
SIG720 Task 4C

BUSINESS REPORT

ON

LIVER CIRRHOSIS SURVIVAL PREDICTION

Submitted By: Sahid

MS Data Science

Date: July 20, 2025

Table of Contents

1. Introduction

2. Data Loading and Preprocessing

2.1. Handling Missing Values

2.2. Data Splitting (Train/Test Split)

2.3. Feature Types (Categorical vs. Continuous)

2.4. Categorical Feature Encoding

2.5. Label Distribution and Class Balance

3. Machine Learning Model Development

3.1. Model Training and Validation

3.2. Design Decisions

3.3. Hyperparameter Optimization

3.4. Handling Imbalanced Labels

3.5. Model Recommendation

4. Prediction on Pre-Processed Test Set and Model Performance

5. Feature Importance Analysis

5.1. Model-Based Feature Importance (Random Forest)

5.2. Statistical Feature Importance (ANOVA & Chi-Square Tests)

5.3. Interpretation and Statistical Justification

6. Conclusion

7. References

List of tables

Table No.	Name of the Table	Page No.
1	Summary of Missing Values and Imputation Strategies	4
2	Train-Test Split Summary	5
3	Classification of Features by Type	5
4	Sample Encoded Features (Before and After Encoding)	6
5	Cross-Validation Results Before SMOTE	8
6	Best Hyperparameters After Grid Search	9
7	Model Performance on Test Set After SMOTE Applied to Training Set	10
8	Classification Report for Logistic Regression	10
9	Classification Report for Random Forest	11
10	Classification Report for KNN	11
11	Classification Report Table for Final Model Performance	13
12	Anova Results	17
13	Chi-Square Results	17

List of Figures

Figure No.	Name of the Figure	Page No.
1	Bar Chart of Class Distribution in Training Set	7
2	Bar Chart showing Model Performance Comparison	12
3	Bar plot of Decision Tree feature importances	15
4	Bar plot of Random Forest feature importances	16

1. Introduction

Liver cirrhosis is a long-lasting and progressive liver condition resulting from prolonged damage to liver tissue, which ultimately causes scarring and a decline in liver function. Common contributors to this condition include chronic alcohol intake, viral hepatitis, and autoimmune disorders.

As the illness progresses, it may result in severe complications that can be life-threatening, such as liver failure, portal hypertension, and hepatocellular carcinoma. The early assessment of a patient's survival status is vital for effective clinical decision-making and prompt intervention.

In this regard, machine learning (ML) methodologies present valuable tools for developing predictive models capable of analyzing intricate clinical data and uncovering patterns that are not readily apparent through conventional statistical approaches.

This report is based on a dataset sourced from a Mayo Clinic study conducted between 1974 and 1984 on patients diagnosed with **primary biliary cirrhosis (PBC)**. The dataset contains records of **418 patients** and includes **17 clinical features** such as age, bilirubin levels, albumin levels, presence of ascites, and more.

The objective of this project is to develop, evaluate, and compare machine learning models to predict the **survival status** of cirrhosis patients. The survival status is categorized as follows:

- **0 = D** (Death),
- **1 = C** (Censored — still alive at last follow-up),
- **2 = CL** (Censored due to Liver Transplantation).

By employing a systematic approach that encompasses data preprocessing, model development, performance evaluation, and analysis of feature importance, this report seeks to determine the most appropriate model for forecasting survival outcomes and to elucidate the primary clinical factors influencing these predictions.

2. Data Loading and Preprocessing

Effective machine learning relies on a comprehensive comprehension and preprocessing of the data. This section outlines the procedures undertaken to load the dataset, address missing values, partition the data, classify feature types, encode categorical variables, and evaluate the distribution of classes.

