

1. Introduction

We leverage machine learning techniques to predict outcomes for patients with liver cirrhosis based on clinical data. By analysing a dataset comprising various features related to liver disease and patient characteristics, this study identifies key predictors of cirrhosis outcomes and develops predictive models that can aid clinicians in decision making and management of the condition.

Our report explains our process for data pre-processing, exploratory analysis, model development, and evaluation, yielding useful models to predict cirrhosis outcomes to assist clinicians. Our modelling approach will also be applicable to larger datasets on liver health, so our pilot study reported here can be scaled up to improve performance further.

Liver cirrhosis is a significant global health burden, characterized by progressive deterioration of liver function due to long-standing liver injury. The disease progression is often silent and can culminate in life-threatening complications without timely intervention. Early and accurate prediction of cirrhosis outcomes can significantly enhance clinical decision-making and patient care. We identify which medical markers are predictive of liver cirrhosis outcomes so these can be gathered and used with our model to predict care needs and disease outcomes.

The objectives of this project were:

- a) Exploratory data analysis to uncover patterns and insights within the liver cirrhosis dataset.
- b) Develop and evaluate machine learning models, specifically focusing on Gradient Boosting and Random Forest, to predict liver cirrhosis outcomes.
- c) Identify key features influencing cirrhosis prognosis and assess the models' predictive performance through various metrics.
- d) Provide actionable recommendations for clinicians based on the analysis findings and suggest what future real-world data should be gathered.

2. Data Description and Pre-processing

Our dataset is synthetic^[1] and based on the Mayo Clinic Liver Cirrhosis Study (1974-84)^[2]. The data features include many medical measures and demographic variables for hypothetical patients with liver cirrhosis. Each record represents a patient's clinical profile and outcome. Our training data is complete and consists of data entries for 7905 patients across 19 medical features with their liver cirrhosis outcomes. Outcomes were measured after N days under observation, with three possible outcomes: '**C**' - **patient is still alive**, '**CL**' - **patient had liver transplant and is still alive**, or '**D**' - **patient died**. The data is imbalanced with 62.8% C, 3.5% CL, 33.7% D, making CL a minority class. A separate holdout test dataset consists of 5271 entries with outcomes removed. Our task is to predict the probability of each outcome class for each entry in this test dataset, and score our model's predictions against their true values. We developed our model entirely on the training dataset so that our final score against the test dataset will fairly measure how our model might perform on new patient data in a clinical setting.

2.1 Data Features Description: Each medical data feature indicates the patient's condition and are fed into our model training and analysed by our trained model to predict liver cirrhosis outcomes and thereby assist effective management of the disease. The dataset includes the following key features: Age, Sex, Bilirubin, Albumin, Cholesterol, Drug, Ascites, Hepatomegaly, Spiders, Edema, Copper, Alk Phos (Alkaline Phosphatase), SGOT, Triglycerides, Platelets, Prothrombin Time, and Stage. The target variable, "Status," categorizes patient outcomes into C,CL,D referring to Alive with Cirrhosis ('C'), Alive with Transplant ('CL'), or Dead ('D'). Transplant recipients who nonetheless die soon afterwards are assigned D, so CL is a label for successful transplant recipients. Appendix D explains each feature and its medical significance to liver cirrhosis.

2.2 Pre-processing Steps

1. **Data cleaning:** Initially, we confirmed our synthetic dataset was completed with no missing values or inconsistent data entries. Later, we simulated real world data through a deletion and imputation study that simulated missing data values and addressed these through appropriate imputation strategies, ensuring a complete dataset for analysis.
2. **Feature encoding:** Categorical variables such as Sex, Drug, Ascites, Hepatomegaly, Spiders, and Edema were encoded using one-hot encoding to convert them into a numerical format suitable for machine learning algorithms.
3. **Feature scaling:** Continuous features were standardised to have zero mean and unit variance. This scaling ensures that features with larger magnitudes do not unduly influence the model's performance.
4. **Data splitting:** Initially, we divided the dataset into training and validation sets following an 80/20 split. Stratified sampling was employed to maintain the proportion of each outcome class in both subsets, solving the risk of class imbalance being exacerbated in the split. Later, we will use cross-validation techniques to overcome possible overfitting to biases introduced by a given choice of 80/20 split.

