

# Liver Cirrhosis Prediction using Machine Learning Approach

Divya Sai Sree Chintala  
Group 12  
A20561001  
dchintala@hawk.iit.edu

## ABSTRACT

Liver Cirrhosis is a chronic liver disease that occurs when scar tissue replaces healthy liver tissue, preventing the liver from functioning normally. Liver cirrhosis, often undetectable in its initial stages due to the absence of symptoms, poses a significant challenge for early diagnosis and treatment. This challenge escalates as the condition advances, complicating further the process of diagnosing and treating the disease. This project introduces an artificial intelligence (AI) system, utilizing machine learning algorithms, designed to aid healthcare providers in the early detection of liver cirrhosis. With the goal of predicting the onset of this condition, this study has developed various machine learning algorithms, including Support Vector Machine, Decision Tree Classification, and Random Forest Classification. Among these, the Random Forest algorithm emerged as the most effective, achieving high accuracy rate. The open-access Liver Cirrhosis data dataset is used in the model development. The accuracy percentage of the models employed in this study are substantially greater than in earlier research, showing that the models utilized in this study are more dependable. Several model comparisons have shown their robustness, and the scheme may be determined from this analysis.

## I. INTRODUCTION

Liver cirrhosis is a widespread problem, especially in North America due to high alcohol intake. Our goal is to predict liver cirrhosis in patients based on their lifestyle and health conditions. In today's world, more than a million people are diagnosed with liver disease each year. Liver cirrhosis, hepatitis (A, B, C), and liver cancer are common liver diseases. Almost 71 million individuals worldwide are chronically unwell because of this. Although overall death is declining because of improvements in advanced treatment and maintaining a sound lifestyle. However, liver-related fatalities accounted for 3.5% of all deaths this century. According to research conducted by the WHO, 3–4 million individuals get infected with this each year. The global burden of this disease is substantial, with approximately 71 million people living with chronic infections and close to 399,000 fatalities reported in 2016 [1]. The World Health Organization identifies liver cirrhosis as a prevalent global health issue, with an annual infection rate of 3–4 million individuals. This virus is predominantly found in less

affluent regions of Asia and Africa, which bear a disproportionate burden compared to more developed countries in Europe and North America [2–4]. The insidious nature of liver cirrhosis is such that the majority of those infected, about 80%, show no early symptoms, often leading to advanced liver damage and increased mortality rates. Currently, no vaccine effectively combats the liver cirrhosis virus, underscoring the importance of accurate assessment of liver damage for proper clinical management and control of the disease's spread [5–7]. Technological advancements in artificial intelligence (AI) are enhancing the promptness and accuracy of disease diagnosis and management. Studies have demonstrated AI's ability to diagnose diseases with a proficiency comparable to, and occasionally surpassing, that of medical professionals, particularly those with less experience [8–12]. The lack of early symptoms often results in a late diagnosis, with 75 to 80 percent of infections advancing to severe stages before detection. Consequently, the disease may progress to the point where liver function is irreparably impaired [13].

## II. PROBLEM STATEMENT

Our paramount objective is to unearth insightful discoveries and trends for identifying links to Liver Cirrhosis. Moreover, we aim to develop a classification system to predict liver damage in individuals. This endeavor is crucial for early awareness and prompt action, empowering individuals to take necessary steps in mitigating the impact of liver-related conditions. This life-threatening disease is manageable if diagnosed in the early stages. Liver cirrhosis is often asymptomatic and as a result, diagnosing and treating patients during the early stages of illness is challenging. As the illness progresses to its latter stages, diagnosis and therapy become increasingly challenging. Here the main purpose is to build an intelligence model based on machine learning algorithms that may assist healthcare practitioners to forecast the possibility of a liver cirrhosis infection. Artificial intelligence (AI) advancements aid doctors in the rapid diagnosis and treatment of patients. There has been research comparing AI to human efficiency in illness diagnosis. A study found that AI was similarly capable of diagnosing as humans and, in fact, excelled at human efficiency when compared to less experienced physicians. Hence, here we are building a model that predict if a patient has Liver Cirrhosis or not using Machine Learning Algorithms.

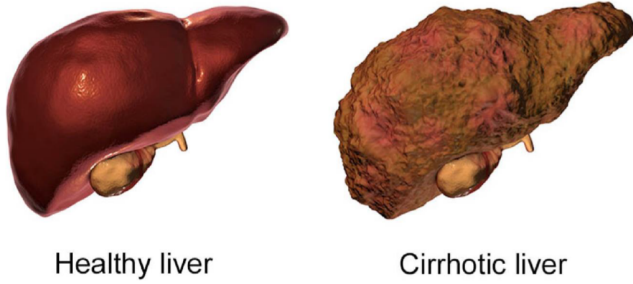


Figure 1: Cirrhotic Liver.

### III. RELATED WORK

Michigan Medicine cites work by Lok [14] indicating that timely intervention, ideally during early stages of the infection, can significantly improve treatment efficacy and outcomes. Lok also notes that symptoms typically manifest at more advanced stages, such as the development of cirrhosis, which escalates the risk of liver cancer, thereby emphasizing the critical need for early diagnostic measures [15]. AI-based disease diagnostics and prediction algorithms may aid in the early detection of acute infections and chronic disorders. Keltch, Lin, and Bayrak (2014) used four distinct kinds of AI algorithms on 424 liver cirrhosis patients' publicly accessible data. By comparing the findings of conventional serum indicators to the outcomes of biopsies, their suggested model aids in predicting the stage of fibrosis. Fresh methodologies and other AI techniques that might aid in the prediction of hepatitis B and C in millions of patients globally without the need for biopsies, thereby benefiting the whole healthcare system. Innovative AI diagnostic tools and predictive models have been explored for their potential in the early detection of liver cirrhosis and other related health conditions. For instance, Keltch, Lin, and Bayrak's study in 2014 applied a variety of AI algorithms to the data of 424 patients, providing insights into fibrosis progression without invasive biopsies, thus proposing a significant improvement in patient care [16]. AI systems have been recognized for their utility in processing structured and unstructured medical data, with machine learning algorithms playing a key role in these advanced systems. The importance of early detection and treatment for better health outcomes is echoed by experts such as Pietrangelo (2018) and Lok (2016), who stress that the initiation of treatment at the earliest sign of infection can prevent complications [17–21]. A particular study made use of 29 algorithmic predictors to craft an AI model for disease diagnosis, utilizing a dataset from Egyptian patients who underwent 18 months of liver cirrhosis treatment, highlighting the potential for AI to provide predictive insights without the need for liver biopsies [22]. The study referenced in [23] employed data mining methods to develop an artificial neural network (ANN) utilizing an extensive socio-medical dataset. This approach proved effective in predicting potential infections of the liver cirrhosis virus. Furthermore, research at the University of