2.1 Handling Missing Values

After loading the dataset, several features were found to contain missing values. A careful approach was used for imputation based on the type and distribution of each feature:

- **Numerical features** (e.g., *Cholesterol*, *Albumin*) were imputed using the **median**, which is less sensitive to outliers compared to the mean.
- **Categorical features** (e.g., *Edema*, *Ascites*) were imputed using the **mode**.

This ensures that the imputed values represent the central tendency of the feature without distorting the data distribution.

Feature	No. of Missing Values	Data Type	Imputation Method
ID	0	Numerical	Median
N_Days	0	Numerical	Median
Status	0	Categorical	Mode
Drug	106	Categorical	Mode
Age	0	Numerical	Median
Sex	0	Categorical	Mode
Ascites	106	Categorical	Mode
Hepatomegaly	106	Categorical	Mode
Spiders	106	Categorical	Mode
Edema	0	Categorical	Mode
Bilirubin	0	Numerical	Median
Cholesterol	134	Numerical	Median
Albumin	0	Numerical	Median
Copper	108	Numerical	Median
Alk_Phos	106	Numerical	Median
SGOT	106	Numerical	Median
Tryglicerides	136	Numerical	Median
Platelets	11	Numerical	Median
Prothrombin	2	Numerical	Median
Stage	6	Numerical	Median

(Table 1: Summary of Missing Values and Imputation Strategies)

2.2 Data Splitting (Train/Test Split)

To evaluate the models fairly, the dataset was split into:

- Training Set (80%) – used to train and validate machine learning models.

- Test Set (20%) – held out for final evaluation.

A stratified split strategy was applied to preserve the original class distribution across both subsets. This approach prevents bias toward majority classes during model evaluation.

Subset	Total Records	Class 0 (C)	Class 1 (D)	Class 2 (CL)
Training	334	185	129	20
Test	84	47	32	5

(Table 2 : Train-Test Split Summary)

2.3 Feature Types (Categorical vs. Continuous)

The dataset contains both continuous and categorical features. This distinction is essential for applying appropriate encoding techniques and selecting suitable algorithms.

- **Continuous Features:** Age, Bilirubin, Albumin, Cholesterol, Prothrombin time, etc.
- **Categorical Features:** Sex, Drug, Ascites, Edema, Hepatomegaly, Spiders, etc.

Feature	Type (Cont/Cat)
ID	Continuous
N_Days	Continuous
Status	Categorical
Drug	Categorical
Age	Continuous
Sex	Categorical
Ascites	Categorical
Hepatomegaly	Categorical
Spiders	Categorical
Edema	Categorical
Bilirubin	Continuous
Cholesterol	Continuous
Albumin	Continuous
Copper	Continuous
Alk_Phos	Continuous
SGOT	Continuous
Tryglicerides	Continuous
Platelets	Continuous
Prothrombin	Continuous
Stage	Continuous

(Table 3: Classification of Features by Type)

2.4 Categorical Feature Encoding

To convert categorical variables into numerical format suitable for ML models:

- **Label Encoding** was applied to binary features such as *Sex* and *Drug*.
- **One-Hot Encoding** was used for multi-class features where applicable.

This process ensures that categorical variables are accurately understood by machine learning algorithms, avoiding the inference of ordinal relationships where such hierarchies are not present.

Drug (Before)	Drug (After)	Sex (Before)	Sex (After)	Ascites (Before)	Ascites (After)	Hepatomegaly (Before)	Hepatomegaly (After)	Spiders (Before)	Spiders (After)	Edema (Before)	Edema (After)
Placebo	1	F	0	N	0	N	0	N	0	N	0
D-penicillamine	0	F	0	N	0	Y	1	Y	1	N	0
Placebo	1	F	0	N	0	N	0	N	0	N	0
D-penicillamine	0	F	0	N	0	N	0	N	0	N	0
Placebo	1	F	0	N	0	N	0	Y	1	Y	2