Through these pre-processing steps, the dataset was transformed into a clean, well-structured format conducive to effective machine learning model development. These steps are crucial for ensuring that subsequent analyses are based on reliable and consistent data.

3. Exploratory Data Analysis

Prior to model development, we conducted an exploratory data analysis (EDA) to gain insights into the dataset's characteristics and uncover any inherent patterns or anomalies. This analysis helped us understand of the data's structure, distributions, and relationships between features.

Key findings:

- **Descriptive statistics:** Summary statistics provided an overview of the central tendency, dispersion, and shape of the dataset's distributions. Notably, variables such as Bilirubin and Albumin exhibited skewed distributions, indicating the presence of outliers or non-normal distribution patterns.
- **Correlation analysis:** A correlation matrix was generated to identify potential relationships between features. This analysis helped in understanding which variables might have collinear relationships, aiding in the feature selection process to avoid multicollinearity in the predictive models.
- **Class distribution:** The distribution of the target variable, 'Status', was examined, revealing a class imbalance with a higher prevalence of the 'Alive Cirrhosis (C)' outcome compared to 'Dead (D)' and 'Alive Transplant (CL)'. This insight guided the choice of stratified sampling and evaluation metrics sensitive to class imbalance.
- **Visual explorations:** Histograms, box plots, and scatter plots were used to visualise the distributions of individual features and their interactions. For instance, box plots highlighted the presence of outliers in features like Bilirubin, suggesting potential issues with extreme values influencing model performance.
- **Feature explorations:** Through scatter plots and pair plots, relationships between key features were explored. Notable observations included the potential relationships between Bilirubin levels and liver disease outcomes, reinforcing clinical understanding that higher bilirubin levels might be associated with worse prognosis.

The exploratory data analysis guided our subsequent data pre-processing and model development stages. We are confident our feature handling is careful enough and encouraged that our feature variables will predict liver cirrhosis outcomes, laying a solid foundation for our predictive modelling.

4. Imputation methods

For our study we chose median, K Nearest Neighbors (KNN) and multiple imputation of chained equations (MICE) as our methods. We used random deletions to remove 10% of the data values. Our data values were chosen by choosing randomly across the columns and rows. This was done to simulate the random incompleteness of real-world data where the diagnostic suite is not always fully available. KNN estimates the values of the missing data based on the similarity to its neighbors while the median imputation takes the median of the remaining values of the attribute in the dataset and estimates it as the values for the missing data. MICE replaces the missing values with multiple imputations using a regression model iteratively until the convergence criteria is met. Mean-Squared Error (MSE) was chosen as a metric to compare the results of the different imputations and MICE turned out to be the best amongst the three methods with an average MSE of 0.000901. We also used XGBoost to compare the imputation methods and MICE again turned out to be the best amongst the three methods with a log loss of 0.510.

5. Tackling imbalanced class data

Our dataset is imbalanced with only 3% CL compared to 63% C and 34 % D. So, we address this class imbalance to avoid bias towards overrepresented classes C and D and neglect of the minority class CL. We considered two techniques, 'weight adjustment' and an oversampling technique called 'SMOTE'. Weight adjustment simply gives more weight to CL data in inverse proportion to its under-representation in the dataset. Whereas our oversampling technique SMOTE tried to balance the data by inferring additional CL data entries based on the existing data (Appendix E). SMOTE offered easy gains for the non-optimized models. Although in the case of the optimized model, we found that class weight adjustment could converge on a slightly better log loss in the end. The result for the fully optimized gradient boosting model + class weighting was a validation log loss of 0.46 compared to 0.53 with SMOTE.

6. Methodology

This section outlines the systematic approach employed in developing and evaluating predictive models for liver cirrhosis outcomes. The methodology encompasses model selection, training, evaluation, and validation processes, underpinned by a rigorous data analysis framework. For model selection, given the nature of the outcome variable (categorical with three levels), Gradient Boosting and Random Forest classifiers were chosen for their robustness and efficacy in handling multi-class classification problems. These models are well-suited for dealing with imbalanced datasets and offer interpretable insights into feature importance. Neural Networks were also explored.