Michigan [24] highlights the importance of reducing the cost of managing liver cirrhosis. Researchers at this institution have formulated a predictive analytics technique aimed at pinpointing high-risk patients, potentially facilitating early and effective treatment. The authors posit that this new method offers improved accuracy over previous research. Additionally, the researchers in [25] examined the impact of factors like gender and obesity on the prevalence of liver cirrhosis across different demographics. Their study emphasizes the critical role these factors play in the design of any diagnostic system, whether it employs artificial intelligence or traditional methods, to ensure the precision of results and optimize treatment strategies for healthcare professionals. Notable variations were observed among patients with regards to variables such as sex, body mass index (BMI), bilirubin, alanine aminotransferase (ALT), among others. It was found that average BMI values in male patients over 60 years old were lower compared to female patients under 60. Additionally, it was noted that higher BMI values correlate with an increased likelihood of earlier onset of complications related to liver cirrhosis in individuals infected with the virus [26]. In [27], the authors offer an extensive evaluation of three data mining techniques: decision trees, naive Bayes, and neural networks, for their utility in forecasting infections caused by the liver cirrhosis virus. The work highlighted in [28] addresses advancements in machine learning and artificial intelligence for anticipating esophageal varices, a condition that can result from chronic liver cirrhosis. This particular research identified 9 out of 24 factors as being most critical for assessment with the methodology they implemented. Furthermore, research indicated in [29] emphasizes the efficacy of decision tree learning methods in identifying individuals at heightened risk for severe liver fibrosis due to the liver cirrhosis virus, suggesting a potential to diminish or even supplant the need for liver biopsies, which are more invasive and carry certain disadvantages. Extensive research on the prediction and in-depth analysis of liver cirrhosis has been conducted, as detailed in studies [30–36], with additional examples cited in [37–40]. These cases highlight the substantial impact of liver cirrhosis and underscore the necessity for further research to effectively combat this disease. Consequently, this paper proposes a predictive model based on a dataset that is publicly available. This life-threatening disease is manageable if diagnosed in the early stages. Liver cirrhosis is often asymptomatic and as a result, diagnosing and treating patients during the early stages of illness is challenging. As the illness progresses to its latter stages, diagnosis and therapy become increasingly challenging. Here the main purpose is to build an intelligence model based on machine learning algorithms that may assist healthcare practitioners to forecast the possibility of a liver cirrhosis infection.

### IV. DATASET

**Data Collection:** Mayo Clinic conducted a trial on Primary Biliary Cirrhosis (PBC) of the liver from 1974 to 1984. The Dataset contains the information collected from the Mayo Clinic trial in primary biliary cirrhosis (PBC) of the liver con-

ducted between 1974 and 1984. It has been taken from the UCI ML Repository [45]. PBC patients participated in a randomized placebo-controlled trial of the drug D-penicillamine. First 312 cases participated in the randomized trial with comprehensive data. Additional 112 cases didn't join the trial but consented for basic measurements and survival. 106 cases from this group, along with the 312 participants, are included in the dataset. Clinical context in Fleming and Harrington's "Counting Processes and Survival Analysis" (1991). Additional insights in Dickson et al.'s "Hepatology" (1989) and Markus et al.'s "N Eng J of Med" (1989).

**DATASET:** Liver Cirrhosis Dataset [41], which comprises 615 entries across 20 columns. In this dataset, the 'Target' field outputs either '1' or '0', indicating the presence (1) or absence (0) of Liver Cirrhosis (LD) in patients. The output field more frequently displays a zero, signifying a higher number of non-cirrhotic cases compared to cirrhotic ones. Data preprocessing is implemented to balance the dataset distribution. This step occurs before preprocessing to ensure a proportional representation of both LD-afflicted and healthy data points

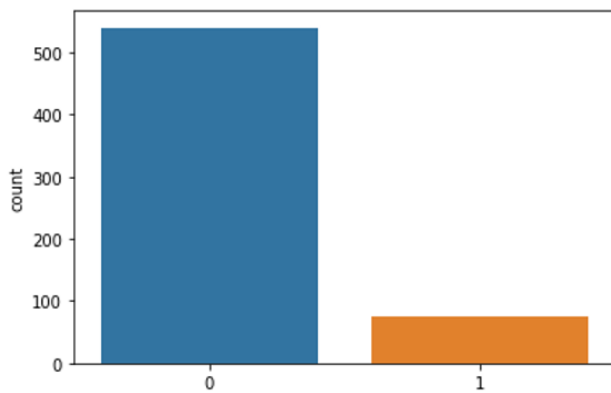


Figure 2: Visualization of Target Column

The dataset for LD prediction displays significant imbalance, with 615 entries, of which only 75 indicate LD presence and 540 suggest its absence. Addressing data imbalance is critical to avoid skewed results that could render the predictions unreliable. The initial step in developing an effective model, therefore, involves rectifying this imbalance, achieved here through the application of the SMOTE technique. Using SMOTE technique, Dataset has been balanced with equal number of presence (1) and absence (0) of Liver Cirrhosis (LD) in final dataset.

**Data Description:** The dataset contains a total 20 columns of which there are 13 Numerical Features and 7 Categorical features and their description is the following:

**ID:** Unique identifier given to each patient.

**N\_Days:** No. of days from registration to Death

**Status:** Status of the patient during study C, CL, D.

**Drug:** Type of drug the patient is consuming

**Age:** Patient age in days.

**Sex:** Gender of Patient (M or F).

**Ascites:** Presence of Ascites (N or Y).

**Hepatomegaly:** Presence of it (N or Y).

**Spiders:** Presence of spiders (N or Y).

**Edema:** Presence of edema (N, S, or Y).

**Bilirubin:** Serum bilirubin in mg/dl.

**Cholesterol:** Serum cholesterol in mg/dl.

**Albumin:** Albumin in gm/dl.

**Copper:** Urine copper in ug/day.

**Alk\_Phos:** Alkaline phosphatase in U/liter.

**SGOT:** Serum glutamic oxaloacetic in U/ml.

**Triglycerides:** Triglycerides in mg/dl.

**Platelets:** Platelets per cubic ml/1000.

**Prothrombin:** Prothrombin time in seconds (s).

**Stage:** Stage of disease converted to (Yes, No).

**Data Preprocessing:** Cleaning and preprocessing the raw dataset to remove artifacts and noise is essential. The data preprocessing steps included:

- Reclassification of Stage values from four categories into a binary 'Yes' or 'No'.
- Imputation of missing numerical variables with respective mean and missing categorical variables with mode.
- Conversion of the Age feature from days into years.
- Removed the ID column as no duplicates were found.
- Identification and removal of outliers to refine the modeling process.
- Transformation of categorical values into numerical form for compatibility with the machine learning model.
- Exclusion of the 'Status' and 'N\_Days' variables to prevent data leakage and potential overfitting of the model.