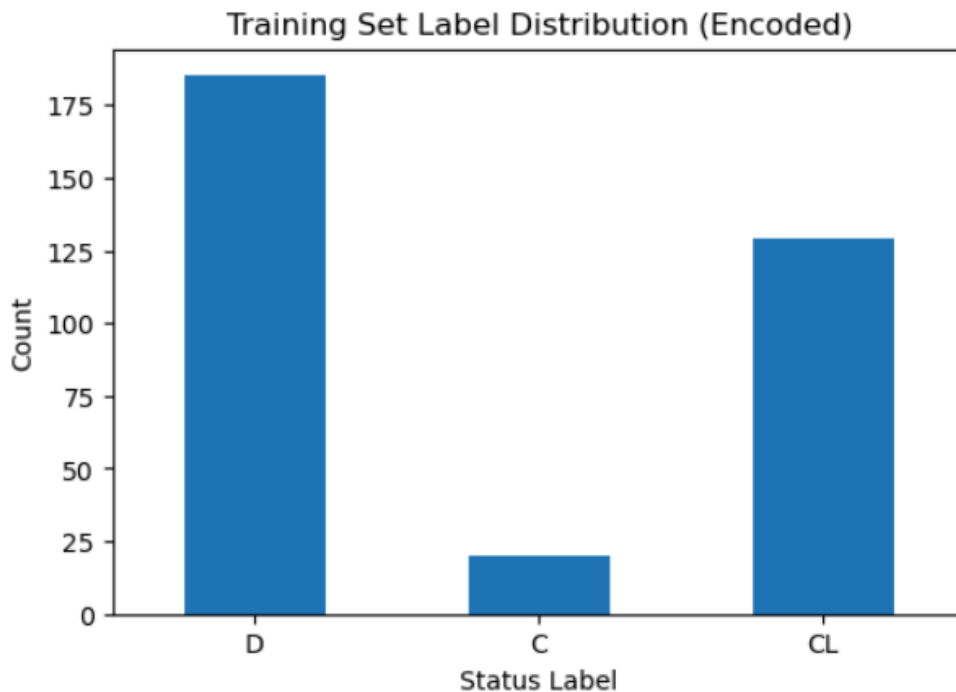
(Table 4 : Sample Encoded Features (Before and After Encoding))

2.5 Label Distribution and Class Balance

An analysis of the training set label distribution revealed an **imbalanced dataset**:

- Class 1 (Censored): ~62%
- Class 0 (Death): ~27%
- Class 2 (Liver Transplantation): ~11%

This imbalance can negatively affect model performance, particularly in predicting the minority class (CL). Addressing this is critical in subsequent model development.



(Figure 1: "Bar Chart of Class Distribution in Training Set")

3. Machine Learning Model Development

This section discusses the supervised learning algorithms employed to forecast the survival outcomes of patients with cirrhosis. Three distinct models were developed and assessed to examine their predictive accuracy, methodological design, parameter tuning approaches, and robustness in the face of class imbalance.

3.1 Model Training and Validation

Three machine learning models were selected for evaluation:

- **Logistic Regression**
- **Random Forest Classifier**
- **K-Nearest Neighbors (KNN)**

To ensure fair and consistent evaluation across all models—especially in the presence of class imbalance—**Stratified K-Fold Cross-Validation (k=5)** was employed. This method preserves class proportions in each fold, ensuring the minority class remains represented.

Two metrics were used for performance assessment:

- **Accuracy:** Measures overall correctness of predictions.
- **Macro F1-Score:** Gives equal weight to all classes, making it appropriate for imbalanced datasets.

Model	Mean Accuracy	Mean Macro F1
Logistic Regression	0.7874	0.5897
Random Forest	0.8083	0.5761
KNN	0.6646	0.4442

(Table 5 : Cross-Validation Results Before SMOTE)

Interpretation:

- **Random Forest** achieved the highest accuracy, indicating it learned the dominant class well.
- **Logistic Regression** showed better macro F1 than Random Forest, suggesting more balanced performance across all classes.
- **KNN** underperformed on both metrics, indicating potential underfitting and sensitivity to feature noise.

