

## Data Collection and Preprocessing Phase

Date	20 <sup>th</sup> June 2024
Team ID	SWTID1720080161
Project Title	Revolutionizing Liver Care : Predicting Liver Cirrhosis Using Advanced Machine Learning Techniques
Maximum Marks	2 Marks

### Data Quality Report:

The Data Quality Report will summarize data quality issues from the selected source, including severity levels and resolution plans. It will aid in systematically identifying and rectifying data discrepancies.

Data Source	Data Quality Issue	Severity	Resolution Plan
Kaggle	<ul style="list-style-type: none"> <li>Too many NULL values</li> </ul>	Moderate	Filling the Numerical Columns with Mean and removing abnormalities from categorical column.
Kaggle	<ul style="list-style-type: none"> <li>Wrong Data Type</li> </ul>	Moderate	<p>Changing the data type by type casting. Ex-</p> <pre>df["TG"] = df["TG"].astype(float) df["LDL"] = df["LDL"].astype(float) df["Total Bilirubin (mg/dl)"] = df["Total Bilirubin (mg/dl)"].astype(float)</pre>

Kaggle	<ul style="list-style-type: none"> <li>Ambiguous string entries in multiple column</li> </ul>	Low	<p>Dropping the below rows</p> <pre>df = df[df['TG'] != '130LDL'] df = df[df['LDL'] != 'HDL'] df = df[df['Total Bilirubin (mg/dl)'] != 'o.4']</pre>
Kaggle	<ul style="list-style-type: none"> <li>Highly imbalanced outcome</li> </ul>	Moderate	<ul style="list-style-type: none"> <li>Synthetically generating 300 rows with minority class</li> </ul> <pre># Select row where Outcome is "NO" df_no_outcome = df[df["Predicted Value(Out Come Patient suffering from liver cirrhosis or not)"] == "no"]  # Check if df_no_outcome is empty and handle it accordingly if df_no_outcome.empty:     print("No rows found with 'NO' outcome. Cannot generate synthetic data.") else:     # Generate synthetic rows based on existing "NO" outcome rows     synthetic_rows = []     num_synthetic_samples = 300 # Number of synthetic samples to generate      for _ in range(num_synthetic_samples):         # Randomly select a row from df_no_outcome         selected_row = df_no_outcome.sample(n=1).iloc[0]          # Create a new row with random values based on the selected row         synthetic_row = selected_row.copy()          # Iterate over df_no_outcome.columns:         for column in df_no_outcome.columns:             if df_no_outcome[column].dtype in ["int64", "float64"]:                 # Add noise within a 10% range of the selected row's value                 noise_factor = 0.1 * selected_row[column]                 synthetic_row[column] = random.uniform(selected_row[column] - noise_factor, selected_row[column] + noise_factor)             elif df_no_outcome[column].dtype == "object":                 synthetic_row[column] = random.choice(df_no_outcome[column].unique())          synthetic_row.append(synthetic_row)      # Convert synthetic rows to a DataFrame and append to the original DataFrame     synthetic_df = pd.DataFrame(synthetic_row, columns=df_no_outcome.columns)     df = pd.concat([df, synthetic_df], ignore_index=True)</pre> <ul style="list-style-type: none"> <li>Checking the value counts.</li> </ul> <pre>df["Outcome"].value_counts()</pre> <pre>Outcome yes    859 no      18 Name: count, dtype: int64</pre> <ul style="list-style-type: none"> <li>Filling the null values with MODE.</li> </ul> <pre>df["Outcome"].fillna("yes", inplace=True)</pre>
Kaggle	<ul style="list-style-type: none"> <li>Many unimportant Features</li> </ul>	Moderate	<ul style="list-style-type: none"> <li>Removing the features by analyzing the importance scores.</li> </ul>

```
from sklearn.ensemble import RandomForestClassifier

model = RandomForestClassifier(n_estimators=100)
model.fit(X, y)
importances = model.feature_importances_

# Print feature importances
for feature, importance in zip(X.columns, importances):
    print(f'{feature}: {importance:.4f}')
```

```
Age: 0.0006
Gender: 0.0000
Duration of alcohol consumption(years): 0.1940
Quantity of alcohol consumption (quarters/day): 0.0206
Type of alcohol consumed: 0.0000
Hepatitis B infection: 0.0000
Hepatitis C infection: 0.0000
Diabetes Result: 0.0044
Blood pressure (mmhg): 0.0001
Obesity: 0.0000
Family history of cirrhosis/ hereditary: 0.0001
TCH: 0.0001
TG: 0.0001
LDL: 0.0002
HDL: 0.0003
Hemoglobin (g/dl): 0.0011
PCV (%): 0.0007
RBC (million cells/microliter): 0.0282
MCV (femtoliters/cell): 0.0007
MCH (picograms/cell): 0.0194
MCHC (grams/deciliter): 0.0534
Total Count: 0.0010
Polymorphs (%) : 0.0104
Lymphocytes (%): 0.0058
Monocytes (%): 0.0025
Eosinophils (%): 0.0000
Basophils (%): 0.0074
Platelet Count (lakhs/mm): 0.0203
Total Bilirubin (mg/dl): 0.1604
Direct (mg/dl): 0.1125
Indirect (mg/dl): 0.0092
Total Protein (g/dl): 0.0024
Albumin (g/dl): 0.0800
Globulin (g/dl): 0.0003
A/G Ratio: 0.0518
AL.Phosphatase (U/L): 0.0204
SGOT/AST (U/L): 0.0199
SGPT/ALT (U/L): 0.0114
USG Abdomen (diffuse liver or not): 0.1605
```

- Dropping all the unnecessary columns.

```
for col in drop_col:
    if col in X.columns:
        X.drop(columns=[col],inplace=True)
```