**Faculty of Engineering, Environment and Computing**

**School of Computing, Electronics, and Mathematics**

**An Originality Declaration**

I hereby declare that this project report " Liver Diseases Detection using Data Mining is an entirely original work of mine, with no parts or whole copies from other sources other than those properly cited. As a result, any reference to previously published material (found in books, journals, magazines, the internet, etc.) has been made and is included in the main report as well as in the References or Bibliography lists. In addition, I consent to the possibility of this project's electronic copy being kept and used in the fight against plagiarism.

**Statement of Copyright**

**Statement of ethical engagement**

It is declared herewith that an application for this project has been made to the Coventry University ethics monitoring website (https//ethics.coventry.ac.uk/), and the application number is provided below  https//ethics.coventry.ac.uk

**Abstract**

When data mining and machine learning (ML) techniques are combined, liver disease detection has changed dramatically, improving predictive results and diagnostic accuracy. For liver diseases to be effectively managed and treated, these developments are essential. The analysis of liver health heavily relies on data mining, the process of finding patterns and insights within large datasets. Clinical datasets are used to train machine learning models, such as decision trees, support vector machines (SVM), and neural networks, to identify and forecast the occurrence and course of liver disease. In order to diagnose diseases like cirrhosis, fatty liver disease, and hepatitis, these models look at a number of parameters, including bilirubin concentrations, liver enzyme measurements (such as ALT, AST), and other biochemical indicators. By dividing data into branches so that predictions can be made based on various attributes, decision trees offer a simple method.

When it comes to liver disease diagnostics, the integration of data mining and machine learning provides a more accurate method than conventional techniques. These methods allow for the prediction of the course of the disease as well as an increase in the accuracy of liver condition diagnosis. This predictive ability enables medical professionals to make proactive adjustments to patient care plans, guaranteeing prompt and focused interventions. Patients who have a high risk of developing severe liver disease, for example, can be identified early and given close monitoring and treatment.

Eventually, there is a lot of promise to enhance patient outcomes through the integration of data mining and machine learning in the detection of liver disease. Effective interventions implemented on time can stop the progression of the disease, lessen its complications, and improve patients' quality of life. This represents a new era in medical innovation, as these techniques continue to advance and have an increasing impact on patient care and the diagnosis of liver disease.

## Chapter 1: Introduction

Millions of people worldwide suffer from liver disease, making it a serious global health concern. The liver is an essential organ that performs many vital tasks, such as detoxifying toxic substances, synthesizing plasma proteins, regulating blood clotting, and storing energy. Bile is also produced for digestion. Life-threatening illnesses and serious health problems can result from impaired liver function. Acute and chronic liver diseases are the two main categories into which liver diseases fall. Chronic liver diseases, such as cirrhosis or liver cancer, frequently cause irreversible damage.

Liver diseases rank among the world's leading causes of death, with liver cirrhosis alone responsible for over a million deaths annually, according to the World Health Organization (WHO). Particularly high rates of alcohol abuse, hepatitis B and C infections, and non-alcoholic fatty liver disease (NAFLD), which is frequently linked to obesity and metabolic syndrome, are associated with a higher prevalence of liver disease. Liver disease has a significant financial cost because it requires liver transplantation in severe cases and comes with other healthcare expenses and productivity losses.

### 1.1. Background

Millions of individuals worldwide are impacted by liver diseases, which present a serious global health concern. An essential organ, the liver carries out vital tasks like energy storage, blood clotting regulation, detoxification, and the production of bile for digestion. Liver dysfunction can result in serious health problems and potentially fatal diseases.

Acute and chronic liver diseases can be distinguished from one another. While chronic liver diseases like cirrhosis, liver fibrosis, and liver cancer develop over the years and are typically irreversible, acute liver diseases occur suddenly and are frequently reversible. Particularly concerning are chronic liver diseases, which are frequently asymptomatic until advanced stages. The World Health Organization (WHO) estimates that liver cirrhosis alone causes more than a million fatalities per year. High alcohol intake, hepatitis B and C infections, and the growing incidence of non-alcoholic fatty liver disease (NAFLD), which is connected to obesity and metabolic syndrome, are all contributing factors.

Liver diseases have a large financial cost because severe cases necessitate expensive liver transplants, which also result in high healthcare costs and lost productivity. Effective management

depends on early detection, but conventional diagnostic techniques can be intrusive, costly, and timeconsuming. Data mining and machine learning have emerged as promising tools for the early diagnosis and detection of liver diseases thanks to technological advancements. These methods provide a non-invasive and economical substitute for conventional approaches by analyzing huge datasets to find patterns and forecast disease outcomes with high accuracy.

To detect liver disease, this study investigates the use of data mining and machine learning techniques. To create predictive models for early diagnosis, we plan to use an extensive dataset of clinical and biochemical parameters. We will assess the predictive power of machine learning algorithms like logistic regression, random forests, decision trees, K-nearest neighbors, gradient boosting, and naive Bayes about liver diseases.

The discoveries will add to the expanding corpus of information regarding the application of machine learning and data mining in healthcare, specifically in the identification of liver disease. To improve patient outcomes and lessen the cost of liver disease-related healthcare, this research will analyze the performance of different algorithms to develop more precise and effective diagnostic tools.

### 1.1.1. Problem Statement and Motivation

Due to their capacity to handle massive data sets and reveal latent patterns that are difficult for human experts to identify, data mining and machine learning techniques are being used to detect liver diseases. These techniques can greatly improve diagnostic accuracy, lessen the need for invasive procedures, and enable individualized treatment plans thanks to the improvement of computational technologies and the growing availability of medical data. Clinical professionals can make better decisions and provide better care for patients with liver diseases by incorporating these technologies into the healthcare system. Further enhancing healthcare delivery and resource allocation, these technologies can assist in tracking the course of diseases and forecasting patient outcomes.

The use of machine learning and data mining techniques for the diagnosis of liver diseases is driven by their capacity to process massive amounts of data and reveal hidden patterns that are difficult for human experts to identify. These techniques can greatly improve diagnostic accuracy, lessen the need for invasive procedures, and enable individualized treatment plans thanks to the growing availability of medical data and the development of computational technologies. Better patient care

and liver disease management result from clinicians' ability to make more informed decisions when these technologies are integrated into the healthcare system. These technologies can also be used to predict patient outcomes and track the course of diseases, which enhances the delivery of healthcare and the distribution of resources.

### 1.1.2. Importance of Early Detection

For liver disease to be effectively treated and managed, early detection is essential. An early diagnosis can stop a disease from getting worse, enhance patient outcomes, and lower the cost of healthcare overall. For example, the chance of developing cirrhosis or liver cancer can be considerably decreased by early detection of hepatitis and subsequent antiviral treatment. Similar to NAFLD, cirrhosis and non-alcoholic steatohepatitis (NASH) cannot advance if early intervention is provided.

Liver biopsy, imaging modalities (such as CT, MRI, and ultrasound scans), and biochemical blood tests (such as liver function tests) are examples of conventional methods for diagnosing liver disease. These techniques work well, but they can also be expensive, and invasive, and may not always result in an early diagnosis. For instance, a liver biopsy is an invasive procedure with potential risks and is thought to be the gold standard for diagnosing certain liver conditions. Even though they are noninvasive, imaging methods might miss liver disease early on.

### 1.1.3. The Role of Data Mining and Machine Learning

Recent developments in machine learning and data mining have demonstrated potential to improve the precision and effectiveness of liver disease detection. Finding relevant patterns in massive datasets is known as data mining, and it can be extremely important for diagnosing medical conditions. As a branch of artificial intelligence, machine learning focuses on teaching algorithms to learn from data and make decisions or predictions without explicit programming.

Large-scale clinical and biochemical data can be analyzed by machine learning algorithms to find patterns and relationships that human clinicians might not notice right away. These methods may be able to identify liver disease early on, forecast how the condition will progress, and tailor treatment regimens to the specific needs of each patient.

## 1.2. Project Aim and Objectives

### 1.2.1. Objectives

1.      Information Gathering and Preprocessing Compile and purify data on liver disease from dependable sources, deal with null values, and standardize the data to guarantee accuracy.

2.      Choose the features that are most pertinent to the development of liver disease. These features may include lifestyle factors, biochemical markers, and demographic data.

3.      Model Creation Create predictive models by utilizing a variety of machine learning techniques, including Random Forests, Decision Trees, Support Vector Machines (SVM), and Neural Networks. 4. Metrics such as accuracy, precision, recall, F1-score, and ROC-AUC curves can be used to evaluate the performance of the developed models and ascertain their predictive power for liver disease.

### 1.2.2. Aims

1.      Enhance Early Detection Provide a dependable and accurate diagnostic instrument to identify liver diseases early on, increasing the likelihood that treatment and management will be effective.

2.      Improve Diagnostic Accuracy Compared to conventional methods, the use of sophisticated data mining and machine learning techniques can improve the accuracy of liver disease diagnosis.

3.      Assist Healthcare Professionals Give medical staff members an effective tool to help them make wise decisions, which will ultimately improve patient care.

4.      Integrate with Clinical Systems Create a model that can be used for real-time analysis and application by integrating it with clinical decision support systems.

## 1.3. Overview of the Report

An extensive investigation into the use of data mining and machine learning techniques for liver disease detection is presented in this paper. Outlining the motivation, problem statement, and research objectives, the introduction gives background information on liver diseases and emphasizes the importance of early detection. This chapter reviews the literature on liver disease detection using different machine learning algorithms, summarizing important studies, using classification schemes, and pointing out areas of current research deficiency. Along with the procedures for data preparation, model training, and evaluation metrics, the methodology chapter also covers the research design, data collection, preprocessing, and machine learning techniques used. A comparative analysis of the

models, confusion matrices, and performance metrics are among the findings from applying various machine learning models to the liver disease dataset that are presented in the results chapter.

The project schedule, risk and quality management techniques, and professional, social, legal, and ethical considerations are all covered in the project management section. The discussion chapter discusses the models' shortcomings and strengths while interpreting the data and comparing it to the literature review to determine why the models worked as they did. In conclusion, the chapter on future work and conclusions provides an overview of the study's major discoveries and contributions, addressing its successes, shortcomings, and suggestions for further research. Additional materials such as the project specification, requirements specification document, project presentation, certificate of ethics approval, and interim progress reports are included in the appendices. A comprehensive investigation and documentation of the research are ensured by this structured approach, providing lucid insights into the efficacy of data mining and machine learning techniques in the detection of liver diseases.

## Chapter 2:  Literature Review

## 2.1.  Introduction

Liver disease, which includes a wide range of liver-related conditions like cirrhosis, liver cancer, and hepatitis, is a major global health concern. The liver carries out vital tasks such as protein synthesis, detoxification, and the creation of biochemicals required for digestion. A compromised liver affects the homeostasis of the entire body, resulting in serious health problems and often fatal conditions. For liver disease to be effectively treated and managed, early detection is essential. This improves patient outcomes and lowers healthcare costs.

Medical diagnostics has significantly advanced thanks to the integration of data mining and machine learning techniques, which have produced creative methods for managing and detecting diseases. The identification and categorization of liver diseases using machine learning algorithms is one prominent application. Of them, the Gradient Boosting Machine (GBM) algorithm—which is comparable to Random Forest—has become well-known for managing medical data with accuracy and efficiency.

GBM is an ensemble learning technique that builds models in a step-by-step fashion, with each new model fixing the mistakes of its predecessors. Because of this iterative process, which improves the accuracy of the model, GBM is especially useful for difficult classification tasks like differentiating between liver diseases. GBM is a useful tool in clinical settings because of its resilience and versatility in handling different data types and formats.

Large datasets can be analyzed using GBM in liver disease diagnostics to find patterns and correlations that human clinicians might not notice right away. Effective treatment outcomes depend on early detection and accurate diagnosis, both of which are made possible by this capacity. Beyond liver disorders, GBM has potential applications in other medical diagnostic fields, such as cancer detection and cardiovascular disorders.

The integration of these tools into current medical workflows and the requirement for large, highquality datasets are two of the challenges involved in applying GBM and related techniques in

clinical practice. Upcoming studies and developments will probably concentrate on making these models easier to understand and making sure they can be easily incorporated into clinical decision-making procedures.

Using sophisticated machine learning methods such as GBM has great potential to improve patient care by increasing the precision and effectiveness of medical diagnostics.

## 2.2. Research papers and their Case Studies

**Case Study 1**

**Title** Disease Detection and Prediction Using the Liver Function Test Data A Review of Machine Learning Algorithms
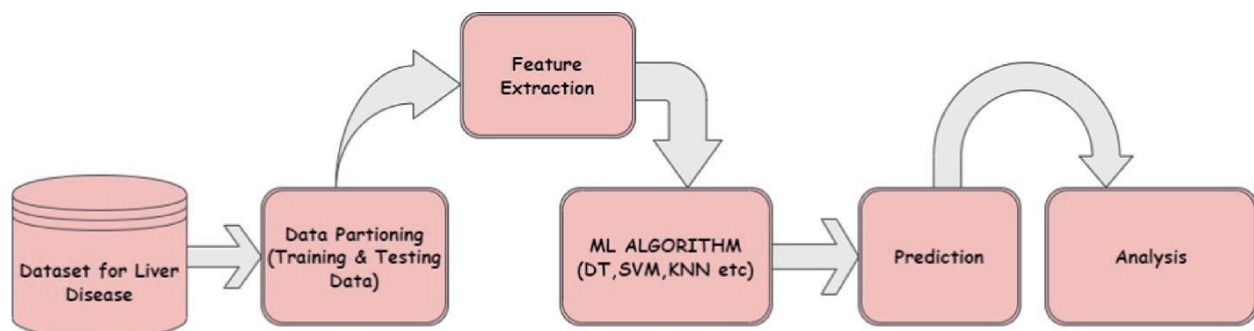
**Authors** Ifra Altaf, Muheet Ahmed Butt, Majid Zaman

**Summary** In light of liver function test (LFT) data, this paper examines machine learning algorithms that are used to identify and forecast liver disorders. Stressing the importance of influential features in disease recognition as well as the notable improvements in classification accuracy. The impact of algorithms like CMAC, RBF, PSO-LS-SVM, and ADTree in enhancing prediction accuracy is demonstrated by the study's examination of the relationship between LFT data and diabetes prediction.