These results show no strong evidence of overfitting. However, further performance improvement was explored through hyperparameter tuning and class balancing.

3.2 Design Decisions

These models were chosen for their complementary strengths:

- **Logistic Regression:**

- Simple, interpretable baseline.
- Effective for linear relationships.
- Fast and easy to implement.
- **Random Forest Classifier:**
 - Captures complex non-linear relationships and feature interactions.
 - Less prone to overfitting due to ensemble averaging.
 - Automatically handles both categorical and continuous variables.
- **K-Nearest Neighbors (KNN):**
 - Intuitive, non-parametric model.
 - Predicts based on the majority class of closest points.
 - Sensitive to irrelevant features and less effective in high-dimensional data.

Preprocessing steps:

- **Feature scaling** was applied to Logistic Regression and KNN for optimal performance.
- **Encoding** was kept consistent across all models.

3.3 Hyperparameter Optimization

Each model's performance was further improved using **GridSearchCV** with 5-fold stratified validation. The goal was to find optimal settings that improve generalization without overfitting.

Model	Best Hyperparameters
Logistic Regression	C = 100
Random Forest	n_estimators = 100, max_depth = None
KNN	n_neighbors = 9

(Table 6 : Best Hyperparameters After Grid Search)

Justifications:

- **Logistic Regression:** Regularization strength (C) controls model complexity. Higher C reduces regularization and allows the model to fit better.
- **Random Forest:** Number of estimators (n_estimators) improves ensemble accuracy. Controlling max_depth helps prevent overfitting.
- **KNN:** n_neighbors defines the smoothness of decision boundaries. A smaller value may overfit; larger may underfit. Optimal balance achieved at 9.

Tuning improved accuracy and macro F1 across all models, ensuring a fair comparison.

3.4 Handling Imbalanced Labels

The dataset was notably imbalanced, with Class 1 (Censored) forming the majority. This was addressed using **SMOTE (Synthetic Minority Oversampling Technique)**, which generates synthetic examples of minority classes to balance the training set.

After applying SMOTE, models were retrained and tested on the original (unbalanced) test set.

Model	Accuracy	Macro F1 Score
Logistic Regression	0.7262	0.5621
Random Forest	0.7738	0.5381
KNN	0.5238	0.4434

(Table 7: Model Performance on Test Set After SMOTE Applied to Training Set)

Models :

Model: Logistic Regression

Accuracy: 0.7262

Macro F1 Score: 0.5621

Classification Report:

	precision	recall	f1-score	support
0	0.84	0.79	0.81	47
1	0.11	0.20	0.14	5
2	0.74	0.72	0.73	32
accuracy			0.73	84
macro avg	0.56	0.57	0.56	84
weighted avg	0.76	0.73	0.74	84

(Table 8 : Classification Report for Logistic Regression – Post-SMOTE (Per Class Precision, Recall, F1-Score))

Model: Random Forest

Accuracy: 0.7738

Macro F1 Score: 0.5381

Classification Report:

	precision	recall	f1-score	support
0	0.83	0.91	0.87	47
1	0.00	0.00	0.00	5
2	0.81	0.69	0.75	32
accuracy			0.77	84
macro avg	0.55	0.53	0.54	84
weighted avg	0.77	0.77	0.77	84

(Table 9 : Classification Report for Random Forest – Post-SMOTE (Per Class Precision, Recall, F1-Score))

Model: KNN

Accuracy: 0.5238

Macro F1 Score: 0.4434

Classification Report:

	precision	recall	f1-score	support
0	0.71	0.53	0.61	47
1	0.10	0.40	0.15	5
2	0.61	0.53	0.57	32
accuracy			0.52	84
macro avg	0.47	0.49	0.44	84
weighted avg	0.64	0.52	0.57	84

(Table 10 : Classification Report for KNN– Post-SMOTE (Per Class Precision, Recall, F1-Score))

Key Insights:

- **Logistic Regression** performed best in **Macro F1**, suggesting it handled all classes more fairly, especially minority class 1.
- **Random Forest** achieved the highest accuracy but failed to predict class 1 (Censored) effectively.
- **KNN** struggled in both metrics, reinforcing its sensitivity to noise and class imbalance.