Upon finalizing the data preparation and addressing the imbalance within the dataset, the next phase involves constructing the model. The dataset is partitioned into training and testing subsets, adhering to a 70% to 30% ratio. Various classification methods are then employed to train the model based on this division. The classification strategies implemented in this study include Random Forest, Decision Tree Classification, and Support Vector Machine.

## V. DESCRIPTIVE ANALYSIS

Over 92% of the total missing data is coming from the patients that did not participate in the clinical trial i.e., *No Drug*. Since the *No Drug* group makes up about a quarter of the total data, removing all it is not a possible option. 6 cases are missing *Stage* data, so those are removed.

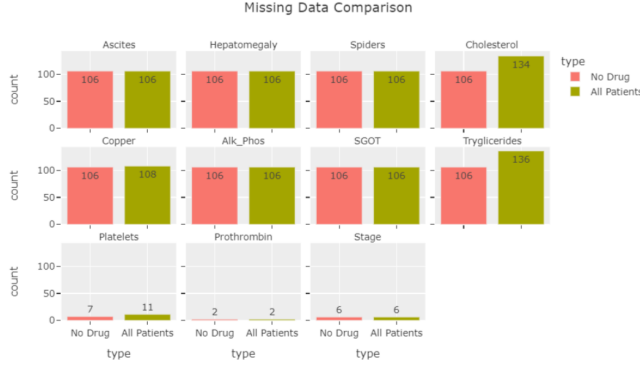


Figure 3: Missing Data Comparison

## VI. EXPLORATORY DATA ANALYSIS

The main purpose of EDA is to help look at data before making any assumptions. It can help identify obvious errors, as well as better understand patterns within the data, detect outliers or anomalous events, find interesting relations among the variables. EDA can help answer questions about standard deviations, categorical variables, and confidence intervals. Once EDA is complete and insights are drawn, its features can then be used for more sophisticated data analysis or modeling, including machine learning.

In the context of Ascites, there is a notable correlation between the severity of Ascites and an elevated risk of disease. Additionally, the presence of spiders appears to be associated with an increased likelihood of developing the disease.

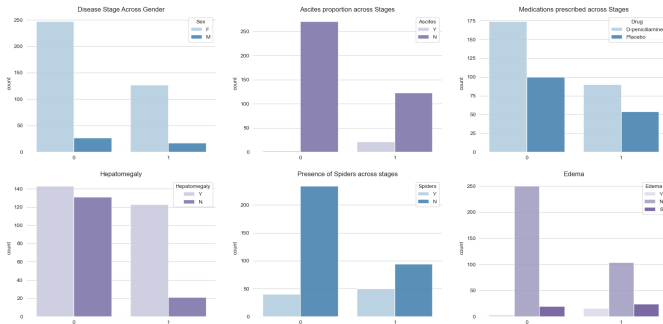


Figure 4: Analysis of Disease Characteristics across Stages

Our analysis indicates that variables such as Age, Prothrombin, and Copper exhibit a positive correlation with disease risk; as their values rise, the probability of disease also increases.

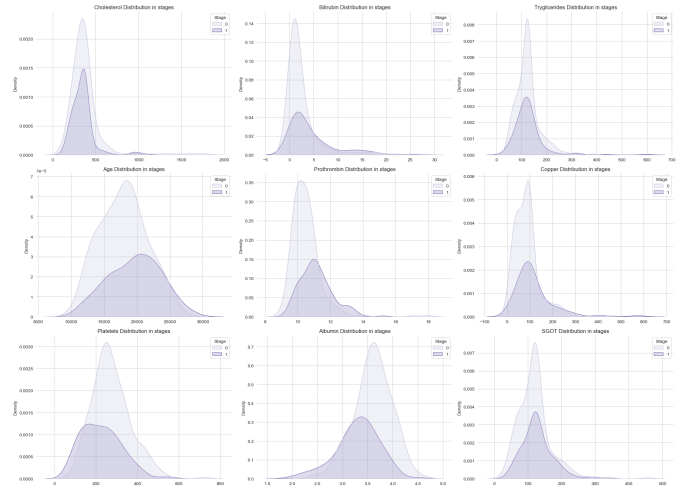


Figure 5: Distribution of Biochemical Markers in Disease Progression

Upon normalizing the data, it becomes evident from the plots that numerous outliers are present. A significant number of these outliers are situated between 2 and 3 standard deviations away from their respective means. The plots suggest a significant presence of outliers.

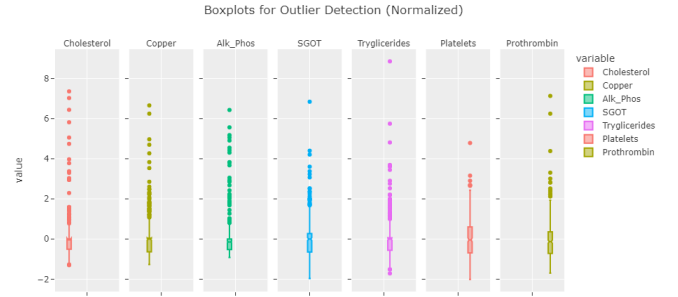


Figure 6: Boxplot for Outlier Detection

Some outliers exceed the 3rd standard deviation. These will be eliminated during data preprocessing, but we should first visualize them. Out of the total, 49 values are considered outliers, constituting roughly 10% of our dataset. However, removing them will enhance our model's predictive accuracy.