For training, the dataset was partitioned into training (80%) and validation (20%) sets using stratified sampling to maintain outcome class proportions, addressing the issue of class imbalance. This partitioning ensures that the models are evaluated on unseen data, providing a fair estimate of their generalisation performance. Training data is then oversampled with SMOTE. The Gradient Boosting Classifier was trained using a range of hyperparameters to optimise performance, with a focus on learning rate, number of estimators, and tree depth. Hyperparameter tuning was conducted using grid search with cross-validation to identify the optimal configuration. Like Gradient Boosting, the Random Forest model underwent hyperparameter tuning with a focus on the number of estimators, max depth, and min samples split. The class weight parameter was adjusted to address class imbalance, enhancing the model's sensitivity to minority classes.

For evaluation, we favoured LogLoss in our development but also measured accuracy, precision, recall, and F1-score metrics for our models, offering a comprehensive view of their performance, particularly on the minority transplant CL. Given the imbalanced nature of the dataset, particular emphasis was placed on recall and F1-score to ensure that the models effectively identify the minority class CL.

Results after fine-tuning (Log-loss on validation set):

1. SMOTE + Random Forest : 0.489
2. SMOTE + Gradient Boosting : 0.525
3. ClassWeights + Gradient Boosting : 0.459
4. SMOTE + XGBoost : 0.576

Results of the same fine-tuned models on Kaggle Log-loss (Public, Private Scores):

1. SMOTE + Random Forest : 0.49, 0.50
2. SMOTE + Gradient Boosting : 0.49, 0.48
3. **ClassWeights + Gradient Boosting : 0.43, 0.42**
4. SMOTE + XGB : 0.53, 0.50

Conclusion: Basing our evaluation on our Kaggle result, we declare **ClassWeights+Gradient Boosting** as our best model.

To assess the models' stability and generalization capability, k-fold cross-validation was employed. This technique provided insights into the models' performance variance across different data subsets, reinforcing the robustness of the evaluation. Post-evaluation, an error analysis was conducted to identify and scrutinize instances where the models' predictions deviated from the actual outcomes. This analysis provided deeper insights into the models' limitations and potential areas for improvement. Our evaluation established a clear understanding of the models' capabilities and limitations, guiding future efforts to refine and deploy them effectively.

Results of cross-validations on the same fine-tuned models:

1. SMOTE + Random Forest : 0.332
2. SMOTE + Gradient Boosting : 0.272
3. ClassWeights + Gradient Boosting : 0.447
4. SMOTE + XGB : 0.265

However, we concluded that these results are implausible and symptomatic of overfitting to the training data, which explains why these models perform so poorly on the test dataset. That is, except for the ClassWeights + Gradient Boosting model, whose 0.45 log loss result here is plausible and consistent with its other results, including on the test dataset.

7. Model Development and Evaluation

This section explains how we developed and evaluated the Gradient Boosting and Random Forest models, detailing each step from initial training to the final evaluation.

For initial model training, the Gradient Boosting and Random Forest models were initially trained using the default hyperparameters. This preliminary training served as a baseline to assess the models' natural performance on the dataset without any tuning. For Gradient Boosting, the model demonstrated promising results, capturing complex patterns in the data. However, there was room for improvement, particularly on the minority class CL. We also wished to reduce overfitting and enhance its ability to generalise. Similarly, our simplest Random Forest model provided a robust baseline performance. Its inherent randomness and ensemble design offered a strong start, yet it was apparent that hyperparameter tuning could further optimise its performance.

Hyperparameter Tuning:

Hyperparameter tuning was conducted using a grid search approach, focusing on key parameters that influence model complexity and learning dynamics. For Gradient Boosting, parameters like ``n_estimators``, ``learning_rate``, and ``max_depth`` were tuned. The optimal configuration significantly improved the model's accuracy and reduced the likelihood of overfitting. In the case of Random Forest, tuning focused on ``n_estimators``, ``max_depth``, and ``min_samples_split``. Adjusting the ``class_weight`` parameter proved crucial in addressing the class imbalance, enhancing the model's performance across all classes. For XGBoost Classifier, tuning focused on controlling overfitting with ``n_estimators``, ``max_depth`` and ``subsample``. The scoring metric for GridSearchCV was log loss and we measured precision too. During development, accuracy and `neg_log_loss` had been tested as scoring method but both result in poorer log loss values for validation. K-fold cross-validation further validated the robustness of the tuned models, providing insights into their stability across different data segments. Both models showed consistent performance across folds, affirming their reliability and generalisation capability.