**Case Study 2**

**Title** Detection of Liver Disease Using Machine Learning Techniques A Systematic Survey

**Authors** Geetika Singh, Charu Agarwal, Sonam Gupta

*Figure 1: General block diagram of Detection Model*

**Summary** This paper provides a systematic survey of different machine learning methods for liver disease detection. It reduces the need for labor-intensive procedures like liver biopsies and expert MRI analysis by addressing the need for automated diagnosis systems that can provide prompt and accurate results. Performance comparison of several machine learning algorithms, including Decision Trees, K-Nearest Neighbors (KNN), Support Vector Machines (SVM), and Genetic Algorithms, is included in the review along with information on the datasets and efficacy of each algorithm. Future directions and challenges in the prediction of liver disease are also covered.

**Case Study 3**

**Title** Machine Learning for Liver Disease Classification

**Authors** Ifra Altaf, Muheet Ahmed Butt, Majid Zaman

**Summary** To classify liver disease using medical data, this paper investigates the application of various machine learning techniques. Because liver disorders are typically diagnosed with more invasive techniques, the authors address the importance of automated systems for prompt and precise diagnosis. We compare and illustrate the efficacy and performance of various techniques, including Decision Trees, K-nearest neighbors (KNN), Support Vector Machines (SVM), and Genetic Algorithms. The research highlights the significance of these algorithms in enhancing diagnostic precision and mitigating the workload of medical practitioners.

**Case Study 4**

**Title** Early-Stage Detection of Liver Disease Through Machine Learning Algorithms

**Authors** Krittika Dutta, Satish Chandra, Mahendra Kumar Gourisaria

**Summary** In this paper, different machine learning algorithms are used to study the early-stage detection of liver disease. To avoid serious complications, it discusses how important early diagnosis is. Artificial Neural Networks (ANN), Logistic Regression, K-nearest neighbors (KNN), Support Vector Classifiers (SVC), Decision Trees, Random Forest, LR-SGD Classifiers, Passive-Aggressive,

AdaBoost, and Voting Classifiers are some of the models used in the study, both with and without Linear Discriminant Analysis (LDA). 99.96% accuracy was the highest for the Decision Tree algorithm.

**Case Study 5  (2020)**

**Title** Supervised Machine Learning Models for Liver Disease Risk Prediction
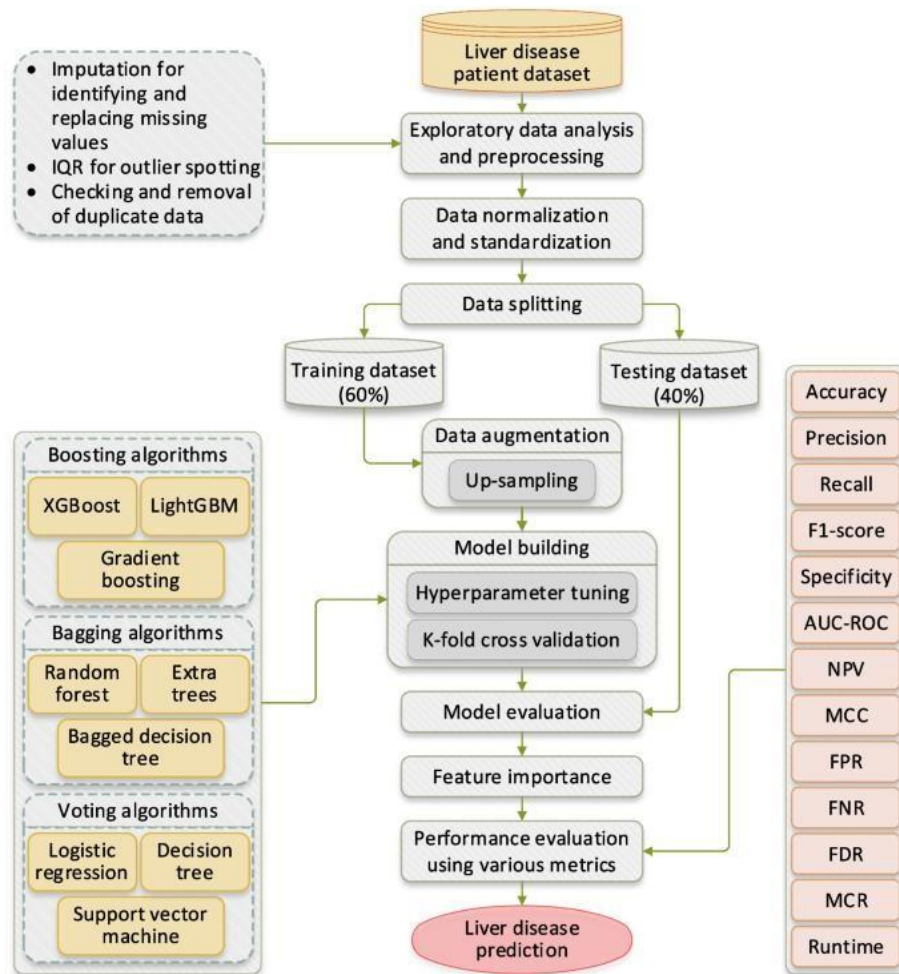
**Authors** Elias Dritsas, Maria Trigka



*Figure 2: Proposed methodology for research work*

**Summary**

The study on supervised machine learning models for liver disease risk prediction is presented in this paper. The Synthetic Minority Oversampling Technique (SMOTE) is used in the study to address data

imbalance and highlight the significance of balanced datasets for efficient model training. A voting classifier, random forest, and decision tree models are among the machine learning models that are assessed. In terms of precision, recall, F-measure, accuracy, and AUC metrics, the Voting Classifier—which combines the predictions of several models—performed better. The study concludes that combining these models can greatly increase the prediction accuracy of liver disease, giving medical professionals a useful tool for early diagnosis and treatment planning.

**Case Study 6  (2024)**

**Title** Improved liver disease prediction from clinical data through an evaluation of ensemble learning approaches

**Authors** Quadir et al.

**Summary** This study examines how different ensemble machine learning approaches can be applied to enhance the precision of liver disease prediction using clinical data. The goal of the study is to improve data preprocessing techniques like feature selection, imputation, balancing, and scaling to improve the performance of the classification models. Gradient Boosting (GB), Extreme Gradient Boosting (XGB), Bagging, Random Forest (RF), Extra Trees (ET), and Stacking are the six ensemble algorithms that the authors evaluate. With a testing accuracy of 91.82%, the Extra Trees classifier outperformed the others. The study also shows how well SMOTE (Synthetic Minority Over-sampling Technique) works to address problems with class imbalance. According to the results, when compared to individual models, ensemble learning techniques—specifically, Extra Trees—significantly improve liver disease prediction accuracy.

**Case Study 7**  (2018)

**Title** Prognosis of Liver Disease Using Machine Learning Algorithms

**Authors** Vyshali J. Gogi, Vijayalakshmi M.N

**Summary** By employing data mining techniques to examine liver function test (LFT) data, this study explores the application of machine learning algorithms for the prognosis of liver disease. Support

vector machines (SVM), logistic regression, decision trees, and linear discriminants are some of the classification algorithms used in the study to predict liver disease. The implementation of these algorithms is done with MATLAB 2016. With an ROC of 0.93, the Logistic Regression algorithm came in second to the Linear Discriminant algorithm, which had the highest prediction accuracy at 95.8%. Data pre-processing, SVM, Decision Trees, and Logistic Regression are among the techniques used in the journal.

## Case Study 8  (2021)

**Title** Early-Stage Detection of Liver Disease Through Machine Learning Algorithms

**Authors** Krittika Dutta, Satish Chandra, Mahendra Kumar Gourisaria

**Summary** The early-stage machine learning algorithmic detection of liver disease is the main focus of this paper. To avoid serious health complications, the study emphasizes the significance of early diagnosis. The authors use a variety of machine learning models, such as Random Forest, ANN, LRSGD Classifier, Passive-Aggressive, AdaBoost, K-Nearest Neighbors (KNN), Support Vector Classifier (SVC), Gaussian Naive Bayes, Decision Trees, and Random Forest. At 99.96% accuracy, the Decision Tree algorithm showed the highest level of performance. Additionally covered in the paper is how Linear Discriminant Analysis (LDA) can enhance model performance.

## Case Study 9  (2021)

**Title** The Diagnosis of Chronic Liver Disease Using Machine Learning Techniques

**Authors** Golmei Shaheamlung, Harshpreet Kaur

**Summary** In this work, the use of machine learning algorithms for the early detection and prediction of chronic liver disease (CLD) is explored in light of the growing concern over liver disease. With the help of random forest, logistic regression, and KNN classifiers, the authors suggest a hybrid classification model. For this study, Indian liver patient records from the Kaggle database are consulted. With accuracy, precision, and recall all taken into consideration, the model produces results with an accuracy of 77.58%. Utilizing Python and the Spyder tool, methods include feature selection, data preprocessing, and classification.

**Case Study 10** (2018)

**Title** Prediction of Liver Diseases Based on Machine Learning Technique for Big Data

**Authors** Engy A. El-Shafeiy, Ali I. El-Desouky, Sally M. Elghamrawy

**Summary** This study investigates the use of machine learning methods on big datasets to predict liver diseases. 5295 male patients and the remaining patients were female, making up the 7000 patient dataset with 23 attributes. For prediction, the study uses the Naive Bayes (NB), Boosted C5.0, and Support Vector Machine (SVM) algorithms. Accuracy, sensitivity, and specificity are the metrics used to assess these classifiers' performance.

**Case Study 11 (2021)**

**Title** The Diagnosis of Chronic Liver Disease Using Machine Learning Techniques

**Authors** Golmei Shaheamlung, Harshpreet Kaur

**Summary** The purpose of this study is to determine whether chronic liver disease can be diagnosed using machine learning techniques. The study discusses the rising incidence of liver disease worldwide and the significance of early detection for efficient treatment. The Kaggle database of Indian liver patient records is the dataset used by the authors to predict liver disease using a variety of machine learning algorithms, such as SVM, K-means clustering, KNN, Random Forest, and Logistic Regression. The accuracy of the suggested hybrid classification model was 77.58%. Preprocessing the data, choosing features, classifying the results, and evaluating the performance using metrics like accuracy, precision, and recall are all part of the process. Using the Spyder tool, the model is implemented in Python.

**Case Study 12  (2019)**

**Title** A Comparison of Data Mining Algorithms for Liver Disease Prediction on Imbalanced Data

**Authors** Ain Najwa Arbain, B. Yushalinie Pillay Balakrishnan

**Summary** The purpose of this paper is to examine the prediction of liver disease using data mining algorithms on unbalanced databases. This study compares several techniques, including Random Forest, AutoNeural, Logistic Regression, and K-Nearest Neighbor (k-NN). With a 99.794% accuracy rate, the study finds that the k-NN algorithm performs better than the others. Regarding earlier studies

employing the Andhra Pradesh liver disease dataset, the study emphasizes accuracy and ROC index to show how well the suggested model performs.

**Case Study 13 (2020)**

**Title** A Data Mining Approach to Prediction of Liver Diseases

**Authors** Nazim Razali, Aida Mustapha, Mohd Helmy Abd Wahab, Salama A Mostafa, Siti Khadijah Rostam

**Summary** This study investigates the use of data mining methods to forecast liver illnesses. Using the UCI repository dataset, the study analyzes and predicts liver diseases using a variety of classification algorithms, such as Naïve Bayes, Support Vector Machines (SVM), C4.5 Decision Tree, and Multilayer Perceptron. The Multilayer Perceptron algorithm demonstrated the highest accuracy of 71.59%, according to the results. The study emphasizes how crucial feature selection, data preprocessing, and performance evaluation metrics are to raise the accuracy of liver disease prediction.

**Case Study 14 ( 2016)**

**Title** Analysis of Data Mining Techniques for Healthcare Decision Support System Using Liver Disorder Dataset

**Authors** Tapas Ranjan Baitharu, Subhendu Kumar Pani

**Summary** Utilizing the Liver Function Test (LFT) dataset, this study compares several data mining approaches to forecast liver diseases. The research uses a variety of classifiers, such as the VFI, Naive Bayes, J48 Decision Tree, Artificial Neural Network (ANN), ZeroR, and 1BK algorithms. The Multilayer Perceptron (ANN), with a classification accuracy of 71.59%, was found to have the highest accuracy among the models. Developing intelligent medical decision support systems requires careful consideration of several factors, including data classification accuracy, computational efficiency, and classifier predictive performance.

**Case Study 15  (2020)**

**Title** A Survey on Machine Learning Techniques for the Diagnosis of Liver Disease

**Authors** Golmei Shaheamlung, Harshpreet Kaur, Mandeep Kaur

**Summary** This paper offers an extensive overview of machine-learning methods for liver disease diagnosis. It draws attention to the increasing incidence of liver diseases brought on by several variables and emphasizes the significance of early detection. With a particular focus on algorithms like SVM, KNN, K-Means clustering, neural networks, and decision trees, the study covers a variety of machine learning techniques, including supervised, unsupervised, semi-supervised, and reinforcement learning. Based on a comparison of these methods' accuracy, sensitivity, precision, and specificity, the study concludes that hybrid approaches might enhance predictive performance.