## VII. MODEL DEVELOPMENT

Once the data has been preprocessed and analyzed, it will be used to build a model. A preprocessed dataset and machine learning algorithms are required for model creation. Decision Tree Classification, Random Forest Classification, and Support Vector Machine (SVM) are chosen to use in this proposed system. The accuracy metrics Accuracy Score, Precision Score, Recall Score, and F1 Score are used to evaluate the models once it is created. Figure depicts the proposed system's block diagram.

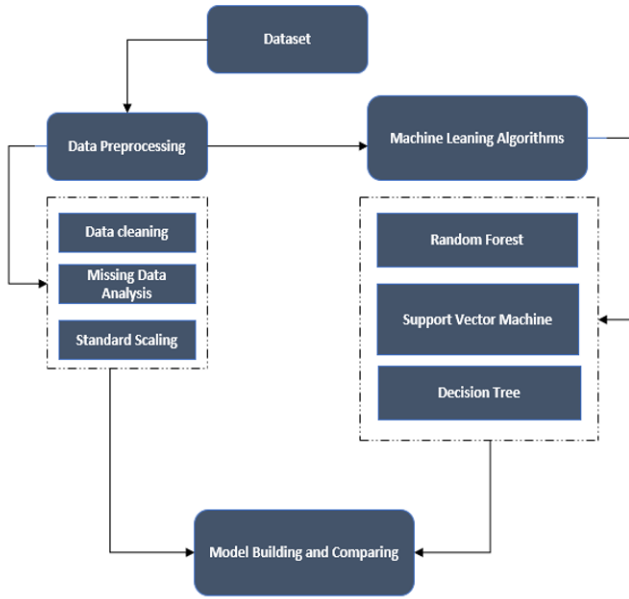


Figure 7: Proposed Method for Liver Cirrhosis Prediction

Initially, established a foundational model employing the Logistic Regression technique, which yielded a promising initial accuracy rate of 71%. Then experimented with the Decision Tree algorithm, which surprisingly produced a slightly lower accuracy of 69%, despite the usual expectation of tree-based models outperforming basic models. After evaluating several different models like SVM Classifier, finally settled on the Random Forest Model, which delivered the highest accuracy at 77%—a notable achievement given the dataset’s modest size of only 400 entries.

## VIII. APPLIED ALGORITHMS

Liver cirrhosis, a condition increasingly observed in medical studies, is a prevalent subject of research due to its rising incidence. This paper evaluates machine learning algorithms aimed at forecasting the recurrence of liver cirrhosis, utilizing a publicly accessible dataset dedicated to liver cirrhosis prediction. The algorithms explored in this study include:

- Logistic Regression (Base Line Model)
- Support Vector Machine
- Decision Tree
- Random Forest

**Random Forest:** Random Forest Classification was used as the algorithm for classification. It operates by constructing multiple decision trees, each trained on a randomly selected subset of the data, collectively forming a Random Forest. During the training process, these trees are developed, producing individual outputs. The algorithm then applies a voting mechanism

to determine its final prediction, where each decision tree casts a vote for one of two possible outcomes ('LD' or 'Healthy'). The outcome receiving the majority of votes from the trees is selected as the ultimate prediction of the Random Forest approach. Figure 7 depicts a random forest classification diagram.

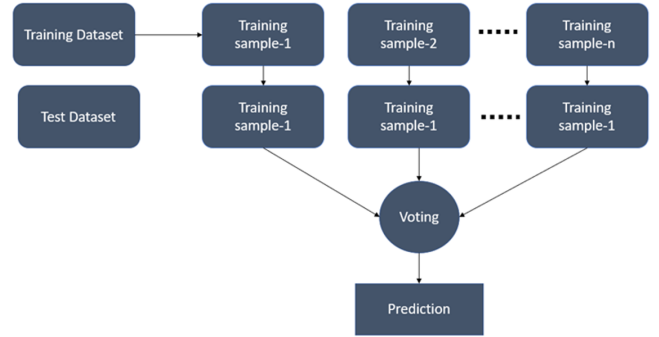


Figure 8: Working Process of Random Forest Classifier

The versatility of the Random Forest algorithm stands out as one of its key advantages. It is capable of handling tasks related to both recurrence prediction and classification effectively. This method allows for easy identification of the most influential features due to its feature importance metric. Moreover, its default hyperparameters often set well-defined expectations, contributing to its strategic effectiveness. A thorough comprehension of these hyperparameters is crucial, given their limited number. Although overfitting is a prevalent challenge in machine learning, the Random Forest classifier is generally less prone to this issue, largely because the presence of numerous trees in the forest helps to mitigate the risk of overfitting the model.

**Decision Tree:** The use of Decision Trees can be applied to both regression and classification problems. Because the input variables have a related output variable, this is a supervised learning approach. A tree can be seen in its resemblance. A specific parameter is used to divide the data constantly in this method.

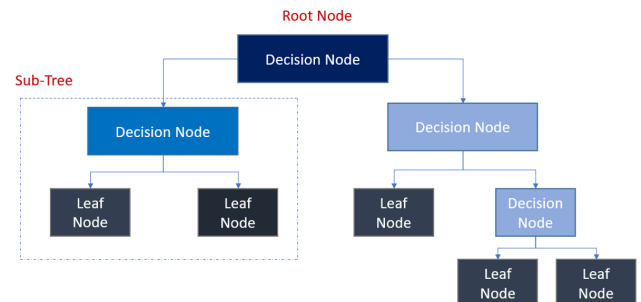


Figure 9: Structure of Decision Tree Classifier

Both the Decision and Leaf nodes form the core of a decision tree. Nodes 1 and 2 are responsible for dividing data, whereas nodes 3 and 4 are responsible for producing the final output.



Using Figure 8, we can see the Decision Tree Classifier's basic structure and operation.

The Decision Tree serves as an intuitive visual guide to the decision-making pathway, making it straightforward and accessible. It significantly simplifies the resolution of decision-making challenges by laying out all conceivable solutions to a given problem. Furthermore, it has the advantage of requiring less frequent data cleaning compared to alternative methods.

**Support Vector Machine:** The Support Vector Machine (SVM) model employs spatial points to delineate examples, creating a distinct separation between different categories with the widest possible margin. This spatial representation ensures that new instances are classified into one of the categories, based on which side of the divide they fall. Figure 9 shows the block diagram of the support vector machine.