Feature Importance Analysis:

Feature importance analysis revealed critical predictors for each model, offering insights into the factors most influential in determining liver cirrhosis outcomes. In the Gradient Boosting model, features like Bilirubin (0.42), Days With Cirrhosis Within Trial (0.15), and Prothrombin (0.10) stood out as top predictors, aligning with clinical expectations regarding their significance in liver disease progression. The Random Forest model gave less importance to each feature, with top importance's for Bilirubin (0.14), Age (0.11,) Copper (0.08), and Days With Cirrhosis Within Trial (0.08). So the Random Forest model also underscored the relative importance of demographic factors, illustrating the multifaceted nature of cirrhosis prognosis.

Error Analysis:

Error analysis pinpointed specific instances where the models faltered, often in predicting less prevalent classes. Our investigation of these errors highlighted potential biases and areas where additional data or feature engineering might improve model accuracy, although we leave that as future work. Using the optimised Gradient Boost model on validation set, we analyzed the distribution of misclassifications (see Appendix B). Our evaluation established a clear understanding of the models' capabilities and limitations, guiding future efforts to refine and deploy them effectively.

8. Key Insights and Interpretations

This section synthesises the findings from the model development and evaluation stages, translating the technical results into interpretable insights and actionable intelligence for stakeholders.

- **Predictive power of clinical features:** The feature importance analysis shows the clinical relevance of certain features (e.g. Bilirubin, Albumin) in predicting cirrhosis outcomes. Their prominence in the models aligns with the prevailing medical understanding, which reinforces that our models are valid, and allows our models to offer data-driven affirmation to clinicians who use our models.
- **Model Performance and Real-World Application:** Both Gradient Boosting and Random Forest models demonstrated robust performance in predicting cirrhosis outcomes. Their ability to handle incomplete datasets, as often encountered in real-world scenarios, makes them particularly valuable tools for healthcare practitioners. When applied to new, incomplete patient data, our models can still provide reliable outcome predictions, aiding clinical decision-making.
- **Scalability with more data:** Our modelling approach is designed to scale effectively with additional data. Should clients provide access to larger, even if incomplete, historical datasets, the models can be retrained or fine-tuned to enhance their predictive accuracy further. This adaptability ensures that our models remain relevant and useful as more data become available, continually improving their predictive capabilities.

- **Implications with healthcare providers:** By integrating these predictive models into their workflow, healthcare providers can gain foresight into potential cirrhosis outcomes, enabling more personalised and timely interventions. This proactive approach could significantly improve patient management and outcomes in liver cirrhosis care.

9. Recommendation and Future Work

Recommendation:

- We recommend using our Gradient Boost model to predict patient outcomes in a clinical setting from any new real-world data gathered by those clinicians.
- We also recommend our approach to modelling and encourage training similar models on larger datasets to improve performance in future.

Future Work:

- **Data Collection:** We encourage continuous data collection for further development, especially for the underrepresented transplant outcome class and the most significant features, in order to reduce bias and improve model performance. Incomplete data should not be a deterrent; our models are designed to accommodate and make the best use of available information.
- **Model Refinement:** As more data become available, iterative retraining and refinement will allow our models adapt to new patterns and information, optimising their predictive accuracy.
- **Integration into Clinical Practice:** Clinicians will find our models most useful upon the development user-friendly interfaces for healthcare providers to leverage the predictive models efficiently, ensuring that the insights are accessible and actionable within the clinical workflow.
- **Future Research:** Further research should explore the integration of additional data types, such as genetic or lifestyle factors, into the models to enrich the predictive context and enhance outcome prediction.

10. Conclusion

Our models demonstrate significant potential for predicting liver cirrhosis outcomes in a clinical setting. Clinicians may gather new real-world data then use our model to predict liver cirrhosis outcomes for their patients to guide their care. The models' ability to discern critical patterns and provide accurate predictions, even with incomplete data, underscores their utility in a clinical setting. As we expand our datasets and refine our models, we anticipate even greater utility for liver disease management, ultimately advancing patient care and outcomes.