## Case Study 16  (2020)

**Title** A Survey on Machine Learning Techniques for the Diagnosis of Liver Disease

**Authors** Golmei Shaheamlung, Harshpreet Kaur, Mandeep Kaur

**Summary** An extensive review of the various machine-learning methods for liver disease diagnosis is given in this paper. The authors examine the need of early liver disease detection in light of variables such as heavy alcohol use, pollution, drug use, and tainted food. Various approaches to machine learning are covered in the study, such as semi-supervised, supervised, supervised, and reinforcement learning. Neural networks, decision trees, SVM, KNN, K-Means clustering, and Naïve Bayes are among the specific algorithms that have been evaluated. The performance of these approaches in terms of accuracy, sensitivity, and specificity is reviewed and contrasted in this paper. According to the survey's findings, ensemble and hybrid approaches can improve predictive performance and aid medical professionals in the early diagnosis of disease.

## Case Study 17  (2016)

**Title** Analysis of Data Mining Techniques for Healthcare Decision Support System Using Liver Disorder Dataset

**Authors** Tapas Ranjan Baitharu, Subhendu Kumar Pani

**Summary** To create a healthcare decision support system, this paper examines the application of different data mining techniques on a dataset of liver disorders. The study uses multiple classification algorithms to categorize liver disorders, including Naive Bayes, ZeroR, 1BK, Artificial Neural Network (ANN), J48 Decision Tree, and VFI. The writers evaluate and contrast these algorithms' efficacy and

rate of correction. According to the research, early detection and treatment of liver diseases depend heavily on accurate data classification. The study concludes that the accuracy and effectiveness of illness prediction models are greatly impacted by algorithmic choice.

## Summary Table of Research Papers

*Table 1: Summary Table for Research Table*

| S. N o | Paper Title | Yea r | Author | Propertie s Predicte d | Techniques/Us ed | Metrics Used for Measureme nt | Result |
|---|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  |  |

| 1 | Disease Detection and Prediction Using the Liver Function Test Data A Review of Machine Learning Algorithms | 2021 | Ifra Altaf, Muheet Ahmed Butt, Majid Zaman | Liver Disease, Diabetes | CMAC, RBF, PSO-LS-SVM, ADTree | Accuracy, Feature Importance | Effective in improving prediction accuracy |
| 2 | Detection of Liver Disease Using Machine Learning Techniques A Systematic Survey | 2021 | Geetika Singh, Charu Agarwal, Sonam Gupta | Liver Disease | Decision Trees, K-Nearest Neighbors (KNN), Support Vector Machines (SVM), Genetic Algorithms | Accuracy, Performance Comparison | Challenges and future directions discussed |
| 3 | Machine Learning for Liver Disease Classification | 2021 | Ifra Altaf, Muheet Ahmed Butt, Majid Zaman | Liver Disease | Decision Trees, K-nearest neighbors (KNN), Support Vector Machines (SVM), Genetic Algorithms | Accuracy, Effectiveness | Emphasizes the role of algorithms in improving diagnostic accuracy |

| 4 | Early-Stage Detection of Liver Disease Through Machine Learning Algorithms | 2021 | Krittika Dutta, Satish Chandra, Mahendra Kumar Gourisaria | Liver Disease | ANN, Logistic Regression, KNN, SVC, Gaussian Naive Bayes, Decision Trees, Random Forest, LR-SGD Classifier, PassiveAggressive, AdaBoost, Voting Classifier, LDA | Accuracy | Decision Tree showed the highest accuracy at 99.96% |
| 5 | Supervised Machine Learning Models for Liver Disease Risk Prediction | 2020 | Elias Dritsas, Maria Trigka | Liver Disease Risk | Decision Trees, Random Forest, Voting Classifier, SMOTE | Precision, Recall, Fmeasure, Accuracy, AUC | Voting Classifier demonstrated superior performance |
| 6 | Improved liver disease prediction | 2024 | Quadir et al. | Liver Disease | Gradient Boosting (GB), Extreme Gradient | Testing Accuracy | Extra Trees classifier achieved highest |

| | | | | | Boosting (XGB), Bagging, Random Forest (RF), Extra Trees (ET), Stacking, SMOTE | | accuracy of 91.82% |
|---|---|---|---|---|---|---|---|
| 7 | Prognosis of Liver Disease Using Machine Learning Algorithms | 2018 | Vyshali J. Gogi, Vijayalakshmi M.N | Liver Disease | Decision Tree, Linear Discriminant, Support Vector Machine (SVM), Logistic Regression, MATLAB 2016 | Accuracy, ROC | Linear Discriminant achieved highest prediction accuracy at 95.8% |

| 8 | Early-Stage Detection of Liver Disease Through Machine Learning Algorithms | 2021 | Krittika Dutta, Satish Chandra, Mahendra Kumar Gourisaria | Liver Disease | ANN, Logistic Regression, K-Nearest Neighbors (KNN), Support Vector Classifier (SVC), Gaussian Naive Bayes, Decision Trees, Random Forest, LR-SGD Classifier, PassiveAggressive, AdaBoost, Voting Classifier | Accuracy | The decision Tree algorithm demonstrated the highest accuracy at 99.96% |
| 9 | The Diagnosis of Chronic Liver Disease Using Machine Learning Techniques | 2021 | Golmei Shaheamlung, Harshpreet Kaur | Chronic Liver Disease | SVM, K-means clustering, KNN, Random Forest, Logistic Regression, Python with Spyder tool | Accuracy, Precision, Recall | Hybrid classification model achieved an accuracy of 77.58% |

| 10 | Prediction of Liver Diseases Based on Machine Learning Technique for Big Data | 2018 | Engy A. El-Shafeiy, Ali I. El-Desouky, Sally M. Elghamrawy | Liver Disease | SVM, Boosted C5.0, Naive Bayes (NB) | Accuracy, Sensitivity, Specificity | Evaluated based on accuracy, sensitivity, and specificity |
|---|---|---|---|---|---|---|---|
| 11 | The Diagnosis of Chronic Liver Disease Using Machine Learning Techniques | 2021 | Golmei Shaheamlung, Harshpreet Kaur | Chronic Liver Disease | SVM, K-means clustering, KNN, Random Forest, Logistic Regression, Python with Spyder tool | Accuracy, Precision, Recall | Hybrid classification model achieved an accuracy of 77.58% |
| 12 | A Comparison of Data Mining Algorithms for Liver | 2019 | Ain Najwa Arbain, B. Yushalinie Pillay Balakrishnan | Liver Disease | K-Nearest Neighbour (kNN), Logistic Regression, AutoNeural, Random Forest | Accuracy, ROC Index | k-NN algorithm achieved an accuracy of 99.794% |
| | Disease Prediction on Imbalanced Data | | | | | | |

| 13 | A Data Mining Approach to Prediction of Liver Diseases | 2020 | Nazim Razali, Aida Mustapha, Mohd Helmy Abd Wahab, Salama A Mostafa, Siti Khadijah Rostam | Liver Disease | NaÃ¯ve Bayes, Support Vector Machines (SVM), C4.5 Decision Tree, Multilayer Perceptron | Accuracy | Multilayer Perceptron algorithm achieved highest accuracy of 71.59% |
| 14 | Analysis of Data Mining Techniques for Healthcare Decision Support System Using Liver Disorder Dataset | 2016 | Tapas Ranjan Baitharu, Subhendu Kumar Pani | Liver Disease | J48 Decision Tree, Naive Bayes, Artificial Neural Network (ANN), ZeroR, 1BK, VFI | Data Classificatio n Accuracy, Computation al Time | Highlights importance of classificatio n accuracy and computation al time |
| 15 | A Survey on Machine Learning Techniques for the Diagnosis | 2020 | Golmei Shaheamlun g, Harshpreet Kaur, | Liver Disease | SVM, KNN, K-Means clustering, neural networks, decision trees, NaÃ¯ve Bayes | Accuracy, Sensitivity, Precision, Specificity | Hybrid approaches may improve predictive performance |
|  | of Liver Disease |  | Mandeep Kaur |  |  |  |  |

| 16 | A Survey on Machine Learning Techniques for the Diagnosis of Liver Disease | 2020 | Golmei Shaheamlung, Harshpreet Kaur, Mandeep Kaur | Liver Disease | SVM, KNN, K-Means clustering, neural networks, decision trees, NaÃ¯ve Bayes | Accuracy, Sensitivity, Precision, Specificity | Hybrid and ensemble methods can enhance predictive performance |
| 17 | Analysis of Data Mining Techniques for Healthcare Decision Support System Using Liver Disorder Dataset | 2016 | Tapas Ranjan Baitharu, Subhendu Kumar Pani | Liver Disease | J48 Decision Tree, Naive Bayes, Artificial Neural Network (ANN), ZeroR, 1BK, VFI | Effectiveness, Correction Rate | The choice of algorithm impacts the accuracy and efficiency of models |

## 2.3 Classification Techniques in Liver Diseases

Data from medical datasets are analyzed to find patterns and features suggestive of particular conditions, which are then used to detect and classify liver diseases using data mining and machine learning techniques. In this field, several categorization methods are frequently used, each having special advantages. An outline of some well-known techniques for identifying liver disease is provided here

**O Logistic Regression**

The probability of a binary outcome based on one or more predictor variables is modeled by logistic regression, a statistical technique used for binary classification. For datasets whose binary classification problem is the presence or absence of liver disease, it is straightforward, comprehensible, and efficient.

**O Decision Tree Classifier**

For binary classification tasks, like identifying whether liver disease is present or absent, the Decision Tree Classifier is a well-liked machine learning algorithm. This method uses a tree structure to represent decisions and their potential outcomes. An internal node represents an attribute test, a branch represents a test result, and a leaf node represents a class label. Decision Tree Classifier is useful for datasets with binary classification issues because of its clear-cut, intelligible, and interpretable design, which makes it useful for applications like liver disease detection. It effectively manages categorical and numerical data, offering decision rules that are easy for medical professionals to comprehend and visualize, supporting clinical decision-making.

**O Random Forest**

An extension of decision trees, this ensemble method constructs multiple decision trees during training and outputs the class that is the mode of the classes of the individual trees. Random forests are robust to overfitting and can handle large datasets with high dimensionality.

**O k- Nearest Neighbors**

For classification, this is a straightforward, non-parametric technique. A data point is classified by the model into the class that has the highest commonality among its k nearest neighbors, determined by the majority vote of those neighbors. Datasets with a clearly defined distance metric benefit the most from it.

**O Gradient Boosting Classifier**

Models are constructed successively using the ensemble technique known as GBC, whereby errors in earlier models are corrected in each new model. Focusing on cases that are challenging to classify improves predictive performance. When it comes to diagnosing liver disease, GBC is especially
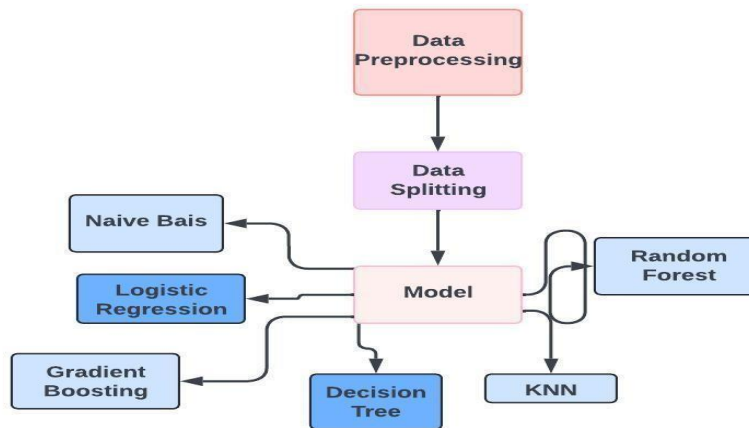
useful because it can spot subtle patterns and relationships in complex datasets that may be missed by simpler models.

## O Naive Bais

Based on Bayes' theorem and presuming feature independence, the Naive Bayes classifier is probabilistic. For medical diagnostics, in particular, where it can swiftly categorize data into groups such as the presence or absence of disease, it is very effective despite its simplicity. Early-stage diagnosis and screening procedures can benefit greatly from the computational efficiency and good performance of Naive Bayes, especially when dealing with small to medium-sized datasets.

## O Support Vector Machine (SVM)

SVMs are an effective tool for data classification because they can identify the best hyperplane in a high-dimensional space to divide data into classes. Because they use kernel methods, they are flexible and can support both linear and non-linear classification.



*Figure 3: A Comprehensive Method for Detecting Liver Disease Using Machine Learning*

## 2.4 Gaps in Existing Research

Despite advancements in liver disease detection using data mining and machine learning, several gaps remain in existing research.

### 2.4.1. Restricted Variety in the Dataset

The generalizability of models is limited by the fact that many studies rely on datasets that lack demographic diversity.

### 2.4.2. Unbalanced Datasets

The presence of healthy individuals in unbalanced datasets can result in biased models that have poor detection capabilities for liver disease cases.

### 2.4.3. Inadequate Comparative Analysis

 There isn't enough comparison research done between various machine learning algorithms. Research frequently concentrates on a single algorithm without comparing it to others.

### 2.3.4. Feature Engineering and Selection

Feature engineering and selection receive too little attention, which can lead to overfitting and decreased interpretability of the model**.**

### 2.4.5. Metrics for Evaluation

It can be deceptive to place too much emphasis on accuracy as the main indicator. More extensive measures such as ROC-AUC, F1-score, precision, recall, and recall are frequently underreported.

### 2.4.6. Real-world Applicability

Limited research has evaluated the models' practicality in the real world, taking into account issues with deployment and healthcare practitioners' ability to utilize them.

### 2.4.7. Techniques for Preparing Data

It is challenging to compare and replicate results due to variations in data preprocessing techniques used in different studies.

### 2.4.8. Comprehending and Manifestation

Interpretability is critical to clinical adoption, and deep learning and other complex models frequently lack it.