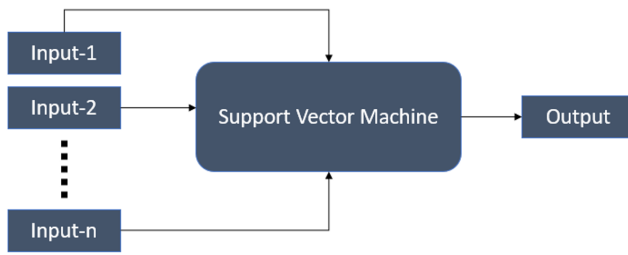


Figure 10: Structure of Support Vector Machine

Alternative classification approaches, including nonlinear classification through the utilization of the "kernel trick," achieve comparable success by implicitly transforming inputs into higher-dimensional functional spaces.

**Voting Classifier:** Voting Classifier using the "hard" option, this model combines predictions from multiple base classifiers, including Random Forest and SVM (Support Vector Machine), to make final predictions based on a simple majority vote system.

Each classifier independently makes a prediction for each data point, and the class that receives the most votes across all

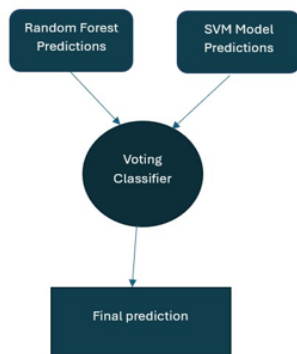


Figure 11: Voting Classifier

classifiers is chosen as the final output. This ensemble technique is particularly powerful for improving model performance and stability by leveraging the strengths of diverse underlying algorithms, reducing the likelihood of overfitting to the training data, and often achieving better accuracy than any single classifier alone. This method is straightforward yet robust, making it suitable for a variety of classification tasks.

**Evaluation Matrix:** The confusion matrix or evaluation matrix is seen in Figure 12 which is a tool for assessing the performance of a machine learning classification algorithm.

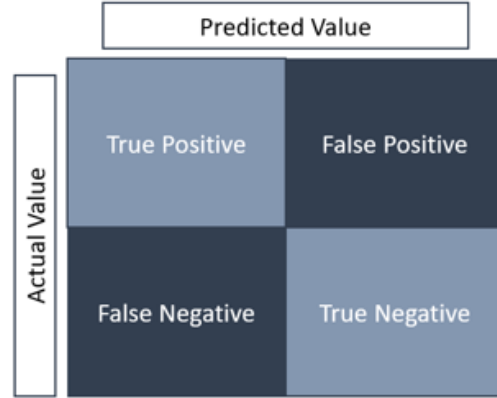


Figure 12: Confusion Matrix

This matrix has been applied to test all models, revealing both their accuracy and inaccuracies. It differentiates between correct predictions, represented by true positives and true negatives, and incorrect predictions, denoted by false positives and false negatives.

## IX. RESULT ANALYSIS

**1. Random Forest:** The final F1 score is 86 percent in this case. Individual F1 scores are at 85 percent for healthy persons and at 86 percent for those who have tested positive for LD. The Random Forest model's forecast is shown in Figure 13.

RandomForest Accuracy: 0.8571428571428571

RandomForest Classification Report:				
	precision	recall	f1-score	support
0	0.90	0.80	0.85	35
1	0.82	0.91	0.86	35
accuracy			0.86	70
macro avg	0.86	0.86	0.86	70
weighted avg	0.86	0.86	0.86	70

Figure 13: Classification Report of Random Forest

Model performance and expected outcomes are shown in a confusion matrix. Here, TP = 28, TN = 32, FP = 7, and FN = 3,

which means there were 60 correct guesses and 10 wrong ones. This model can predict 28 healthy and 32 LD-affected patient records correctly. But it denoted 7 healthy patient records as LD positive and 3 LD positive patients as healthy.

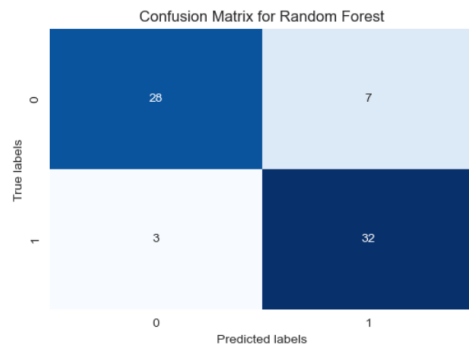


Figure 14: Confusion Matrix of Random Forest

**2. Decision Tree:** The final F1 score is 76 percent in this case. Individual F1 scores are at 75 percent for healthy persons and at 77 percent for those who have tested positive for LD. The Random Forest model's forecast is shown in Figure 14.

Decision Trees Accuracy: 0.7571428571428571

Decision Trees Classification Report:				
	precision	recall	f1-score	support
0	0.78	0.71	0.75	35
1	0.74	0.80	0.77	35
accuracy			0.76	70
macro avg	0.76	0.76	0.76	70
weighted avg	0.76	0.76	0.76	70

Figure 15: Classification Report of Decision Tree

Model performance and expected outcomes are shown in a confusion matrix. Here, TP = 27, TN = 28, FP = 8, and FN = 7, which means there were 55 correct guesses and 15 wrong ones. This model can predict 27 healthy and 28 LD-affected patient records correctly. But it denoted 7 healthy patient records as LD positive and 3 LD positive patients as healthy.

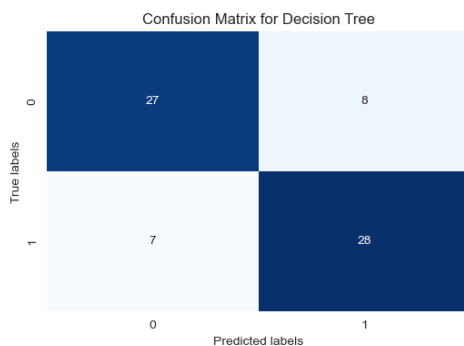


Figure 16: Confusion Matrix of Decision Tree

**3. SVM Model:** The final F1 score is 83 percent in this case.

Individual F1 scores are at 82 percent for healthy persons and at 83 percent for those who have tested positive for LD. The Random Forest model's forecast is shown in Figure 17.

SVM Accuracy: 0.8285714285714286

SVM Classification Report:				
	precision	recall	f1-score	support
0	0.85	0.80	0.82	35
1	0.81	0.86	0.83	35
accuracy			0.83	70
macro avg	0.83	0.83	0.83	70
weighted avg	0.83	0.83	0.83	70

Figure 17: Classification Report of SVM Model

Model performance and expected outcomes are shown in a confusion matrix. Here, TP = 28, TN = 30, FP = 7, and FN = 5, which means there were 58 correct guesses and 12 wrong ones. This model can predict 28 healthy and 30 LD-affected patient records correctly. But it denoted 7 healthy patient records as LD positive and 5 LD positive patients as healthy.