References

- [1] Kaggle. (n.d.). Synthetic Liver Cirrhosis Dataset, Kaggle: Playground Series S3E26. Retrieved March 11, 2024, from <https://www.kaggle.com/competitions/playground-series-s3e26/>
- [2] Arvidsson, J. (n.d.). Mayo Clinic Liver Cirrhosis Dataset 1974-84, Kaggle: Mayo Clinic. Retrieved March 12, 2024, from <https://www.kaggle.com/datasets/joebeachcapital/cirrhosis-patient-survival-prediction>

Appendices

Appendix A – Scores for our candidate models:

A.1. Kaggle scores for our candidate models:

Model	Private Score	Public Score	
XGB+SMOTE	0.53	0.53	
GB+SMOTE	0.49	0.48	
RF+SMOTE	0.49	0.50	
GB+WEIGHTED	0.43	0.42	<– Best

A.2. Classification Report scores for GB+Weighted on the Validation Set:

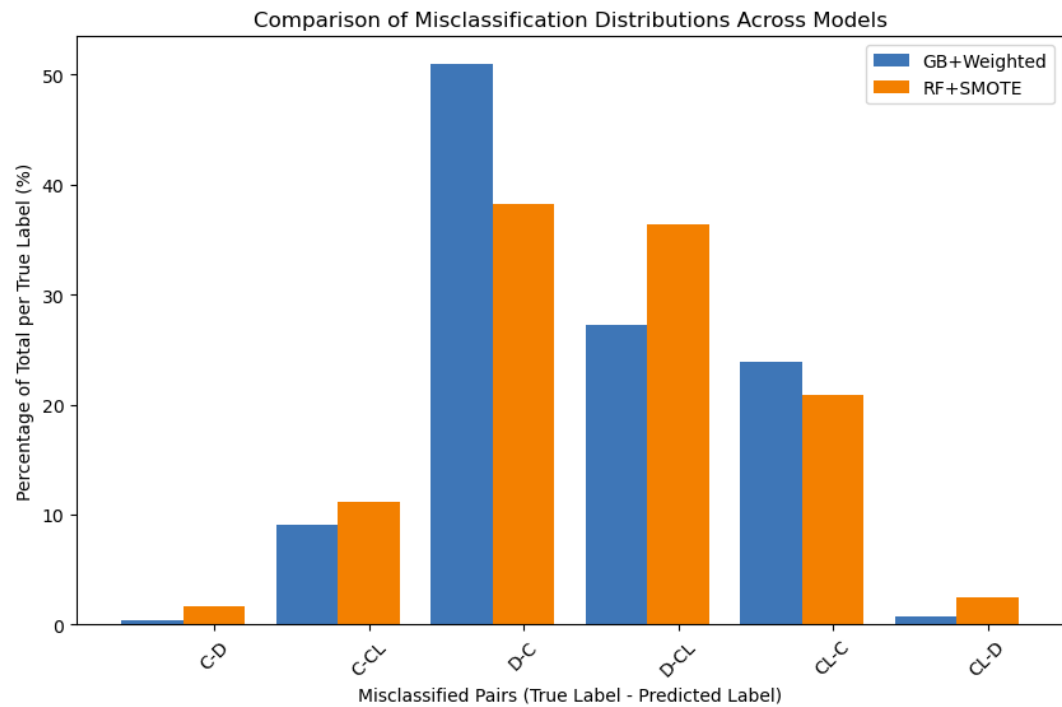
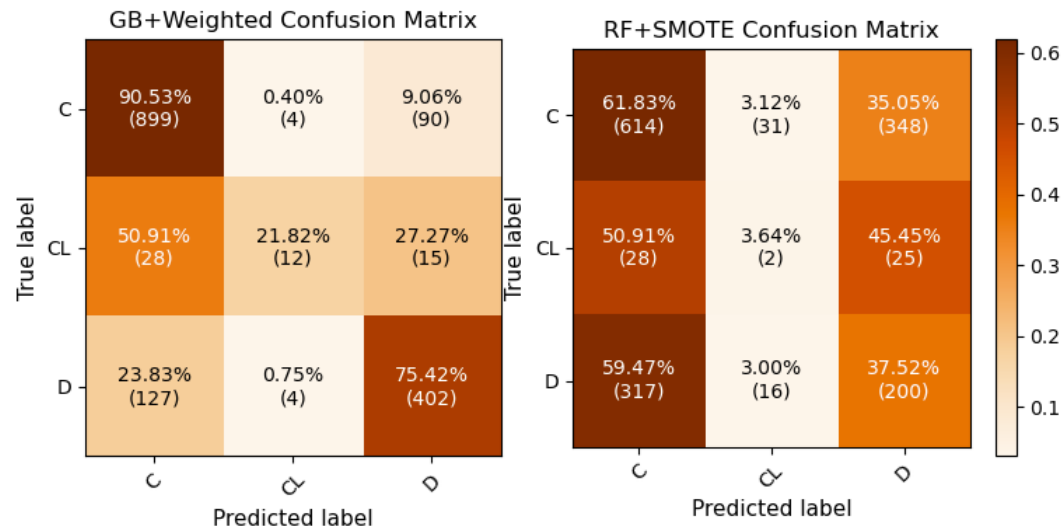
Class	Precision	Recall	F1-Score	Support
C	0.85	0.91	0.88	993
CL	0.60	0.22	0.32	55
D	0.79	0.75	0.77	533
Accuracy			0.83	1581
Macro Average	0.75	0.63	0.66	1581
Weighted Average	0.82	0.83	0.82	1581