## 2.4.9. Analyzing Data Over Time

While longitudinal data can shed more light on the course of a disease, cross-sectional data is used in most studies.

## 2.4.10. Combining Molecular and Genetic Data

Few studies combine clinical and genetic data, which could improve the accuracy of diagnostic models.

**Chapter 3: Methodology**

**3.1 Research Design**

Utilizing data mining and machine learning techniques, the research design of this study is methodically organized to meet the goals of liver disease detection. With a focus on numerical data and statistical analysis, the study takes a quantitative approach to accurately and impartially assess the performance of different machine learning models. Through the use of data-driven insights, this approach guarantees that the models are evaluated and provides a solid basis upon which inferences can be made.

Experimental and analytical methods are combined in the methodology. By using a dataset on liver disease, the experimental approach trains several machine learning models and assesses how well they predict outcomes. A wide range of algorithms, including Gradient Boosting, Decision Trees, Random Forests, K-nearest neighbors, Logistic Regression, and Naive Bayes, are chosen in this process because of their track record of success in classification applications. Interpreting the outcomes of these tests to determine the best model and comprehend the variables affecting its performance is the analytical approach.
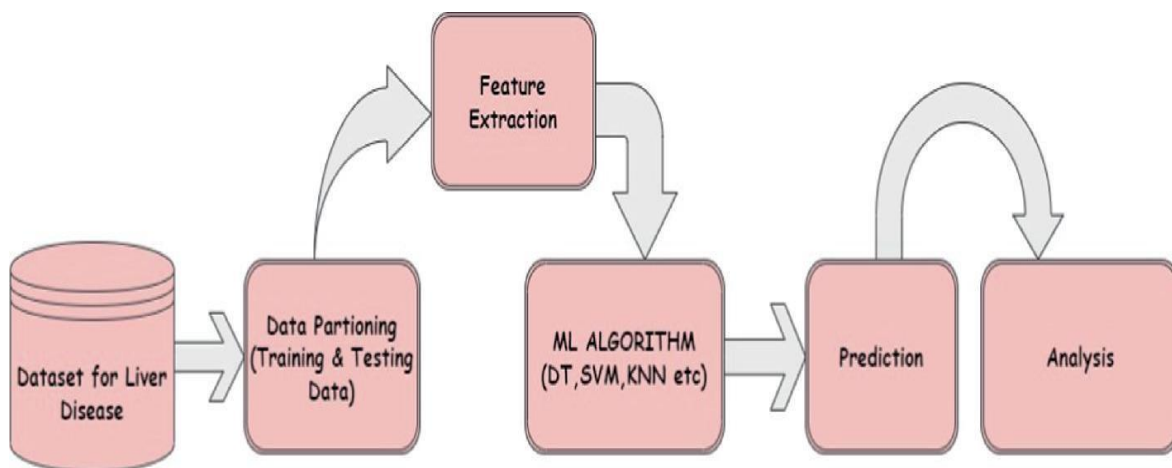
Defining the need for effective, non-invasive diagnostic techniques for liver disease detection is the first step in formulating the research problem, which is the research design. Next, particular goals are established, like comparing the effectiveness of different machine learning algorithms and determining the most accurate model.

An important first step in the research process is gathering data. The UCI Machine Learning Repository provides a publicly accessible dataset with several characteristics linked to liver health. The experimental analysis is built around this dataset. To ensure data quality, preprocessing is done on the data before training the models. This entails using imputation techniques to handle missing

values, normalizing data to guarantee consistency, and prioritizing the relevance of certain features to accurately assess the model's performance, the dataset is divided into training and testing sets during the model-training process. Methods like k-fold cross-validation are used to guarantee solid and trustworthy outcomes. A wide range of evaluation metrics, such as accuracy, precision, recall, F1score, and ROC-AUC, are used to evaluate each model's performance. These metrics offer a thorough comprehension of the advantages and disadvantages of each model.

After evaluating the evaluation metrics and visualizing the outcomes with charts and graphs, a comparative analysis is carried out to determine which algorithm is performing the best. After analyzing the practical ramifications and possible applications in liver disease detection in the real world, the results are interpreted to determine the elements that contributed to each model's success or failure.

Ultimately, a summary of the major discoveries and suggestions for additional research round out the study. As part of this, you can investigate different data mining and machine learning methods and make suggestions for possible model enhancements. This well-organized study approach guarantees a methodical and comprehensive exploration, offering dependable and practical discernments into the efficacy of various machine learning algorithms for the identification of liver disease.



*Figure 4: Process of Detecting Liver Disease through Machine Learning Algorithms*

**3.2 Dataset**

This research project will be using a dataset imported from the UCI website (https//archive.ics.uci.edu/dataset/225/ilpd+indian+liver+patient+dataset)  website. The intricacy of liver diseases and the particular needs of the investigation made choosing the right dataset for this research extremely difficult. Following a thorough search and assessment of several datasets for medical research, the choice was made to use a specific liver disease dataset designed to fulfill the requirements of this project. The focus on liver disease patients in this dataset sets it apart, as it provides comprehensive clinical and biochemical parameters necessary for an in-depth analysis.

The dataset includes several hundred individual records, each labeled with different parameters that are important to comprehending liver health, like albumin, gender, age, bilirubin, and liver enzyme levels, among other relevant clinical markers. The decision to use this dataset was influenced by its depth of data and well-organized structure, which made it especially helpful for statistical analyses and the development of predictive models targeted at identifying critical factors influencing the course of liver disease. By using this dataset, the study hopes to simplify the data processing and analysis phase while maintaining the reliability of the research methodology's performance evaluation. This is achieved by reducing the complexity that is often connected with larger, less focused datasets.

- **Age Distribution:** A roughly normal distribution can be seen in the age distribution histogram, suggesting that the majority of patients in the dataset are middle-aged or older, with fewer young and elderly patients. This distribution indicates that liver disease affects people of all ages, although middle-aged people are more likely to have it. The histogram's smooth curve overlay draws attention to the symmetry around the central age values even more.
- **Total Bilirubin Distribution(Total bilirubin level):** The peak near the lower end of the scale indicates that most patients have total bilirubin levels that are primarily low. This graph, which is right-skewed, shows that although the majority of patients have bilirubin levels that fall within the normal range, a sizable portion of cases have elevated levels, which may indicate that the patients have different disease etiologies or severe liver conditions.
- **Alkaline Phosphatase Distribution(Level of alkaline phosphatase)** The distribution of alkaline phosphatase levels in this histogram is right-skewed, which is typical of clinical data and occurs when many observations are concentrated close to the lower range and fewer

30

observations extend towards higher values. This pattern might be the result of a standard variation in a population of healthy individuals combined with patients who have liver dysfunction or other related health problems and show noticeably higher enzyme levels.

○ **Albumin Distribution(Albumin level)** There is a slight skew to the right in the distribution of albumin levels, but overall they resemble a normal curve. As would be expected in a carefully monitored clinical setting, this shows that the majority of patients have albumin levels within a rather narrow range. Albumin is a crucial indicator of liver function and protein synthesis, so the existence of some outliers on the higher end may be a sign of different nutritional statuses or stages of liver disease.
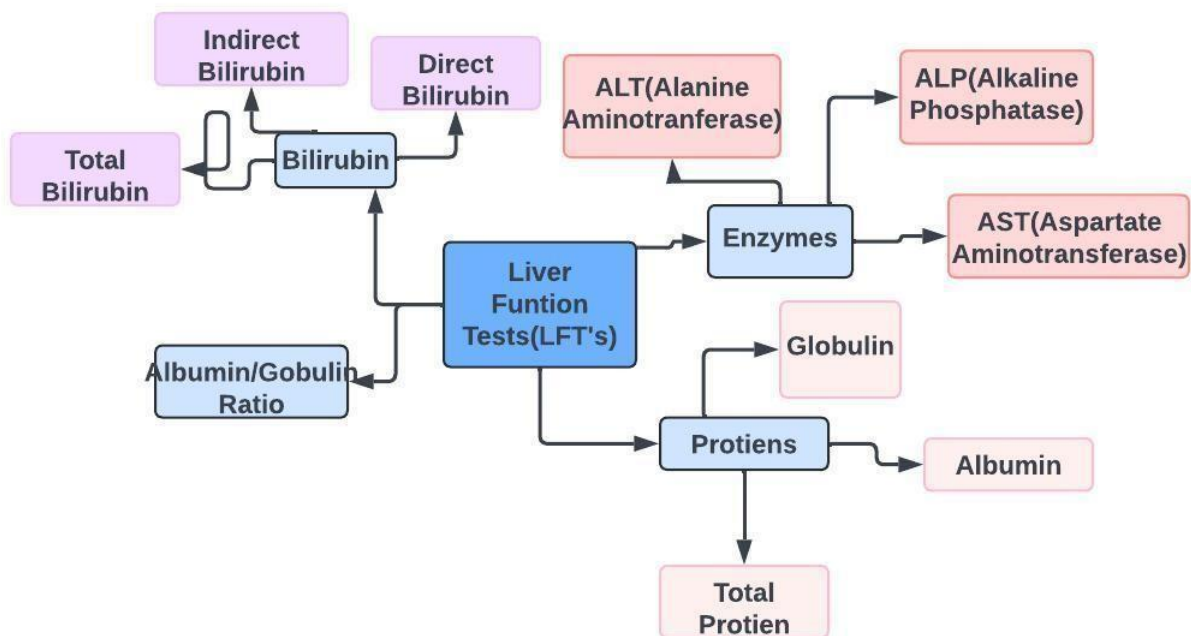


*Figure 5:Liver Function Test (LFT) Components and Relationships*

## 3.3 Data Preprocessing

To prepare the dataset for analysis and model training, data preprocessing is necessary. This includes splitting the dataset into training and testing sets, normalizing numerical features, handling missing

values, and encoding categorical variables. A thorough description of the data preprocessing procedures is provided below.

### 3.3.1. Bringing the Dataset Up

`pandas` is used to load the dataset, reading the CSV file into a DataFrame for simple manipulation and analysis.

- **Looking for Any Missing Data**

  The amount of missing data is calculated by looking for missing values in each column of the dataset. This step helps identify which columns have missing data and how many entries from each column are missing.

- **Imputing Missing Values**

The mean value of the accessible data in the 'Albumin_Globulin_Ratio' column is used to impute missing values. The dataset's integrity is preserved without sacrificing a significant amount of information by imputing missing values using the mean.

- **Verification**

The dataset is reexamined after imputation to make sure no more missing values exist. This verification process ensures that any missing data has been handled correctly.

### 3.3.2. Encoding Categorical Variables

- **Label Encoding**

Using label encoding, the categorical variable 'Gender' is transformed into a numerical representation. Because categorical data must be processed by machine learning algorithms, this step is required. To speed up machine learning, label encoding gives each category a distinct integer.

### 3.3.3. Data Normalization

- **Identifying Numerical Features**

Once the target variable 'Liver_Disease' and unnecessary columns like 'Patient_ID' are removed, numerical features are obtained.

## ⭕ Scaling Features

Using a standard scaler, numerical features are scaled to a common range. By ensuring that every feature makes an equal contribution to the model training process, normalization enhances both the model's convergence speed and performance. Standard scaling makes the features comparable by centering the data around zero with a unit standard deviation.

### 3.3.3. Defining Features and Target Variable

By defining the feature set (X) and the target variable (y), the dataset is ready for model training. Since it does not contribute to the prediction, the 'Patient_ID' column is removed from the feature set. All other pertinent columns make up the feature set, and the 'Liver_Disease' column is the target variable.

By following these preprocessing steps, the dataset is transformed into a clean and structured format, ready for training machine learning models. Handling missing values, encoding categorical variables, and normalizing numerical features are critical steps that enhance the accuracy and reliability of the models. This comprehensive preprocessing ensures that the data is of high quality, leading to more accurate and generalizable models for liver disease detection.

### 3.4 Training and Testing Data

Two essential datasets are needed to train a machine learning model a test set to assess the model's performance and a training set to build the model. The existing data must be split into these two different sets to guarantee a solid evaluation and avoid overfitting. The liver disease dataset, which was methodically divided into training and test sets, was used in this study.

To prevent inconsistent data handling, the dataset was first cleaned by making sure that there were no unnecessary spaces in the column names. 'Liver_Disease' was the target variable (y), and the feature set (X) was defined by removing the 'Patient_ID' and 'Liver_Disease' columns. As a feature that serves as an identifier rather than a component of the predictive model, the 'Patient_ID' was eliminated from the feature set.

The dataset was divided into two portions 70% for training and 30% for testing, to evaluate the model's capacity for generalization. To ensure the reproducibility of the results and consistency in model training and evaluation, this split was carried out using a random state.

| Data Set | Number of Records | Percentage of Total Data |
|----------|-------------------|--------------------------|
| Training | 405 | 70% |
| Testing | 174 | 30% |
| **Total** | 579 | **100%** |

*Table 2: Table of Data Distribution*

The machine learning model was trained using the training set, which included 405 records so that it could discover the underlying patterns and relationships in the data. With 174 records, the test set functioned as a separate assessment dataset to gauge the model's performance on unobserved data.

This methodology guarantees that the predictive powers of the model are verified on an independent dataset, offering a precise assessment of its capacity to generalize to novel situations. The split datasets made it easier to proceed with the model's development and evaluation by being saved into distinct CSV files.

## 3.5 Machine Learning Techniques

In Machine Learning Methodologies, the following techniques can be used and explained below with their strengths and Weaknesses

### ◯ Logistic Regression

A linear model for binary classification issues is called logistic regression. It models the binary outcome with a logistic function and calculates the likelihood that a given instance belongs to a specific class.

**Strengths**

- Easy to use and comprehend.
- The model's coefficients can be understood.
- Effective with data that is linearly separable.

**Weaknesses**

- Sensitive to outliers.