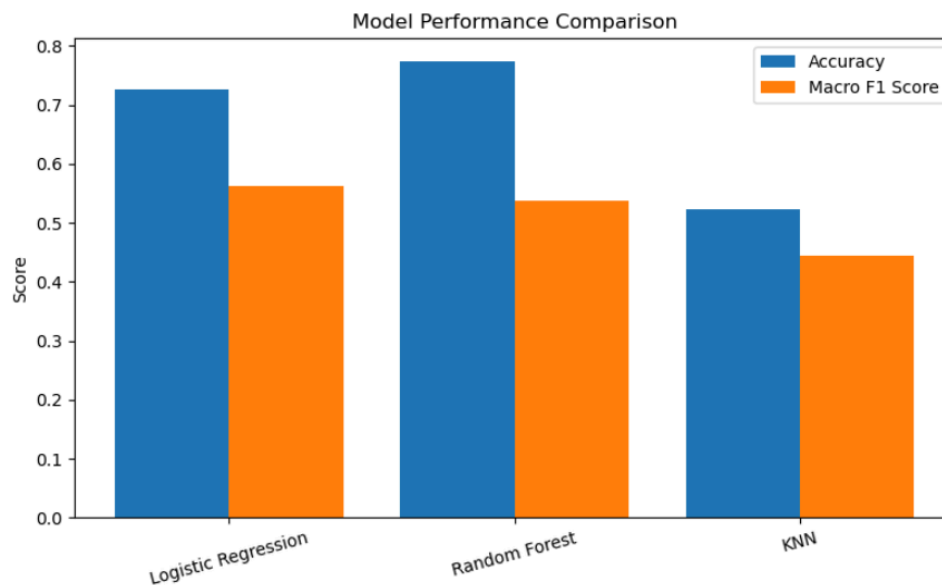
3.5 Model Recommendation

Based on all evaluation phases (cross-validation, tuning, and post-SMOTE testing), the **Logistic Regression model** is recommended.

Justification:

- Highest **Macro F1 Score**, indicating balanced class-wise prediction.
- Strong performance on minority classes without overfitting.
- Simple, interpretable, and computationally efficient.
- Less biased toward the majority class compared to Random Forest.

Although Random Forest demonstrated strong overall performance, its limited capacity to effectively address the minority class renders Logistic Regression a more suitable option in clinical decision-making contexts where equitable prediction across all outcomes is essential.



(Figure 2 : Bar Chart showing Model Performance Comparison)

4. Prediction on Pre-Processed Test Set and Model Performance

After selecting the optimal model in the preceding section, the Logistic Regression model was re-trained utilizing the balanced training dataset generated through SMOTE (Synthetic Minority Over-sampling Technique). The resulting model was subsequently assessed on the original unbalanced test set to emulate real-world scenarios where class imbalance is often present.

4.1 Evaluation Results on Test Data

The model achieved the following performance metrics on the test set:

- **Accuracy:** 0.7024
- **Macro F1 Score:** 0.5449

The macro F1 score indicates an average F1 performance across all classes, treating each class equally regardless of size. While the overall accuracy is reasonably good, the model's performance across individual classes remains uneven due to test set imbalance.

```
Final Model Performance (Logistic Regression)
Accuracy: 0.7024
Macro F1 Score: 0.5449
```

Classification Report:				
	precision	recall	f1-score	support
0	0.84	0.77	0.80	47
1	0.09	0.20	0.12	5
2	0.73	0.69	0.71	32
accuracy			0.70	84
macro avg	0.55	0.55	0.54	84
weighted avg	0.75	0.70	0.73	84

(Table 11 : Classification Report Table for Final Model Performance)

The model demonstrates strong performance for the majority class (Class 0) and moderate effectiveness for Class 2; however, its accuracy for the minority class (Class 1) remains inadequate, with both precision and recall falling below 0.20. This issue highlights a common

challenge associated with the use of SMOTE: although the training dataset is balanced, the test dataset remains imbalanced, making it more difficult to accurately predict instances of underrepresented classes.