Appendix B – Misclassification Study:

B.1. Confusion Matrices and Histogram of Misclassification Pairs

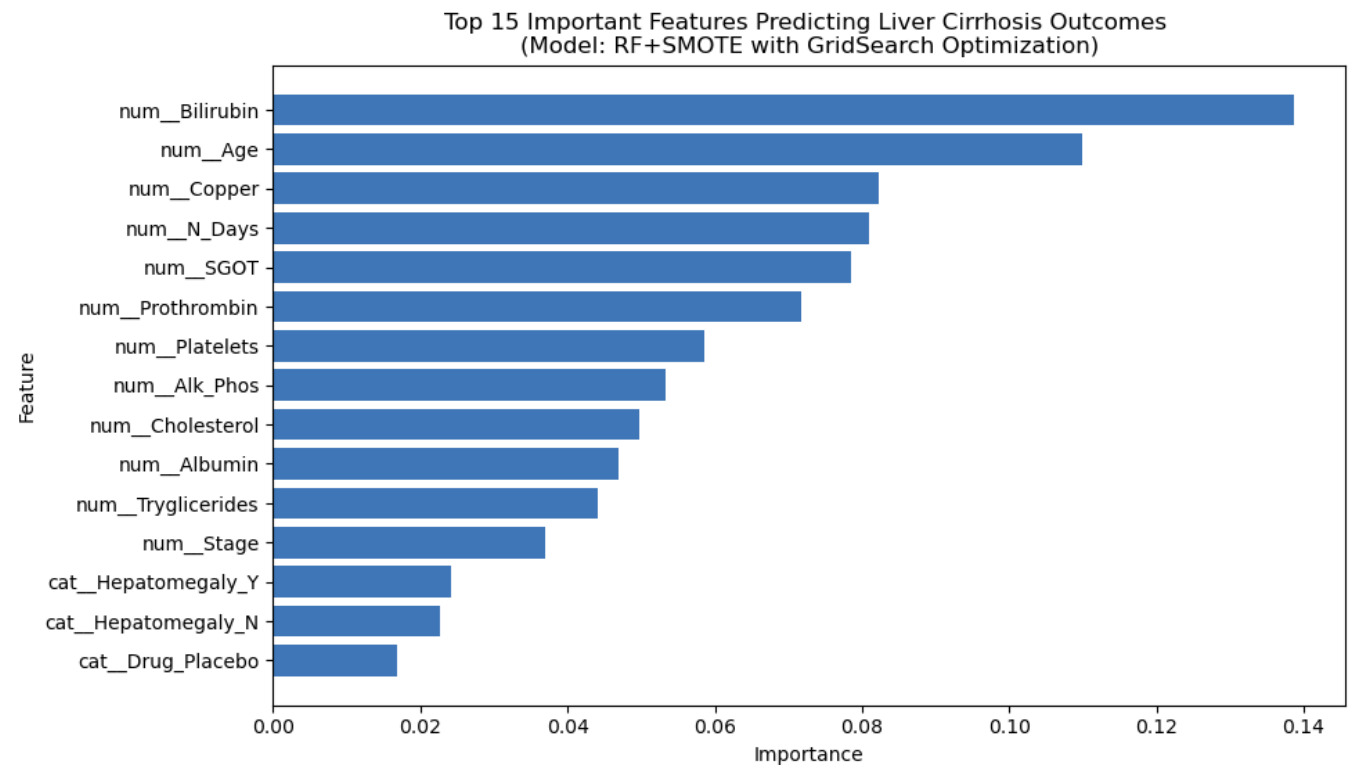
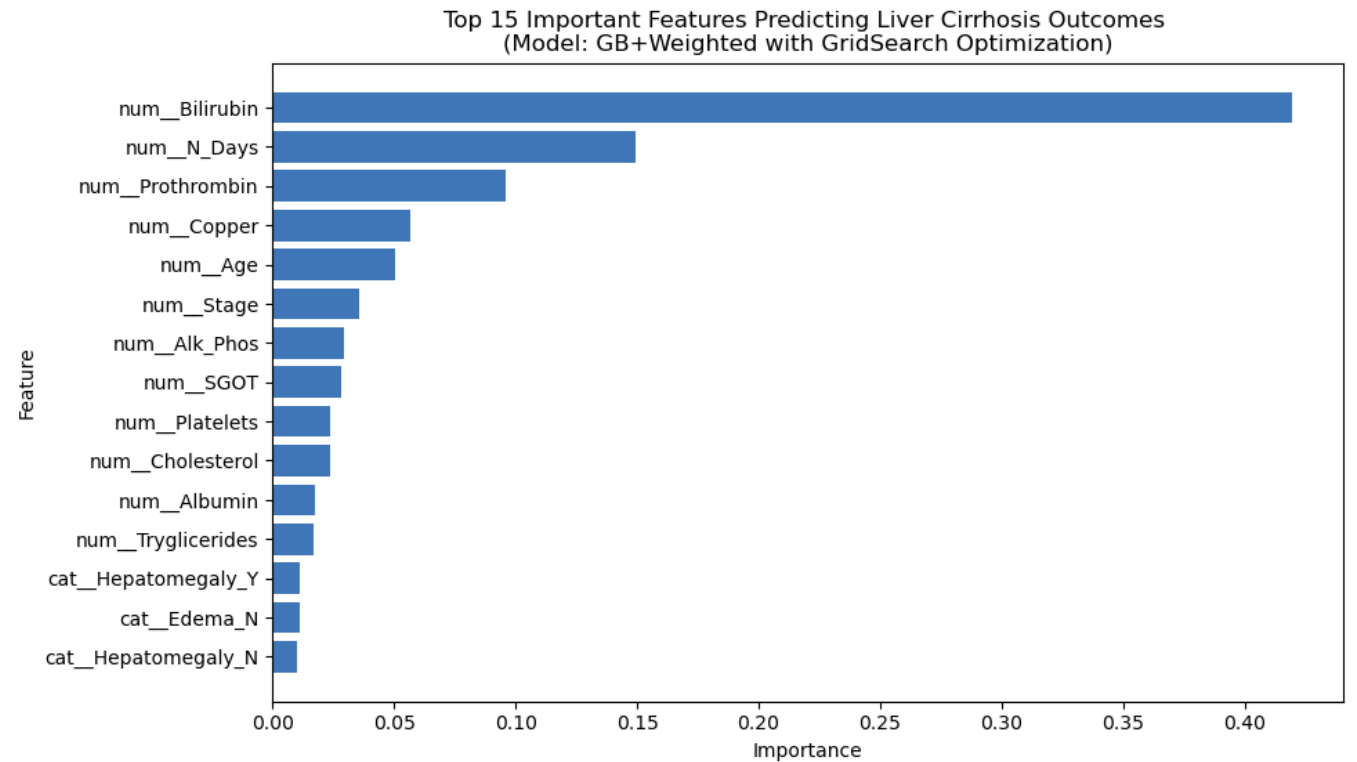
We prepared confusion matrices for both our best models to help us see whether and how they misclassify patients.

- 1. The GB+Weighted model looks good. Of its 20 liver transplant predictions, it suggests 12 out of 20 liver transplants should go to those whose True label is CL.
- 2. The RF+SMOTE model looks bad and little better than guessing. It suggests wasting 31 out of 49 liver transplants on those who would live anyway with True label C.



Appendix C – Feature Importance Study:

We ranked data features by their importance to each of our best two models. This is a measure of how useful data on each feature is to helping the model discriminate and decide between each class C, CL, D. We recommend prioritizing gathering the most important data in a clinical setting in order to allow our model to best assist clinicians with predictions of disease outcome.



Appendix D – Medical Significance of Clinical Data Features

N Days: A count of days suffering from cirrhosis for the patient in the study. Cirrhosis is associated with a chance of deterioration towards death. So `N Days` may predict a death outcome.