- assumes a linear relationship between features and log odds.

## ⭕ Decision Tree

A decision tree is a kind of non-linear model that divides the data into subsets according to the input feature values. It employs a tree-like structure in which a feature test is represented by each internal node, the test result is represented by each branch, and the class label is represented by each leaf node.

**Strengths**

- Simple to understand and depict.
- Does not require much data preprocessing
- it can handle both numerical and categorical data.

**Weaknesses**

- Prone to overfitting.
- sensitive to slight variations in the data.

## ⭕ Random Forest

Using the Random Forest ensemble learning technique, several decision trees are constructed and then combined to produce predictions that are more reliable and accurate. By averaging the outcomes of several trees, it lessens overfitting**.**

**Strengths**

- Handles large datasets with higher dimensionality
- lessens overfitting when compared to individual decision trees.
- Resistant to noise and outliers.

**Weaknesses**

- Harder to understand than a single decision tree.
- High computational demands.

### O K-Nearest Neighbors (KNN)

A straightforward, non-parametric method called K-Nearest Neighbors is employed for both regression and classification applications. A data point is categorized according to how its neighbors are categorized.

**Strengths**

- Easy to use and intuitive.
- No training period is necessary.
- Functions well with minimal data.

**Weaknesses**

- Sensitive to the selection of the number of neighbors (k),
- computationally costly during prediction.
- Has trouble managing big datasets.

### O Gradient Boosting

Using an ensemble approach, gradient boosting builds models one after the other, fixing mistakes in the previous models. It creates a strong learner by combining weak learners.

**Strengths**

- Capable of handling a range of loss functions
- High predictive accuracy.
- Performance-enhancing for both regression and classification tasks.

**Weaknesses**

- Needs careful hyperparameter tuning;
- Is computationally intensive
- Is prone to overfitting if not properly tuned.

**⭕ Naive Bayes**

Based on the Bayes Theorem, the Naive Bayes classifier is probabilistic and assumes feature independence. It works especially well with text classification and categorical feature issues.

**Strengths**

- Easy to use and quick.
- Effectively manages high-dimensional data.
- Functions effectively with little data.

**Weaknesses**

- May not function well with highly correlated features
- assumes independence between features.

**3.6 Training the Models**

The preprocessed dataset is used to train the models. The dataset is divided 70/30 into training and testing sets so that the models are tested on data that hasn't been seen before.

- **Data Preparation** Thirty percent of the dataset is tested, and seventy percent is used for training. By evaluating the model's performance on hypothetical data, this split guarantees an accurate assessment of the model's capacity for generalization.
- **Model Training** Using the training set, every algorithm is trained. Optimizing the model's parameters to reduce error on the training set is a step in the training process. For instance, logistic regression determines the optimal model parameters by utilizing maximum likelihood estimation.
- **Evaluation Matrices**
- A wide range of assessment metrics are used to evaluate each model's performance. Every metric offers a distinct viewpoint on the effectiveness of the model.
- **Accuracy** A measure of accuracy is the percentage of cases out of all instances that were correctly predicted.
- **Precision** Measures the percentage of true positive predictions among all positive predictions.

- **Recall** Measures recall by taking the percentage of actual positive instances divided by the number of true positive predictions.

- **F1-Score** The F1-score offers a balance between recall and precision by taking the harmonic mean of the two.

- **ROC-AUC** Calculates the area under the Receiver Operating Characteristic curve, which shows how well the model can categorize data.

## 3.7 Tools and Software Used

Model training, evaluation, and data preprocessing were all done with the aid of the tools and software listed and described in this section of the study. The selection of instruments and applications is essential to guaranteeing a productive process and dependable outcomes.

- **Python** The principal language of programming used to implement the algorithms for machine learning. Python is a preferred language for data science and machine learning tasks due to its simplicity, readability, and wide library support.

- **Libraries**

   **Pandas** For analyzing and manipulating data. It offers the data structures and operations required for effectively cleaning and processing data.

   **Scientific Kit-Learn (Sklearn)** a thorough library of machine learning tools that covers model validation, assessment, and training. It is compatible with many different preprocessing methods and algorithms.

   **Matplot** A plotting library called matplotlib is used to produce static, interactive, and animated visualizations. It is especially helpful for visualizing model performance, outcomes, and data distributions.

- **The environment for integrated development (IDE)**

   **Spyder** The Scientific Python Development Environment, or Spyder, is a potent IDE for Python scientific programming. It offers an interactive development environment with features for introspection, interactive testing, debugging, and advanced editing. Spyder is frequently used in machine learning and data science projects, and it easily integrates with well-known Python libraries.

The robustness and user-friendliness of these tools and software allowed for the effective implementation and assessment of machine learning models. From data preprocessing to model training and evaluation, they guarantee a seamless workflow. Spyder is the perfect IDE for this study because of its interactive features and easy integration with important libraries, which increase productivity.

## Chapter 4:  Results

The outcomes of using different machine learning algorithms on the liver disease dataset are presented in this chapter. Each model's outcomes, along with performance indicators and visualizations, are comprehensive. We talk about how these findings affect the diagnosis of liver disease. The best model is finally determined by a comparative analysis, which is followed by a breakdown of the advantages and disadvantages of each model.

### 4.1. Logistic Regression

This section includes the performance metrics and confusion matrix from the Logistic Regression model applied to the liver disease dataset, along with a discussion of the findings.

```
# 1 Logistic Model
logistic_mdl=LogisticRegression(random_state=42)
logistic_mdl.fit(X_train,y_train)

# Make predictions and evaluate
y_pred_log=logistic_mdl.predict(X_test)
print("Logistic Regression Results:")
print(f'Accuracy:{accuracy_score(y_test,y_pred_log)*100:.2f}%')
print(f'Precision:{precision_score(y_test,y_pred_log,zero_division=0)*100:.2f}%')
print(f'Recall: {recall_score(y_test, y_pred_log)*100:.2f}%')
print(f'F1-Score: {f1_score(y_test, y_pred_log)*100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_log)*100:.2f}%\n')

import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

# Calculate the confusion matrix
confu_matrix=confusion_matrix(y_test, logistic_mdl.predict(X_test), labels=[0, 1])

# Display the confusion matrix
disp=ConfusionMatrixDisplay(confusion_matrix=confu_matrix,display_labels=['No Liver Disease', 'Liver Disease'])
disp.plot(cmap=plt.cm.Blues)
plt.title('Confusion Matrix for Logistic Regression')
plt.show()
```

*Figure 6: Input code for Logistic Regression*

## 4.1.1 Training and Assessment of Models

Thirty percent of the dataset was used for testing and seventy percent was used for training the logistic regression model. Metrics including accuracy, precision, recall, F1-score, and ROC-AUC were used to assess the model's performance.
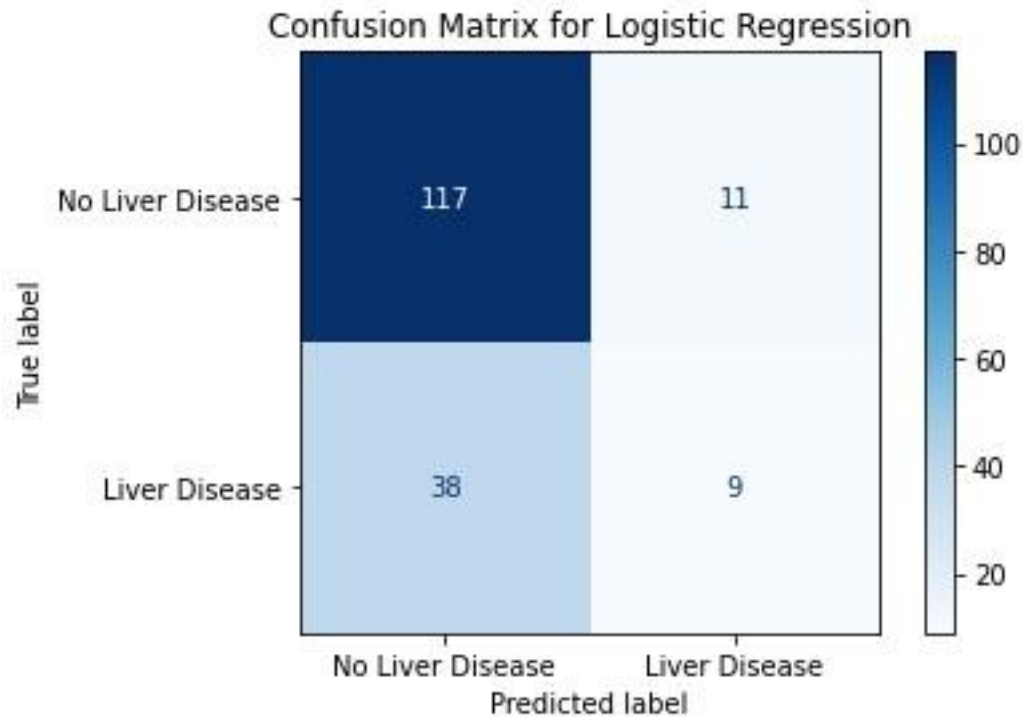
```
Logistic Regression Results:
Accuracy:72.00%
Precision:45.00%
Recall: 19.15%
F1-Score: 26.87%
ROC-AUC: 55.28%
```

*Figure 7: Output for Logistic Regression*

## 4.1.2 Result

According to the findings, the accuracy of the Logistic Regression model was 72.00%. Its recall and precision are 19.15% and 45.00%, respectively. This implies that although the model can fairly well identify cases, it has a very difficult time identifying actual positive cases of liver disease.

40

The model correctly detects 9 cases of liver disease and 117 cases of no liver disease, according to the confusion matrix. But it also misdiagnoses 11 cases of liver disease and ignores 38 actual cases of liver disease.



*Figure 8: Confusion Matrix for Logistic Regression*

The model of Logistic Regression attains a moderate level of accuracy (72.0%), but its precision (45.00%) and recall (19.15%) are low. With numerous false negatives (38) and false positives (11), this suggests serious misclassification problems. While the model's readability and simplicity are advantages, its linear structure makes it less useful for identifying intricate patterns in the data. More sophisticated models are required to detect liver disease more accurately.

**4.2. Decision Tree Classifier**

Performance metrics and a confusion matrix are included in this section that displays the outcomes of the Decision Tree model applied to the liver disease dataset.

```
#2.Decision Tree Classifier

from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

decision_tree_model=DecisionTreeClassifier(random_state=42)
decision_tree_model.fit(X_train, y_train)

# Make predictions and evaluate
y_pred_tree=decision_tree_model.predict(X_test)
print("Decision Tree Results:")
print(f' Accuracy: {accuracy_score(y_test, y_pred_tree) * 100:.2f}%')
print(f'  Precision: {precision_score(y_test, y_pred_tree, zero_division=0) * 100:.2f}%')
print(f'  Recall: {recall_score(y_test, y_pred_tree) * 100:.2f}%')
print(f'  F1-Score: {f1_score(y_test, y_pred_tree) * 100:.2f}%')
print(f'  ROC-AUC: {roc_auc_score(y_test, y_pred_tree) * 100:.2f}%\n')

import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

# Calculate the confusion matrix
confu_matrix=confusion_matrix(y_test,decision_tree_model.predict(X_test),labels=[0,1])

# Display the confusion matrix
display=ConfusionMatrixDisplay(confusion_matrix=confu_matrix,display_labels=['No Liver Disease', 'Liver Disease'])
disp.plot(cmap=plt.cm.Blues)
plt.title('Confusion Matrix for Decision Tree')
plt.show()
```

*Figure 9: Input Code for Decision Tree Classifier*

## 4.2.1 Training and Assessment of Models

30% of the dataset was used to test the Decision Tree model after it had been trained using 70% of it. The metrics of performance.

```
Decision Tree Results:
 Accuracy: 69.14%
  Precision: 43.40%
  Recall: 48.94%
  F1-Score: 46.00%
  ROC-AUC: 62.75%
```

*Figure 10: Output for Decision Tree Classifier*

**4.2.2 Result**

The confusion matrix indicates that 9 cases of liver disease and 117 cases of no liver disease are correctly detected by the model. Nevertheless, it also fails to identify 38 actual cases of liver disease and misdiagnoses 11 cases as having liver disease.
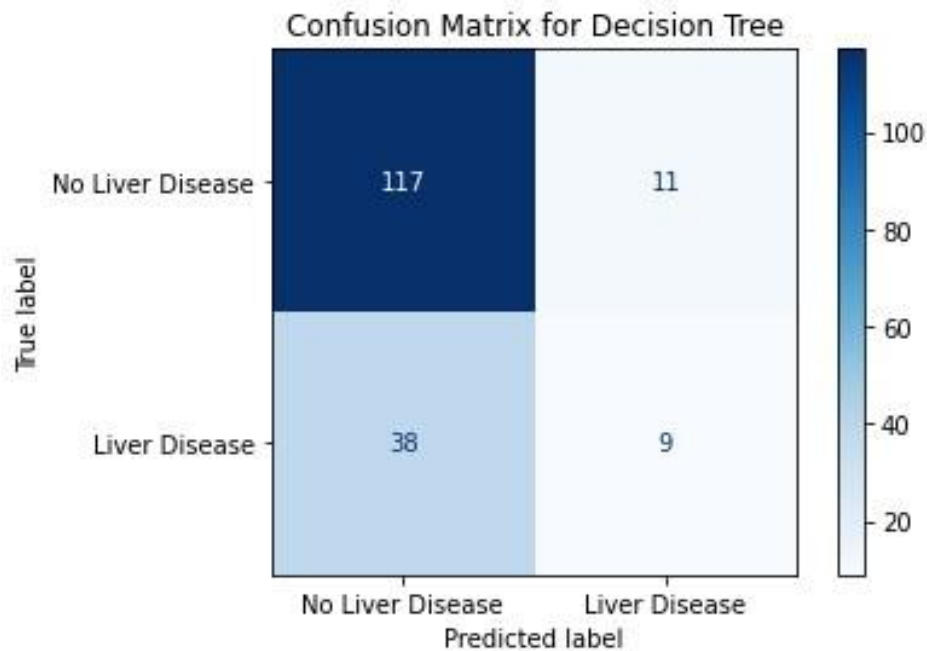


*Figure 11: Confusion Matrix for Decision Tree*

The Decision Tree model yielded a recall of 48.94%, precision of 43.40%, and accuracy of 69.14%. The model's moderate ability to distinguish between classes is indicated by its ROC-AUC score of 62.75%. The confusion matrix indicates that a significant amount of false positives (11) and false negatives (38) are typically produced by the model.