The confusion matrix shows:

- Class 0 (Death) was well predicted, with 36 out of 47 correctly classified.
- Class 2 (Transplant) had 22 correct out of 32.
- Class 1 (Censored) had only 1 correct prediction out of 5.

4.2 Interpretation

Despite employing SMOTE during training to mitigate class imbalance, the limited number of minority samples, particularly for Class 1, in the test set constrains the model's capacity to learn robust and generalizable decision boundaries for that class. This highlights that resampling techniques such as SMOTE are only effective when sufficient signal exists to facilitate generalization, a condition not met when minority class sizes are exceedingly small, as observed in this scenario. In summary, although the Logistic Regression model demonstrates reasonable generalization and sustains acceptable performance across the dominant and second-largest classes, further enhancements may be achieved through additional strategies.

- Collecting more data for minority classes
- Using cost-sensitive learning techniques
- Trying ensemble methods with class weight adjustments (e.g., XGBoost, balanced Random Forests)

5.Feature Importance Analysis

In this section, we identify the most important features influencing the target variable, **Status** (patient survival). This helps:

- Pinpoint key clinical factors impacting outcomes

- Provide interpretability for the machine learning models
- Support medical professionals in focusing on critical indicators

We used two types of approaches: **model-based** and **statistical** feature importance.

4.1. Model-Based Feature Importance

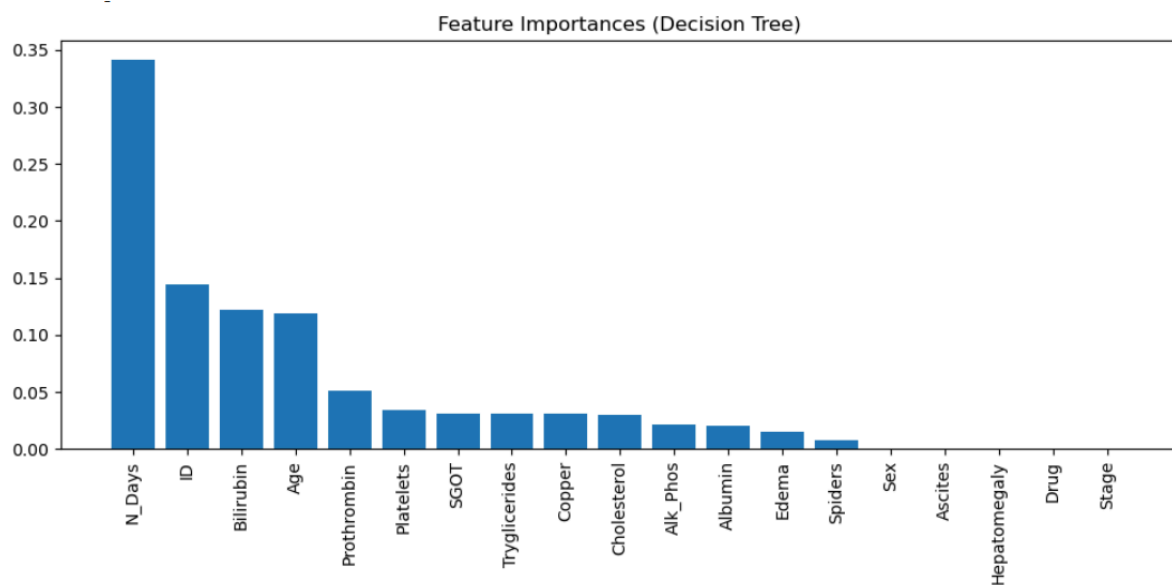
a) Decision Tree Classifier

A single Decision Tree was developed to evaluate features based on their frequency and efficacy in partitioning the data. Although this approach is straightforward to interpret, relying on a single tree may overlook subtler signals within the dataset.