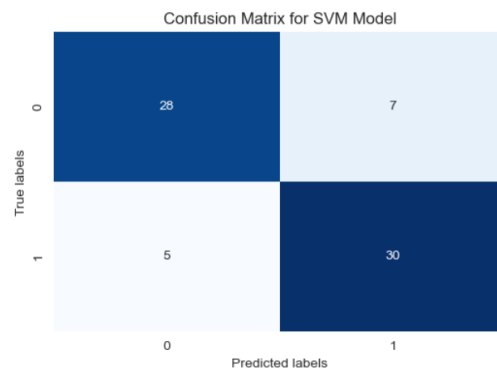


Figure 18: Confusion Matrix of SVM Model

**4. Voting Classifier:** The final F1 score is 87 percent in this case. Individual F1 scores are at 87 percent for healthy persons and at 87 percent for those who have tested positive for LD. The Random Forest model's forecast is shown in Figure 19.

Voting Classifier Accuracy Score: 0.8714285714285714

Classification Report:				
	precision	recall	f1-score	support
0	0.86	0.89	0.87	35
1	0.88	0.86	0.87	35
accuracy			0.87	70
macro avg	0.87	0.87	0.87	70
weighted avg	0.87	0.87	0.87	70

Figure 19: Classification Report of Voting Classifier

Model performance and expected outcomes are shown in a confusion matrix. Here, TP = 31, TN = 30, FP = 4, and FN = 5, which means there were 61 correct guesses and 9 wrong ones.

This model can predict 31 healthy and 30 LD-affected patient records correctly. But it denoted 4 healthy patient records as LD positive and 5 LD positive patients as healthy.

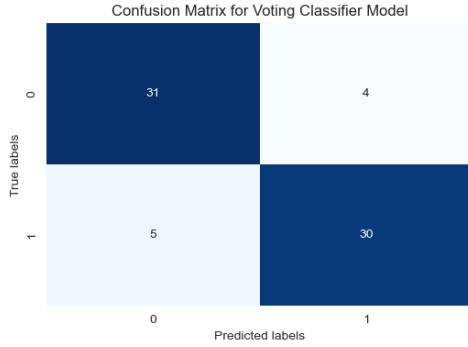


Figure 20: Confusion Matrix of Voting Classifier

The table clearly shows that Random Forest and Support Vector Machine are the best models among the several models in the framework. Using the same method, previous work [29] achieved an accuracy of 75 percent. Also using SVM ref [31] achieved 75 percent accuracy which is lower than our implemented model's result. Voting classifier is introduced in our work which is not present in the previous study ( So, marked as NA ) which achieved accuracy of 87.14% Table 1 shows the result comparison with previously reported results.

Table 1: Result Comparison			
Reference Number	Algorithm Name	Accuracy (%)	Accuracy of my work (%)
31	Random Forest	80.3	85.7
29	Support Vector Machine	75	82.8
NA	Voting Classifier	NA	87.14

Below ROC curve plot shows the trade-off between true positive rate (TPR) and false positive rate (FPR) for RandomForest, SVM, and the Voting Classifier, with respective AUC scores of 0.92, 0.92, and 0.93.

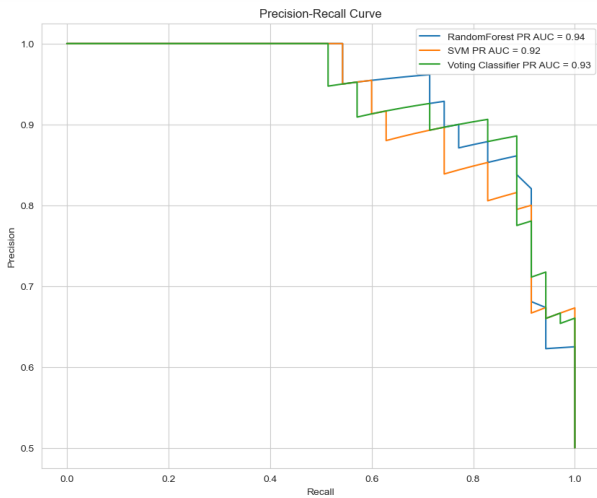


Figure 21: ROC Curves of RF, SVM and Voting Classifier

The Voting Classifier has the highest AUC, indicating better discrimination between classes compared to RandomForest and SVM.

Below plot is Precision-recall curve for RandomForest, SVM, and the Voting Classifier. It shows the relationship between precision and recall for each classifier. The Voting Classifier has the highest Precision-Recall AUC (0.94), indicating better performance in balancing precision and recall compared to RandomForest and SVM, which both have a Precision-Recall AUC of 0.92.

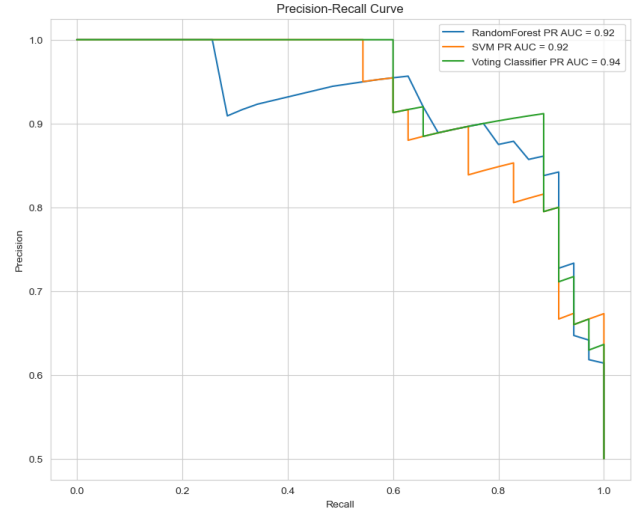


Figure 22: Precision-recall Curve of RF, SVM and Voting Classifier

## X. CONCLUSION & FUTURE WORK

Liver Disease (LD) is a potentially fatal infection that must be treated immediately to avert future complications. The construction of a machine learning model may aid in the detection of LD and may help mitigate its long-term health consequences. The Voting Classifier model demonstrated the effectiveness of ensemble learning, particularly how combining different types of classifiers can enhance predictive accuracy and model robustness compared to individual models. Utilizing hard voting, the classifier successfully integrated diverse of both Random Forest and SVM models, highlighting the strength of ensemble approaches in handling complex, varied datasets. For future work, the project could explore incorporating additional classifiers like Gradient Boosting or Neural Networks to further diversify the decision-making process. Experimenting with hard voting to leverage probability estimates for potentially more nuanced decision-making could also prove beneficial. With about 400 training examples we were able to achieve pretty significant scores, however the scores can massively improve by training our model with a bigger dataset. Additional features like diet, Weight, BMI also can play a significant role in the prediction. Lastly, applying the ensemble method to other complex datasets or in different domains could validate its versatility and robustness across various scenarios.