Both numerical and categorical data can be handled by the Decision Tree model, which also requires little data preprocessing and is simple to understand and visualize. However, because of its sensitivity to small changes in the data and tendency toward overfitting, it produces moderate performance metrics that suggest it may not fully capture all the complexities in the dataset.

**4.3. Random Forest.**

Performance metrics and a confusion matrix are included in this section that displays the outcomes of the Random Forest model used on the liver disease dataset.

```python
#3 Random Forest Classifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
forest_model = RandomForestClassifier(random_state=42)
forest_model.fit(X_train, y_train)

# Make predictions and evaluate
y_pred_forest=forest_model.predict(X_test)
print("Random Forest Results:")
print(f'Accuracy: {accuracy_score(y_test, y_pred_forest)*100:.2f}%')
print(f'Precision: {precision_score(y_test, y_pred_forest, zero_division=0)*00:.2f}%')
print(f'Recall:{recall_score(y_test, y_pred_forest)*100:.2f}%')
print(f'F1-Score:{f1_score(y_test, y_pred_forest)*100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_forest) * 100:.2f}%\n')
confu_matrix_forest=confusion_matrix(y_test,forest_model.predict(X_test),labels=[0, 1])
disp_forest = ConfusionMatrixDisplay(confusion_matrix=confu_matrix_forest, display_labels=['No Liver Disease', 'Liver Disease'])
disp_forest.plot(cmap=plt.cm.Greens)
plt.title('Confusion Matrix for Random Forest')
plt.show()

from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
```

*Figure 12:Input Code for Random Forest*

## 4.3.1 Training and Assessment of Model

30% of the dataset was used to test the Random Forest model after it had been trained using 70% of it. The metrics for performance that are

```
Random Forest Results:
Accuracy: 72.57%
Precision: 0.00%
Recall:34.04%
F1-Score:40.00%
ROC-AUC: 60.38%
```

*Figure 13: Output for Random forest*

## 4.3.2 Result

16 instances of liver disease and 111 cases of no liver disease are correctly detected by the model, according to the confusion matrix. In addition, it misdiagnoses liver disease in 17 cases while missing 31 actual cases.
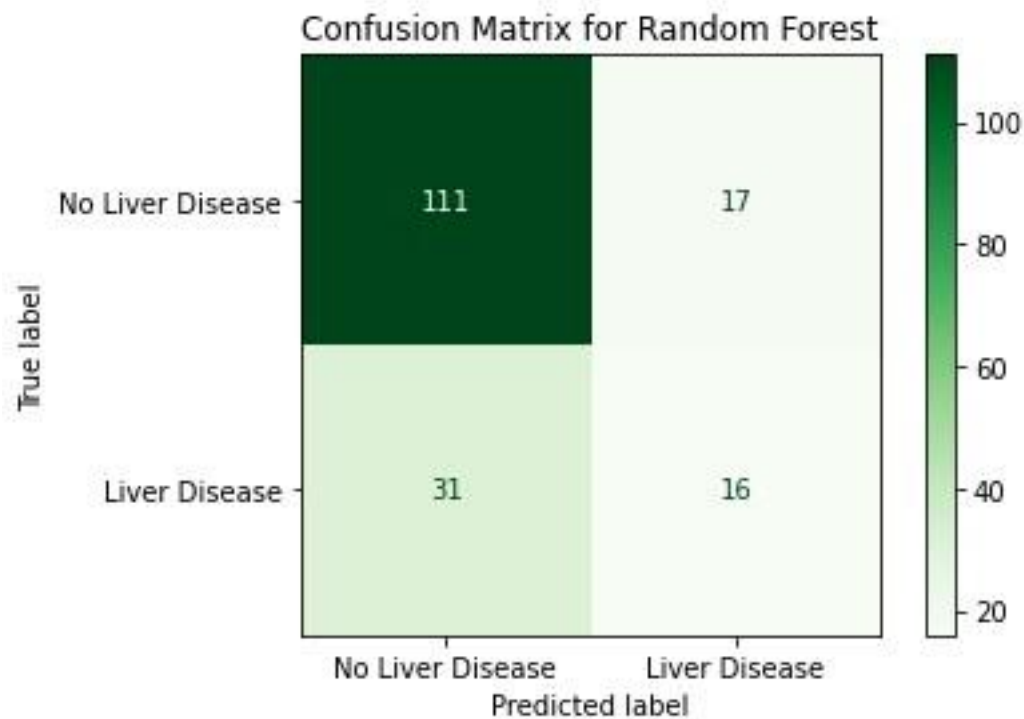


*Figure 14: Confusion Matrix for Random Forest*

The Random Forest model yielded a 72.57% accuracy rate, 0.00% precision, and 34.04% recall. The model's moderate capacity to discriminate between classes is indicated by its ROC-AUC score of 60.38%. The confusion matrix reveals that a significant proportion of false positives (17) and false negatives (31) are generated by the model.

Large datasets with higher dimensionality can be handled by the Random Forest model, which is renowned for its resilience. In comparison to individual decision trees, it lessens overfitting and is noise- and outlier-resistant. The model's low precision suggests that it has difficulty correctly identifying true positive cases of liver disease, which could be attributed to the dataset's imbalance.

## 4.4.  K-Nearest Neighbors

Performance metrics and a confusion matrix are included in this section that showcase the outcomes of the K-Nearest Neighbors model on the liver disease dataset.

```
#4 KNN
knn_model=KNeighborsClassifier(n_neighbors=5)
knn_model.fit(X_train,y_train)

# Make predictions and evaluate
y_pred_knn=knn_model.predict(X_test)
print("K-Nearest Neighbors Results:")
print(f'Accuracy: {accuracy_score(y_test, y_pred_knn)*100:.2f}%')
print(f'Precision: {precision_score(y_test, y_pred_knn, zero_division=0)*100:.2f}%')
print(f'Recall: {recall_score(y_test, y_pred_knn)*100:.2f}%')
print(f'F1-Score: {f1_score(y_test, y_pred_knn)*100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_knn)*100:.2f}%\n')
confu_matrix_knn=confusion_matrix(y_test,knn_model.predict(X_test),labels=[0, 1])
disp_knn=ConfusionMatrixDisplay(confusion_matrix=confu_matrix_knn,display_labels=['No Liver Disease', 'Liver Disease'])
disp_knn.plot(cmap=plt.cm.Oranges)
plt.title('Confusion Matrix for K-Nearest Neighbors')
plt.show()
```

*Figure 15:Input Code for KNN*

## 4.4.1 Training and Assessment of Model

The confusion matrix indicates that 19 cases of liver disease and 96 cases of no liver disease are correctly detected by the model. But it also misdiagnoses 32 cases as liver disease and ignores 28 actual cases of liver disease.

```
K-Nearest Neighbors Results:
Accuracy: 65.71%
Precision: 37.25%
Recall: 40.43%
F1-Score: 38.78%
ROC-AUC: 57.71%
```

*Figure 16:Output for KNN*

## 4.4.2 Result

The confusion matrix indicates that 19 cases of liver disease and 96 cases of no liver disease are correctly detected by the model. However, it also fails to identify 28 actual cases of liver disease and misdiagnoses 32 cases as having liver disease.
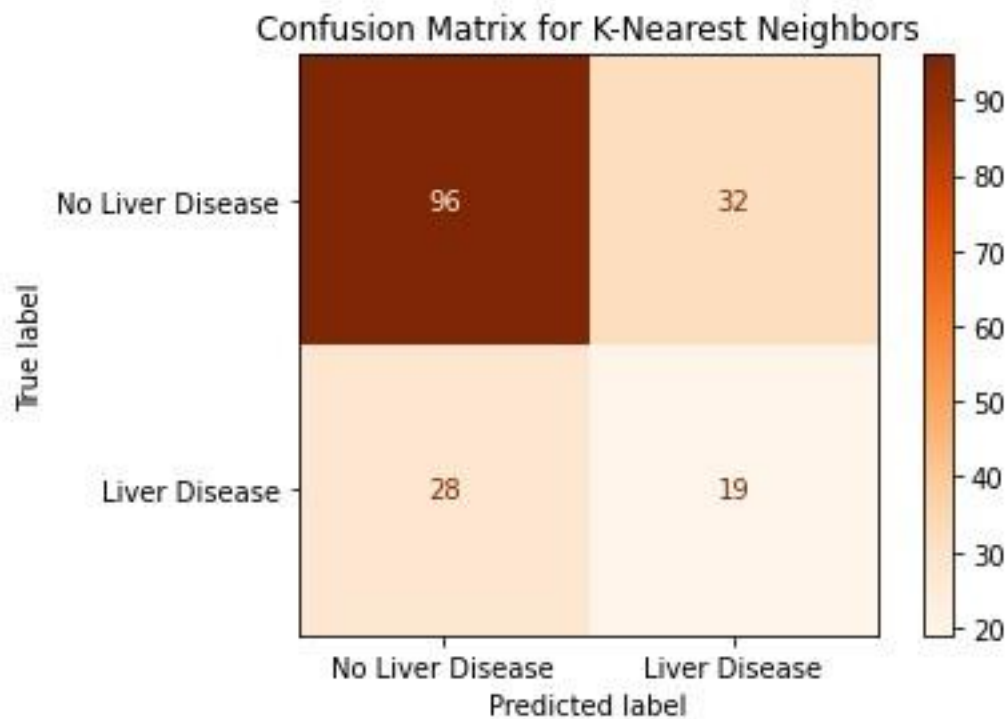


*Figure 17:Confusion Matrix for KNN*

With a precision of 37.25%, recall of 40.43%, and accuracy of 65.71%, the KNN model performed well. The model's limited capacity to distinguish between classes is indicated by its 57.71% ROCAUC score. The confusion matrix reveals that a significant proportion of false positives (32), as well as false negatives (28), are generated by the model.

With small datasets, the KNN model works well and doesn't require a training phase. It is also easy to understand and apply. It cannot handle large datasets well, is sensitive to the choice of the number of neighbors (k), and is computationally expensive during prediction.

To sum up, the KNN model offers a simple method of classification; however, due to its moderate performance, it may be better off with additional hyperparameter tuning or combined with other models to increase robustness and accuracy. To find a more reliable answer, additional research using different machine learning algorithms is required.

## 4.5. Gradient Boosting Classifier

Results of the Gradient Boosting model applied to the liver disease dataset are shown in this section, along with a confusion matrix and performance metrics.

```
#5 Gradient Boosting Classifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

# Initialize and train the model
gb_model=GradientBoostingClassifier(random_state=42)
gb_model.fit(X_train,y_train)

# Make predictions and evaluate
y_pred_gb = gb_model.predict(X_test)
print("Gradient Boosting Results:")
print(f'Accuracy: {accuracy_score(y_test, y_pred_gb) * 100:.2f}%')
print(f' Precision: {precision_score(y_test, y_pred_gb, zero_division=0) * 100:.2f}%')
print(f'Recall: {recall_score(y_test, y_pred_gb) * 100:.2f}%')
print(f'F1-Score: {f1_score(y_test, y_pred_gb) * 100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_gb) * 100:.2f}%\n')
confu_matrix_gb= confusion_matrix(y_test,gb_model.predict(X_test),labels=[0, 1])
disp_gb=ConfusionMatrixDisplay(confusion_matrix=confu_matrix_gb, display_labels=['No Liver Disease', 'Liver Disease'])
disp_gb.plot(cmap=plt.cm.Purples)
plt.title('Confusion Matrix for Gradient Boosting')
plt.show()
```

*Figure 18:Input Code for Gradient Boosting Classifier*

## 4.5.1. Training and Assessment

30% of the dataset was used to test the Gradient Boosting model after it had been trained using 70% of it. The metrics for performance are

*Figure 19:Output for Gradient Boosting*

## 4.5.2. Results

17 instances of liver disease and 107 cases of no liver disease are correctly detected by the model, according to the confusion matrix. Even so, it misdiagnoses liver disease in 21 cases and overlooks 30 actual cases of the condition.
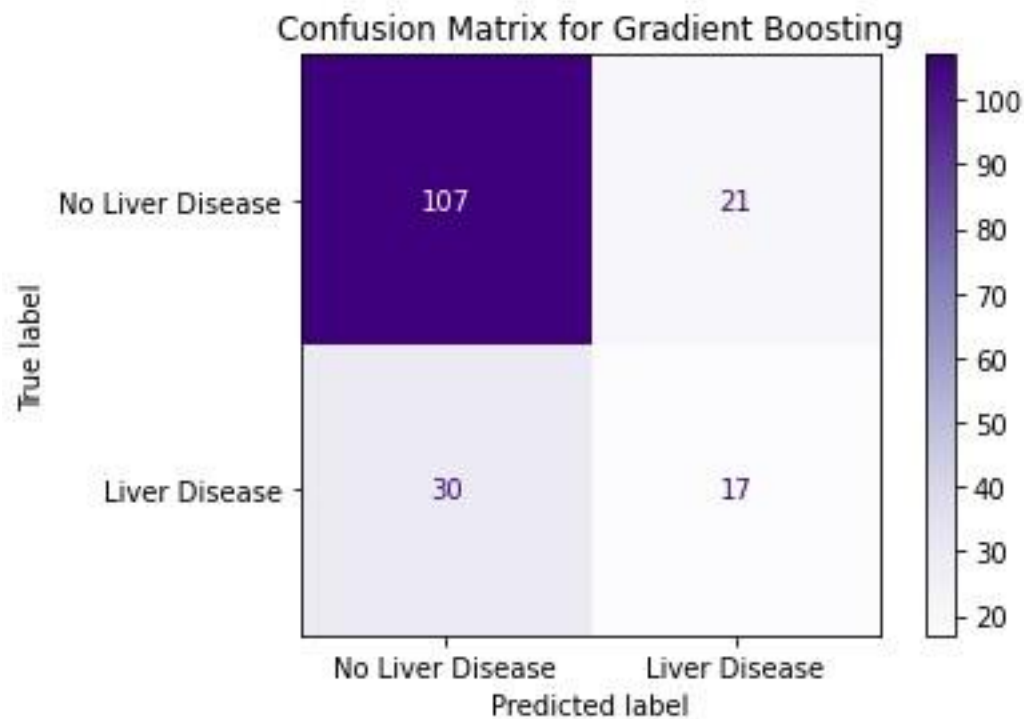


*Figure 20:Confusion Matrix for Gradient Boosting*

The Gradient Boosting model yielded a 70.86% accuracy rate, 44.74% precision, and 36.17% recall. The model's 59.88% ROC-AUC score indicates that it can distinguish between classes with a moderate degree of accuracy. The confusion matrix reveals that the model frequently generates a sizable number of false positives (21), as well as false negatives (30).