Age: The age of the patient can be a critical factor in the progression and management of liver cirrhosis. Older patients might have a different prognosis due to the natural aging process and potential comorbidities.

Sex: There might be sex-specific differences in the progression of liver diseases, with some studies suggesting variations in the susceptibility to liver injury and fibrosis between males and females.

Bilirubin: This is a by-product of red blood cell breakdown, and elevated levels can indicate liver dysfunction, often associated with worse outcomes in cirrhosis.

Albumin: A protein made by the liver, low albumin levels can indicate poor liver function and are associated with more advanced liver disease.

Cholesterol: While less directly related to cirrhosis, abnormal cholesterol levels can be associated with liver dysfunction and other metabolic conditions affecting the liver.

Drug (Treatment Type): The type of treatment can influence the progression of the disease. This could reflect different management strategies for cirrhosis and potentially different outcomes.

Ascites: The accumulation of fluid in the abdomen, indicating advanced liver disease, significantly impacts patient prognosis.

Hepatomegaly: Liver enlargement can be a sign of liver disease progression and is relevant in assessing the disease's severity.

Spiders (Spider Angiomata): These are small, spider-like blood vessels visible on the skin and can be an indicator of liver disease.

Edema: Swelling due to fluid accumulation, often in the lower legs, can indicate worsening liver function.

Copper: Elevated copper levels in the blood can be found in certain types of liver disease, such as Wilson's disease.

Alk Phos (Alkaline Phosphatase): An enzyme related to the bile ducts; high levels can indicate bile duct obstruction or other liver diseases.

SGOT (AST): An enzyme found in the liver, heart, and other tissues. Elevated levels can indicate liver damage.

Tryglicerides: While primarily a cardiovascular risk marker, abnormalities can also reflect metabolic issues related to liver disease.

Platelets: Liver cirrhosis can lead to thrombocytopenia (low platelet count), which is a critical factor in assessing liver function and disease severity.

Prothrombin: Clotting factor made by the liver. Liver damage can reduce levels, indicating liver dysfunction.

Stage: This reflects the severity of cirrhosis, with higher stages indicating more advanced disease and typically worse outcomes.

Status (Outcome): The target variable, indicating the patient's outcome - either alive with cirrhosis (C), alive with transplant (CL), or death (D), each of which the models aim to predict based on the features.

Appendix E – SMOTE Explanation

SMOTE (Synthetic Minority Over-sampling Technique) enhances dataset balance by generating synthetic samples from the minority class in an imbalanced dataset to avoid model bias toward the majority class. SMOTE can help achieve more robust and accurate predictive models, although one should take care to ensure the synthetic datapoints are medically plausible.

How it works:

1. Neighbor Identification: For each sample in the minority class, identify its k-nearest neighbors in the feature space.

Q: How would we sort the neighbors in feature space? There is more than one feature, so there's more than one feature axis, and the k nearest neighbors with respect to each feature are different, so which list should I draw the k nearest neighbors from?

A: Weighted Distance Metric: We weight the distance metric according to the importance of each feature. The Euclidean distance between two points in an n-dimensional space can be weighted with respect to each feature where weight reflects the importance of each feature. Then we sort the nearest neighbors according to this metric in order to respect the relative importance of different features. In the absence of a measure of feature importance, make an educated guess based on domain knowledge or preliminary analysis.

2. Synthetic Sample Creation: Create a synthetic data entry that lies between a chosen data entry belonging to a minority class and one of its k-nearest neighbors. For each data entry belonging to a minority class, select a random neighbor from its k-nearest neighbors. Imagine plotting both these data points and drawing a 'line' between these two points in 'feature space'. Generate a synthetic sample along the line segment connecting the minority class data entry and its selected neighbor. The synthetic sample point we generate will have feature values that are a linear combination of the two original data points, with the combination coefficients randomly chosen between 0 and 1. So the features of this synthetic data entry necessarily lie between those of our two original points.

3. Repetition: Continue this process until the dataset achieves the desired level of balance or until the required number of synthetic samples is generated.

Caution:

- SMOTE can introduce artificial bias, particularly if synthetic points do not well represent the underlying medical scenarios as hoped. One could scrutinize the medical plausibility of synthetic data.
- Validate on unmodified test dataset, to ensure improvements in performance metrics are genuine instead of overfitting to synthetic samples.