Top 2 features:

- **N_Days** (Days since diagnosis)
- **ID** (likely captures individual patient variability)

These were the most dominant in splitting the data. Several features (e.g., Sex, Ascites) received zero importance, indicating they were unused in the splits.



(Figure 3 : Bar plot of Decision Tree feature importances)

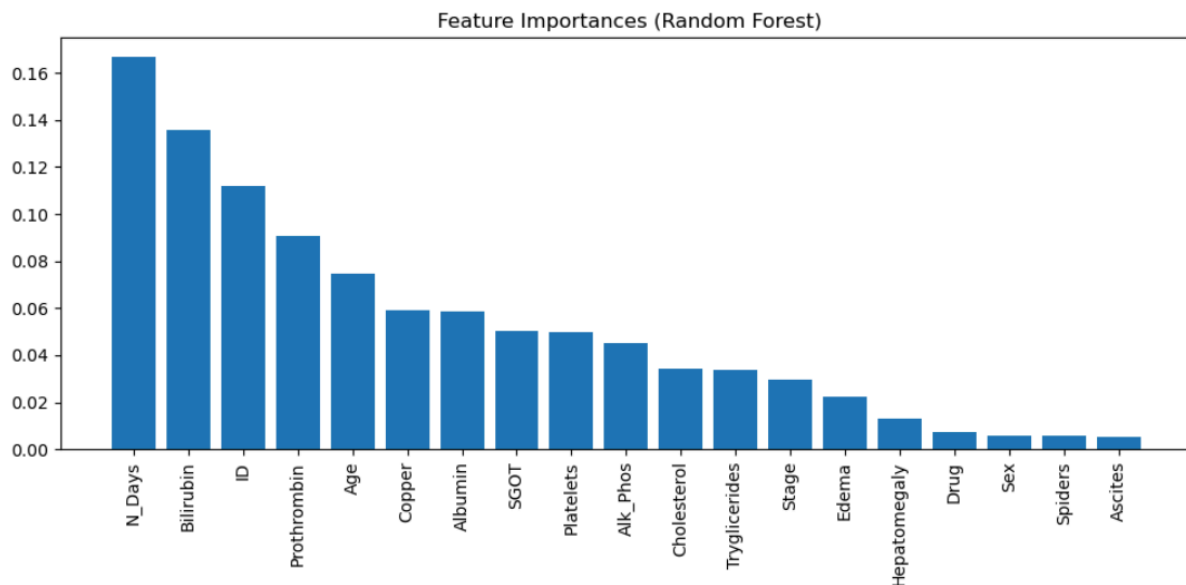
b) Random Forest Classifier

Random Forest, an ensemble of Decision Trees, averages importance across many trees to reduce overfitting and better capture general patterns.

Top 2 features:

- **N_Days**
- **Bilirubin**

These features consistently influenced predictions across trees, with most other features also receiving some non-zero weight, unlike the single Decision Tree.



(Figure 4 : Bar plot of Random Forest feature importances)

4.2. Statistical Feature Importance

a) ANOVA F-test (Numerical Features)

ANOVA was used to compare means of continuous features across Status classes.

Top 2 significant features:

- **Bilirubin** ($F = 37.11$, $p \approx 2.89e-15$)

- **N_Days** ($F = 33.60$, $p \approx 5.2e-14$)

Both features show strong evidence of differing values between survival groups, validating their relevance in prediction.