## References

- [1] World Health Organization, "Liver cirrhosis [Liver cirrhosis]," WHO, 2020.
- [2] R. Stoean, C. Stoean, M. Lupsor, H. Stefanescu, and R. Badea, "Evolutionary-driven support vector machines for determining the degree of liver fibrosis in chronic Liver cirrhosis," *Artificial Intelligence in Medicine*, vol. 51, no. 1, pp. 53–65, 2011.
- [3] A. A. Mohamed, T. A. Elbedewy, M. El-Serafy, N. El-Toukhy, W. Ahmed et al., "Liver cirrhosis virus: A global view," *World Journal of Hepatology*, vol. 7, no. 26, pp. 2676–2680, 2015.
- [4] R. Huang, H. Rao, M. Yang, Y. Gao, J. Wang et al., "Noninvasive measurements predict liver fibrosis well in Liver cirrhosis virus patients after direct-acting antiviral therapy," *Digestive Diseases and Sciences*, vol. 65, no. 5, pp. 1491–1500, 2020.
- [5] Z. Cheng, Y. Zhang and C. Zhou, "QSAR models for phosphoramidate prodrugs of 2'-methylcytidine as inhibitors of Liver cirrhosis virus based on PSO boosting," *Chemical Biology & Drug Design*, vol. 78, no. 6, pp. 948–959, 2011.
- [6] L. Singh, R. R. Janghel and S. P. Sahu, "Classification of hepatic disease using machine learning algorithms," in *Advances in Biomedical Engineering and Technology*. Berlin, Germany: Springer, pp. 161–173, 2021.
- [7] J. Vergniol, J. Foucher, E. Terrebbonne, P. Bernard, B. le Bail et al., "Noninvasive tests for fibrosis and liver stiffness predict 5-year outcomes of patients with chronic Liver cirrhosis," *Gastroenterology*, vol. 140, no. 7, pp. 1970–1979, 2011.
- [8] J. Shen, C. J. P. Zhang, B. Jiang, J. Chen, J. Song et al., "Artificial intelligence versus clinicians in disease diagnosis: Systematic review," *JMIR Medical Informatics*, vol. 21, no. 8, pp. 1–15, 2019.
- [9] Y. Murawaki, Y. Ikuta, K. Okamoto, M. Koda and H. Kawasaki, "Diagnostic value of serum markers of connective tissue turnover for predicting histological staging and grading in patients with chronic Liver cirrhosis," *Journal of Gastroenterology*, vol. 36, no. 6, pp. 399–406, 2001.
- [10] C. Lackner, G. Struber, B. Liegl, S. Leibl, P. Ofner et al., "Comparison and validation of simple noninvasive tests for prediction of fibrosis in chronic Liver cirrhosis," *Hepatology*, vol. 41, no. 6, pp. 1376–1382, 2005.
- [11] C.-T. Wai, J. K. Greenon, R. J. Fontana, J. D. Kalbfleisch, J. A. Marrero et al., "A simple noninvasive index can predict both significant fibrosis and cirrhosis in patients with chronic Liver cirrhosis," *Hepatology*, vol. 38, no. 2, pp. 518–526, 2003.
- [12] P. Halfon, M. Bourlière, G. Pénaranda, R. Deydier, C. Renou et al., "Accuracy of hyaluronic acid level for predicting liver fibrosis stages in patients with Liver cirrhosis virus," *Comparative Hepatology*, vol. 4, no. 1, pp. 6, 2005.
- [13] Healthline, "What are the stages of Liver cirrhosis?-HepatitisC.net," Healthline Media, 2018. [Online]. Available: <https://hepatitisc.net/living/what-are-the-stages-of-hepatitis-c/> [Accessed: 15-Aug-2020].
- [14] MHealth Lab, "Early detection, early treatment for Liver cirrhosis," MHealth Lab, 2016. [Online]. Available: <https://labblog.uofmhealth.org/rounds/early-detection-early-treatment-for-hepatitis-c> [Accessed: 20- Aug-2020].
- [15] E. Gupta, M. Bajpai and A. Choudhary, "Liver cirrhosis virus: Screening, diagnosis, and interpretation of laboratory assays," *Asian Journal of Transfusion Science*, vol. 8, no. 1, pp. 19–25, 2014.
- [16] F. Jiang, Y. Jiang, H. Zhi, Y. Dong, H. Li et al., "Artificial intelligence in healthcare: Past, present and future," *Stroke and Vascular Neurology*, vol. 2, no. 4, pp. 230–243, 2017.
- [17] N. Papadopoulos, S. Vasileiadi, M. Papavdi, E. Sveroni, P. Antonakaki et al., "Liver fibrosis staging with combination of APRI and FIB-4 scoring systems in chronic Liver cirrhosis as an alternative to transient elastography," *Annals of Gastroenterology*, vol. 32, no. 5, pp. 498, 2019.
- [18] Y. Zhao, P. H. Thuraiarah, R. Kumar, J. Tan, E. K. Teo et al., "Novel non-invasive score to predict cirrhosis in the era of Liver cirrhosis elimination: A population study of ex-substance users in Singapore," *Hepatobiliary & Pancreatic Diseases International*, vol. 18, no. 2, pp. 143–148, 2019.
- [19] J. Tani, A. Morishita, T. Sakamoto, K. Takuma, M. Nakahara et al., "Simple scoring system for prediction of hepatocellular carcinoma occurrence after Liver cirrhosis virus eradication by direct-acting antiviral treatment: All kagawa Liver cirrhosis group study," *Oncology Letters*, vol. 19, no. 3, pp. 2205–2212, 2020.
- [20] S. C. R. Nandipati, C. XinYing and K. K. Wah, "Liver cirrhosis virus (LD) prediction by machine learning techniques," *Applications of Modelling and Simulation*, vol. 4, pp. 89–100, 2020.
- [21] T. M. K. Motawi, N. A. H. Sadik, D. Sabry, N. N. Shahin and A. S., "Fahim, rs2267531, a promoter SNP within glypican-3 gene in the X chromosome, is associated with hepatocellular carcinoma in Egyptians," *Scientific Reports*, vol. 9, no. 1, pp. 1–10, 2019.
- [22] M. Reiser, B. Wiebner and J. Hirsch, "Neural-network analysis of socio-medical data to identify predictors of undiagnosed Liver cirrhosis virus infections in Germany (DETECT)," *Journal of Translational Medicine*, vol. 17, no. 1, pp. 1–7, 2019.
- [23] J. Bresnick, "Predictive analytics identify high risk Liver cirrhosis patients," *Health IT Analytics*, 2015. [Online]. Available: <https://healthitanalytics.com/news/predictive-analytics-identify-high-risk-hepatitis-cpatients/> [Accessed: 15-Oct-2020].
- [24] Y. C. Tsao, J. Y. Chen, W. C. Yeh, Y. S. Peng and W. C. Li, "Association between visceral obesity and Liver cirrhosis infection stratified by gender: A cross-sectional study in Taiwan," *BMJ Open*, vol. 7, no. 11, pp. e017117, 2017.
- [25] T. Akiyama, T. Mizuta, S. Kawazoe, Y. Eguchi, Y. Kawaguchi et al., "Body mass index is associated with age-at-onset of LD infected hepatocellular carcinoma patients," *World Journal of Gastroenterology*, vol. 17, no. 7, pp. 914–921, 2011.
- [26] A. A. A. Radwan and H. Mamdouh, "An analysis of Liver cirrhosis virus prediction using different data mining techniques," *International Journal of Computer Science, Engineering and Information Technology*, vol. 3, no. 4, pp. 209–220, 2013.
- [27] S. M. Abd El-Salam, M. M. Ezz, S. Hashem, W. Elakel, R. Salama et al., "Performance of machine learning approaches on prediction of esophageal varices for Egyptian chronic Liver cirrhosis patients," *Informatics in Medicine Unlocked*, vol. 17, no. September, pp. 100267, 2019.

