre> [ ] # Creating a pickle file for the classifier import pickle filename = 'Liver.pkl' pickle.dump(RandomForestClassifier, open(filename, 'wb'))</pre>                                                                                                                                                                                                                                                                                                                             |