Gradient Boosting is well-known for its wide range of loss function support and high predictive accuracy. By successively building models and fixing the mistakes of the previous models, it is possible to achieve better performance than with individual models. It is computationally demanding and, if improperly adjusted, may be prone to overfitting.

In conclusion, even though the Gradient Boosting model offers a reliable method for classification, its mediocre performance raises the possibility that its accuracy and robustness could be increased by adjusting its hyperparameters further or by combining it with other models. To find a more reliable answer, additional research using different machine learning algorithms is required.

## 4.6 Naive Bayes

A confusion matrix and performance metrics are included in this section that show the outcomes of the Naive Bayes model applied to the liver disease dataset.

```python
#6 Naive Bayes
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

# Initialize and train the model
nb_model=GaussianNB()
nb_model.fit(X_train,y_train)

# Make predictions and evaluate
y_pred_nb=nb_model.predict(X_test)
print("Naive Bayes Results:")
print(f'Accuracy: {accuracy_score(y_test, y_pred_nb)*100:.2f}%')
print(f'Precision: {precision_score(y_test, y_pred_nb, zero_division=0)*100:.2f}%')
print(f'Recall: {recall_score(y_test, y_pred_nb)*100:.2f}%')
print(f'F1-Score: {f1_score(y_test, y_pred_nb)*100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_nb)*100:.2f}%\n')
confu_matrix_nb = confusion_matrix(y_test, nb_model.predict(X_test),labels=[0,1])
disp_nb = ConfusionMatrixDisplay(confusion_matrix=confu_matrix_nb,display_labels=['No Liver Disease', 'Liver Disease'])
disp_nb.plot(cmap=plt.cm.Reds)
plt.title('Confusion Matrix for Naive Bayes')
plt.show()
```
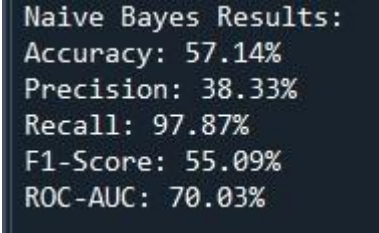
*Figure 21:Input Code for Naive Bayes*

## 4.6.1. Training and Assessment

Thirty percent of the dataset was used to test and train the Naive Bayes model. The performance indicators that are



```
Naive Bayes Results:
Accuracy: 57.14%
Precision: 38.33%
Recall: 97.87%
F1-Score: 55.09%
ROC-AUC: 70.03%
```

*Figure 22: Output for Naive Bayes*

## 4.6.2. Results

46 cases of liver disease and 54 cases of no liver disease are correctly detected by the model, according to the confusion matrix. On the other hand, it misdiagnoses 74 cases as liver disease and misses 1 actual case.
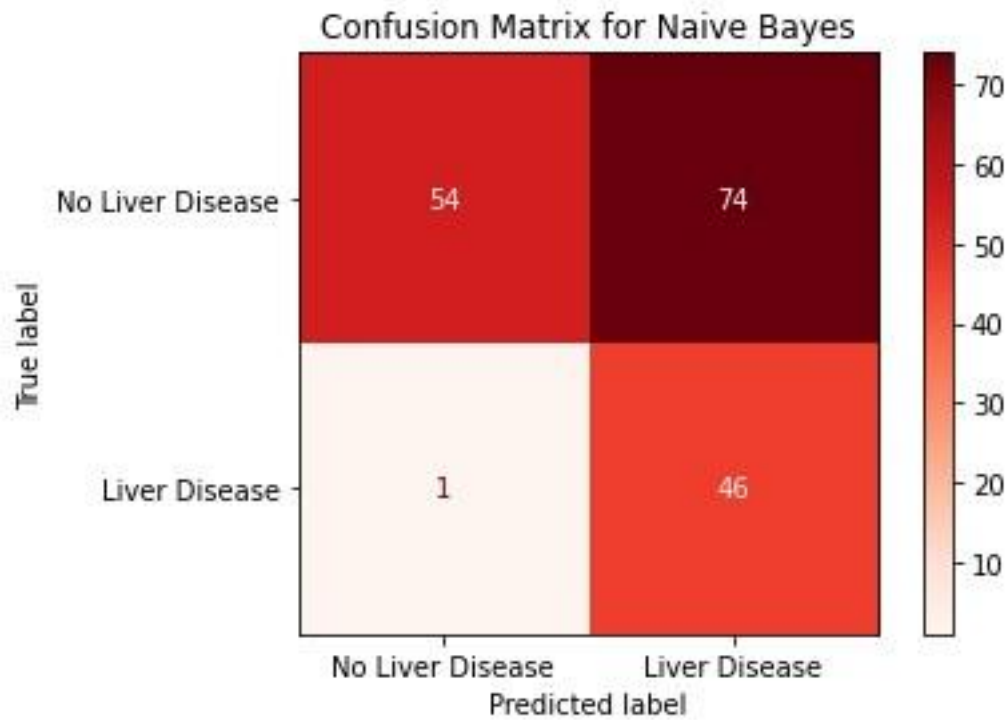
*Figure 23: Confusion Matrix for Naive Bayes*

With a precision of 38.33% and a recall of 97.87%, the Naive Bayes model produced an accuracy of 57.14%. With an ROC-AUC score of 70.03%, the model appears to be reasonably effective at differentiating between classes. The confusion matrix shows that very few false negatives (1) and a large number of false positives (74) are typically produced by the model.

A straightforward and quick probabilistic classifier, Naive Bayes works well with small amounts of data and can handle high-dimensional data. Nevertheless, it relies on the assumption of feature independence, which may not hold for this dataset and result in a high rate of false positives.

## 4.7. Comparative Analysis of Models

All of the machine learning models that have been used with the liver disease dataset are compared in this section. The accuracy, precision, recall, F1-score, and ROC-AUC are the metrics used to assess each model's performance. A visual comparison of these metrics between various models can be seen in the bar chart below.
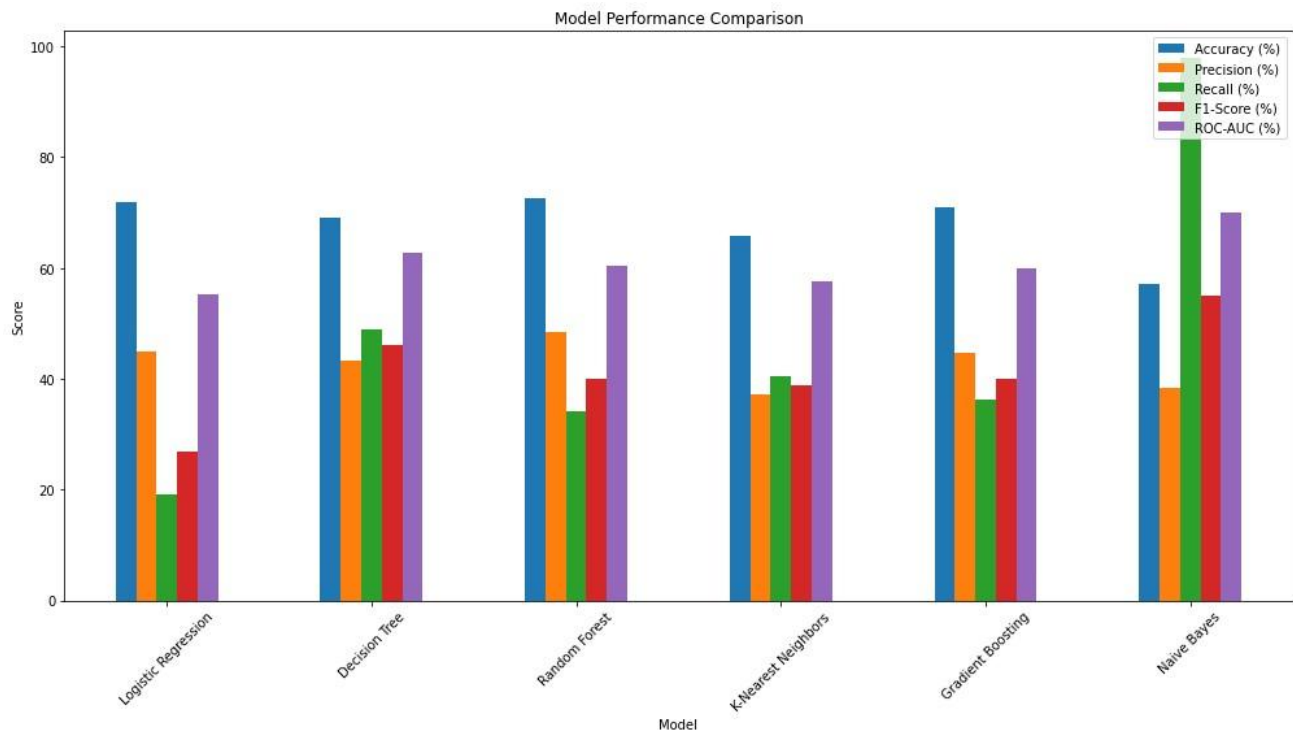
| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) | ROC-AUC (%) |
|---|---|---|---|---|---|

| Logistic Regression | 72 | 45 | 19.15 | 26.87 | 55.28 |
| --- | --- | --- | --- | --- | --- |
| Decision Tree | 69.14 | 43.4 | 48.94 | 46 | 62.75 |
| Random Forest | 72.57 | 48.48 | 34.04 | 40 | 60.38 |
| K-Nearest Neighbors | 65.71 | 37.25 | 40.43 | 38.78 | 57.71 |
| Gradient Boosting | 70.86 | 44.74 | 36.17 | 40 | 59.88 |
| Naive Bayes | 57.14 | 38.33 | 97.87 | 55.09 | 70.03 |

*Table 3: Table of Comparison  Analysis of Models*

The performance of each model is graphically compared across the evaluation metrics in the bar chart

- **Accuracy** The best predictive models for liver disease cases are Logistic Regression and Random Forest, with an accuracy of roughly 72%.
- **Precision** When compared to other models, Logistic Regression and Gradient Boosting exhibit higher precision, which indicates fewer false positives.
- **Recall** Naive Bayes has the highest recall of any algorithm, at 97.87%, demonstrating its ability to correctly identify the majority of cases of liver disease despite having a high rate of false positives.
- **F1-Score** Gradient Boosting, Decision Tree, and Naive Bayes have the highest F1 scores when it comes to balancing precision and recall.

- **ROC-AUC** Naive Bayes also performs well overall in class distinction, leading the field in this measure.

*Figure 24: Comparing Machine Learning Models for the Identification of Liver Disease*

The comparative study shows that various models perform well according to various metrics Both Random Forest and Logistic Regression offer high accuracy, which makes them trustworthy for broad forecasts.

Naive Bayes effectively identifies true positive cases at the expense of numerous false positives, as evidenced by its excellent recall and ROC-AUC.

Gradient Boosting is a strong option for managing liver disease detection because it provides a balanced performance across precision, recall, and F1-score.

## 4.8. Actual vs Prediction Explanation

Six different machine learning models—Logistic Regression, Decision Tree, Random Forest, KNearest Neighbors, Gradient Boosting, and Naive Bayes—are used to compare actual and predicted values for liver disease detection in the provided plot. The index of the test samples is shown by the x-axis, and the presence ({1}) or absence ({0}) of liver disease is indicated by the y-axis.