- [28] S. Hashem, G. Esmat, W. Elakel, S. Habashy, S. Abdel Raouf et al., "Accurate prediction of advanced liver fibrosis using the decision tree learning algorithm in chronic Liver cirrhosis Egyptian patients," *Gastroenterology Research and Practice*, vol. 2016, pp. 1–7, 2016.
- [29] E. W. Abd El-Wahab, H. A. Ayoub, A. A. Shorbila, A. Mikheal, M. Fadl et al., "Noninvasive biomarkers predict improvement in liver fibrosis after successful generic DAAs based therapy of chronic Liver cirrhosis in Egypt," *Clinical Epidemiology and Global Health*, vol. 8, no. 4, pp. 1177–1188, 2020.
- [30] J.-P. Zarski, S. David-Tchouda, C. Trocme, J. Margier, A. Vilotitch et al., "Non-invasive fibrosis tests to predict complications in compensated post-Liver cirrhosis cirrhosis," *Clinics and Research in Hepatology and Gastroenterology*, vol. 44, no. 4, pp. 524–531, 2020.
- [31] K. Fujita, K. Oura, H. Yoneyama, T. Shi, K. Takuma et al., "Albumin-bilirubin score indicates liver fibrosis staging and prognosis in patients with chronic Liver cirrhosis," *Hepatology Research*, vol. 49, no. 7, pp. 731–742, 2019.
- [32] S. Hashem, M. ElHefnawi, S. Habashy, M. El-Adawy, G. Esmat et al., "Machine learning prediction models for diagnosing hepatocellular carcinoma with LD-related chronic Liver cirrhosis," *Computer Methods and Programs in Biomedicine*, vol. 196, pp. 105551, 2020.
- [33] H. M. Fayed, H. S. Mahmoud and A. E. M. Ali, "The utility of retinol-binding protein 4 in predicting liver fibrosis in chronic Liver cirrhosis patients in response to direct-acting antivirals," *Clinical and Experimental Gastroenterology*, vol. 13, pp. 53, 2020.
- [34] O. Hegazy, M. Allam, A. Sabry, M. A. S. Kohla, W. Abogharbia et al., "Liver stiffness measurement by transient elastography can predict outcome after hepatic resection for Liver cirrhosis virus-induced hepatocellular carcinoma," *The Egyptian Journal of Surgery*, vol. 38, no. 2, pp. 313, 2019.
- [35] X. Li, H. Xu and P. Gao, "Fibrosis index based on 4 factors (fib 4) predicts liver cirrhosis and hepatocellular carcinoma in chronic Liver cirrhosis virus (LD) patients," *Medical Science Monitor: International Medical Journal of Experimental and Clinical Research*, vol. 25, pp. 7243, 2019.
- [36] E. Gupta, M. Bajpai and A. Choudhary, "Liver cirrhosis virus: Screening, diagnosis, and interpretation of laboratory assays," *Asian Journal of Transfusion Science*, vol. 8, no. 1, pp. 19–25, 2014.
- [37] M. W. Nadeem, M. A. Al Ghamdi, M. Hussain, M. A. Khan, K. M. Khan et al., "Brain tumor analysis empowered with deep learning: A review, taxonomy, and future challenges," *Brain Sciences*, vol. 10, no. 2, pp. 118, 2020.
- [38] H. Malik, M. S. Farooq, A. Khelifi, A. Abid, J. N. Qureshi et al., "A comparison of transfer learning performance versus health experts in disease diagnosis from medical imaging," *IEEE Access*, vol. 8, pp. 139367–139386, 2020.
- [39] M. W. Nadeem, H. G. Goh, A. Ali, M. Hussain and M. A. Khan, "Bone age assessment empowered with deep learning: A survey, open research challenges and future directions," *Diagnostics*, vol. 10, no. 10, pp. 781, 2020.
- [40] H. Khalid, M. Hussain, M. A. Al Ghamdi, T. Khalid, K. Khalid et al., "A comparative systematic literature review on knee bone reports from MRI, X-rays and CT scans using deep learning and machine learning methodologies," *Diagnostics*, vol. 10, no. 8, pp. 518, 2020.
- [41] Liver cirrhosis Dataset. Available on: <https://www.kaggle.com/datasets/fedesoriano/cirrhosis-prediction-dataset>