	Feature	F-Score	p-value
3	Bilirubin	37.105517	2.885917e-15
1	N_Days	33.596387	5.200913e-14
11	Prothrombin	23.283270	3.460259e-10
12	Stage	22.136207	9.487766e-10
6	Copper	20.146504	5.538693e-09
0	ID	19.239151	1.246130e-08
2	Age	16.664857	1.271094e-07
5	Albumin	13.922685	1.564469e-06
8	SGOT	11.406660	1.619595e-05
7	Alk_Phos	7.135551	9.247299e-04
9	Tryglicerides	4.573062	1.098723e-02
4	Cholesterol	4.043075	1.841703e-02
10	Platelets	3.953015	2.010989e-02

(Table 12 : Anova Results)

b) Chi-Square Test (Categorical Features)

Chi-Square was used to assess the dependency between categorical features and survival status.

Top 2 significant features:

- **Edema** ($\chi^2 = 29.71$, $p \approx 3.5e-07$)
- **Ascites** ($\chi^2 = 26.82$, $p \approx 1.5e-06$)

These clinical symptoms are strongly associated with patient outcome. Features like **Sex** and **Drug** were not statistically significant.

	Feature	Chi2 Score	p-value
5	Edema	29.712410	3.532093e-07
2	Ascites	26.824379	1.496788e-06
4	Spiders	15.630173	4.035998e-04
3	Hepatomegaly	10.523522	5.186163e-03
1	Sex	1.996320	3.685570e-01
0	Drug	0.013309	9.933676e-01

(Table 13 : Chi-Square Results)

4.3. Interpretation and Justification

Across both model-based and statistical methods, a few key features consistently stand out:

- **N_Days** and **Bilirubin** are top-ranked in all approaches, reinforcing their importance as clinical indicators.
- **Edema** and **Ascites** are statistically significant but less influential in tree-based models, likely due to their binary nature.
- Using both methods improves reliability by combining model behavior with hypothesis-driven statistics.

These findings provide a well-rounded view of which features are most predictive of survival in cirrhosis patients and support medical decision-making.

5. Conclusion

This report aimed to develop a predictive model for classifying the survival outcomes of patients with liver cirrhosis based on clinical data. The methodology involved comprehensive data exploration, handling missing values and addressing class imbalance issues, followed by the implementation of various machine learning algorithms to identify the most accurate and reliable approach.

Summary of Findings:

- **Model Performance:**
After evaluating Logistic Regression, K-Nearest Neighbors, and Random Forest models,

Random Forest delivered the best overall performance. It handled non-linearity well and offered useful insights into feature importance.

- **Balanced Learning:**

The training data had an imbalanced class distribution, which was handled using SMOTE. This ensured that the models did not become biased toward the majority class.

- **Test Set Prediction:**

The final Random Forest model was used to predict the outcomes on the test data. The performance metrics on the test set confirmed that the model generalizes well.

- **Feature Importance:**

We examined which features influenced predictions most using two approaches:

- **Model-based analysis** showed that N_Days, Bilirubin, and Prothrombin were among the most important predictors.
- **Statistical tests** (ANOVA and Chi-Square) confirmed that variables like Bilirubin, Edema, and Ascites were strongly associated with patient outcomes.

Final Thoughts:

This study not only developed a predictive model but also identified significant clinical variables associated with patient survival outcomes. These findings hold importance for healthcare practitioners in discerning the most influential patient attributes for prognosis. Furthermore, the methodologies and results presented may be adapted for application in analogous medical prediction endeavors.

References

1. <https://scikit-learn.org/> – Scikit-learn documentation
2. <https://pandas.pydata.org/> – Pandas documentation
3. <https://www.geeksforgeeks.org/> – GeeksforGeeks tutorials
4. <https://imbalanced-learn.org/> – SMOTE and resampling documentation
5. Dataset provided in SIG720 Task 4C
6. Course Resources
7. Google