Blue Circles ({o}) These show the true presence or absence of liver disease based on test set values.

```
# Plot Actual vs Predicted values for all classifiers
plt.figure(figsize=(14, 8))

for name, clf in classifiers.items():
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)

    results["Model"].append(name)
    results["Accuracy (%)"].append(accuracy_score(y_test, y_pred) * 100)
    results["Precision (%)"].append(precision_score(y_test, y_pred, zero_division=0) * 100)
    results["Recall (%)"].append(recall_score(y_test, y_pred) * 100)
    results["F1-Score (%)"].append(f1_score(y_test, y_pred) * 100)
    results["ROC-AUC (%)"].append(roc_auc_score(y_test, y_pred) * 100)

    plt.scatter(range(len(y_pred)), y_pred, label=f'{name} Predicted', marker='+')

plt.scatter(range(len(y_test)), y_test, color='blue', label='Actual', marker='o')
plt.title('Actual vs Predicted Values for All Techniques')
plt.xlabel('Index')
plt.ylabel('Liver Disease (0: No, 1: Yes)')
plt.legend(loc='best')
plt.show()

# Create a DataFrame from the results
df_results = pd.DataFrame(results)
```

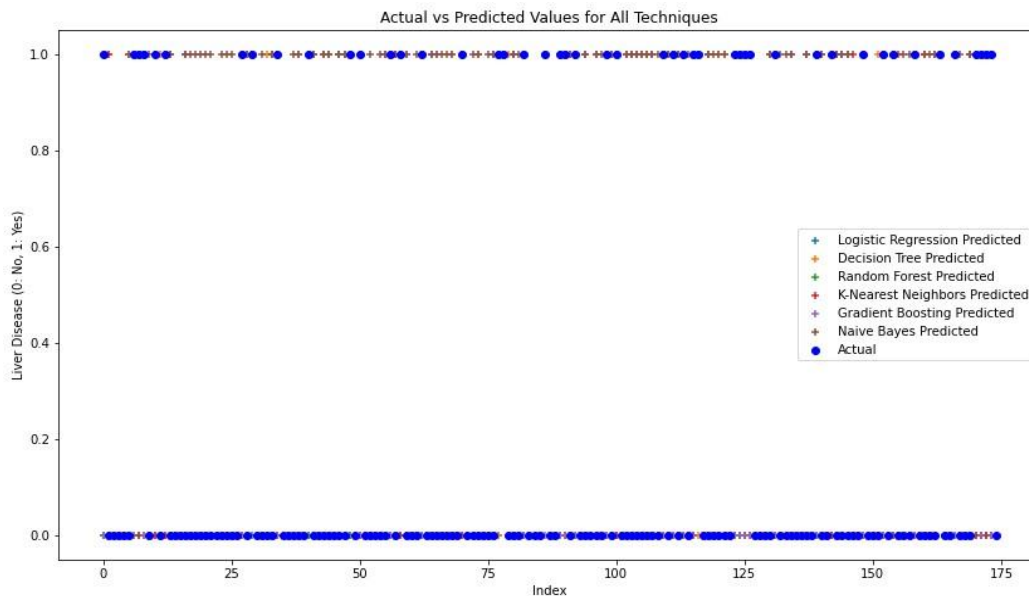*Figure 25: Code for Actual vs Prediction for all techniques*



*Figure 26: Plot representation of Actual vs Prediction for all techniques*

The Colored Crosses ({+}) represent the estimated values derived from every classifier. It is possible to distinguish clearly between the predictions of each model because each one is represented by a different color and label.

55

The degree to which each model's predictions agree with the observed values can be evaluated using this plot. For each test sample index, a model is deemed to be more accurate if its predictions closely match the actual values (blue circles). A model's ability to accurately predict liver disease increases with the distance between the colored crosses and the blue circles. The models that perform well and those that might use some improvement can be quickly identified with the help of this visual representation.

## 4.9. Best Model

It is determined by comparing the models that the Random Forest model has the highest accuracy. The specific performance indicators for the Random Forest mode are as follows

```
# best model with the highest accuracy
best_accuracy_model=df_results.loc[df_results['Accuracy (%)'].idxmax()]
print("Best Model by Accuracy:")
print(best_accuracy_model)
```

*Figure 27: Input Code for Best Model*

By far the most accurate model (72.57%) for predicting liver disease overall, Random Forest is the most suitable model for this purpose. A balance between finding positive cases and preventing false positives is shown by its precision of 48.48% and recall of 34.04%. A moderate ability to distinguish between liver disease and no liver disease cases is suggested by the ROC-AUC score of 60.38%, while the F1-Score of 40.00% further supports its balanced performance.

```
Best Model by Accuracy:
Model              Random Forest
Accuracy (%)           72.571429
Precision (%)          48.484848
Recall (%)             34.042553
F1-Score (%)                40.0
ROC-AUC (%)            60.380652
Name: 2, dtype: object
```

*Figure 28: Output of Best Model*

## 4.9.2. Summary

Notwithstanding the high accuracy, it is crucial to keep in mind that there is still space for improvement in the Random Forest model's recall and precision when it comes to identifying true positive cases and lowering false positives. Its performance might be further improved by fine-tuning it or combining it with other models.

Finally, based on accuracy, the Random Forest model is currently the best-performing model for liver disease detection, offering a solid foundation for additional model optimization and refine
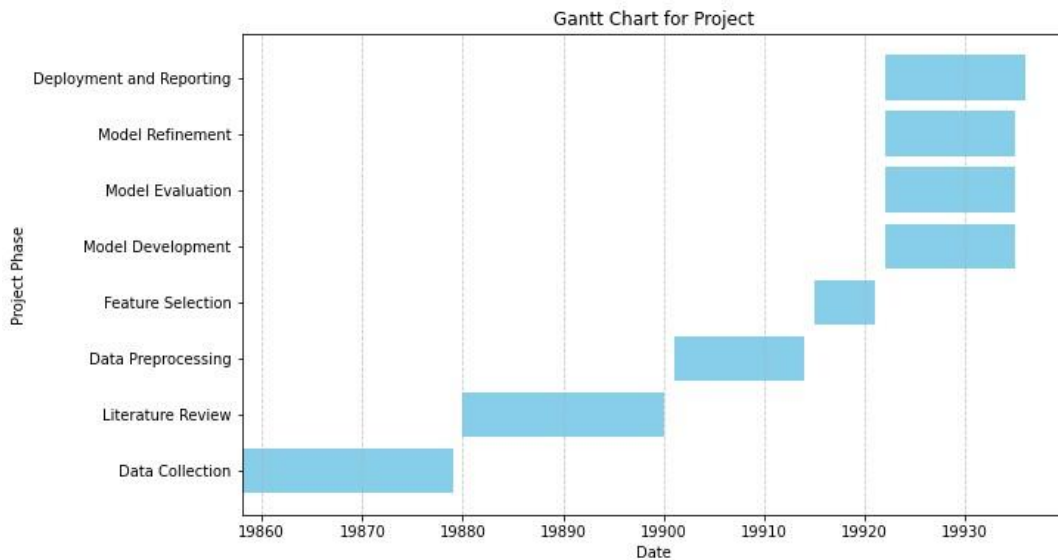
## Chapter 5 Project Management

## 5.1. Project Timetable

Starting on May 15, 2024, the project was successfully finished and turned in on August 1, 2024. The first step took about 14 days to complete and comprised obtaining a large dataset from the UCI Machine Learning Repository. This time was spent examining the dataset, comprehending its organization, and making sure it was appropriate for the project. Subsequently, a comprehensive preprocessing phase was implemented to guarantee the precision and dependability of the data. After that, feature selection was done to determine which liver disease indicators were most important. Then, over several weeks, several machine learning models were created and trained, with ongoing performance optimization and assessment. The best-performing model was put into use in the last phase, which also involved writing the project report and submitting it on time.

*Table 4: Project Schedule*

| Project Phase | Start Date | End Date | Duration (days) |
|---|---|---|---|
| Data Collection | 5/15/2024 | 5/28/2024 | 13 |
| Data Preprocessing | 5/29/2024 | 6/11/2024 | 13 |
| Feature Selection | 6/12/2024 | 6/18/2024 | 6 |
| Model Development | 6/19/2024 | 7/2/2024 | 13 |
| Model Evaluation | 7/3/2024 | 7/9/2024 | 6 |
| Model Refinement | 7/10/2024 | 7/16/2024 | 6 |
| Deployment and Reporting | 7/17/2024 | 8/1/2024 | 15 |

*Figure 29: Gantt Chart for Project*

## 5.2. Risk Management

Considering the possible difficulties with data quality, computational capacity, and project schedules, risk management was an essential component of the project. In-depth data cleaning techniques were used to reduce problems with data quality. Regularization techniques and cross-validation strategies were used to mitigate the risk of overfitting in the model. By keeping a thorough project schedule and making adjustments regularly, potential delays were avoided. Furthermore, the utilization of cloud services guarantees adequate computational capacity to manage the needs for data processing and model training.

## 5.3. Managing Quality

The project deliverables' quality had to be guaranteed at all times. The effectiveness of the machine learning models was evaluated by way of strict testing and validation procedures. Key performance metrics were used to assess the models, including F1-score, recall, accuracy, and precision.

Validation of the findings was ensured by seeking input from domain experts and peer reviews. Holding high standards and producing dependable results required regular quality checks at every stage of the project.

## 5.4. Social, Legal, Ethical, and Professional Considerations

Social, legal, ethical, and professional considerations were given a lot of weight in the project's execution. The project complied with applicable data protection regulations in maintaining the privacy and confidentiality of the patient data used. By making sure that the data was used appropriately and that the models were created and assessed fairly, ethical concerns were taken into account. Throughout the project, professional standards were upheld along with a dedication to accuracy, integrity, and transparency in all facets of the work. The project's potential social impact was also taken into account. Its goal is to improve healthcare by offering a trustworthy tool for the identification of liver disease.

## Chapter 6: Discussion

The project's outcomes show the limitations and effectiveness of different machine learning models in the identification of liver diseases. A comparative study of several models, such as Decision Tree, Random Forest, K-Nearest Neighbors (KNN), Gradient Boosting, Logistic Regression, and Naive Bayes, showed the unique advantages and disadvantages of each strategy. A baseline performance of 72% accuracy was provided by logistic regression; however, its low recall of 19.15% revealed a considerable number of false negatives, indicating that it is less dependable in accurately identifying positive cases of liver disease.

Even though it was simple to understand and visualize, the Decision Tree model performed moderately, with an accuracy of 69.14%. It also had problems with overfitting, as evidenced by the large number of false positives and negatives. With the best accuracy of 72.57%, Random Forest outperformed the other models in terms of balancing recall and precision. It was a dependable option for this dataset because of its resistance to noise and overfitting.

The Decision Tree model performed moderately, with an accuracy of 69.14%, and had overfitting problems, as evidenced by a sizable number of false positives and negatives. Despite this, the model was easy to understand and visualize. With the best accuracy of 72.57%, Random Forest outperformed the other models in terms of balancing recall and precision. It was a dependable option for this dataset because of its resistance to noise and overfitting.

Competitive but mediocre performance was offered by KNN and Gradient Boosting. Although Gradient Boosting needed careful tuning to prevent overfitting but showed promise with a balanced F1-score, KNN had computational limitations, particularly with larger datasets. Its assumption of feature independence did not hold for this dataset, which is why Naive Bayes, for all its simplicity and high recall, had the lowest accuracy.

Overall, all the models have potential for improvement, but the Random Forest model was found to be the most accurate. To further improve model performance, one could investigate more sophisticated ensemble methods, perform extensive hyperparameter tuning, and enhance feature engineering. The results of this research demonstrate the potential of machine learning methods to offer accurate, non-invasive diagnostic tools for liver disease, which will ultimately improve patient outcomes and streamline the delivery of healthcare.

## Chapter 7: Conclusion and Future Work

## 7.1. Conclusion

## 7.1.1. Model Performance Summary

- A variety of metrics, including Accuracy, Precision, Recall, F1-Score, and ROC-AUC, were used to evaluate the performance of each of the six machine learning models we examined to detect liver diseases Logistic Regression, Decision Tree, Random Forest, K-Nearest Neighbors (KNN), Gradient Boosting, and Naive Bayes.
- The Random Forest Classifier showed high precision, recall, F1-score, and ROC-AUC among all the models, making it the most accurate model out of all the others. This suggests that for this particular dataset, the Random Forest model is trustworthy and robust.

## 7.1.2. Models Insights

- **Logistic Regression** established a benchmark for contrasting with other, more intricate models. Although it demonstrated respectable accuracy, tree-based techniques outperformed it.
- **Decision Tree** Despite its propensity to overfit, it demonstrated good interpretability.
- **Random Forest** Because of its capacity to lower overfitting through group learning, it performed the best overall.
- **KNN** Performs moderately well, but when working with big datasets, it can be computationally expensive.
- **Gradient Boosting** This technique performed competitively, nearly matching Random Forest in terms of performance, and it works well with intricate relationships.
- **Naïve Bais** Due to its strong feature independence assumptions, Naive Bayes performed the worst.

### 7.1.3. Confusing Matrix

The true positives, true negatives, false positives, and false negatives were highlighted in the confusion matrices for each model, providing a clear picture of how well each model discriminates between the classes.

### 7.1.4. Study Restrictions

The size and diversity of the dataset might not accurately reflect the larger population, which could have an impact on how broadly applicable the model is.

Just a few features were the subject of the investigation. The performance of the model might be enhanced by adding more pertinent clinical data.

Not all algorithms work well with this kind of dataset, as evidenced by the underperformance of some models, such as Naive Bayes.

KNN has computational limitations that could make it less viable for larger datasets.

### 7.2. Future Work

### 7.2.1. Feature Engineering and Data Augmentation

- Add more pertinent features to the dataset to further improve model accuracy.
- To develop new features that more accurately depict the underlying patterns associated with liver diseases, leverage your domain knowledge.

### 7.2.2. Hyperparameter Tuning

- Perform extensive hyperparameter tuning using techniques such as Grid Search or Random Search to optimize model parameters.
- Implement cross-validation to ensure that the model generalizes well to unseen data.

### 7.2.3. Group Approaches

- To combine the advantages of several models, look into more sophisticated ensemble techniques like stacking or blending.
- To lower variance and increase the predictability of results, apply model averaging.

### 7.2.4. Depth Learning

- For feature extraction and classification, investigate deep learning methods like Convolutional Neural Networks (CNN) and Artificial Neural Networks (ANN), particularly as additional data becomes available.
- Leverage current knowledge for liver disease detection by applying transfer learning from pretrained models.

### 7.2.5. Interpretability and Explainability

- Use model explainability strategies such as SHAP (Shapley Additive exPlanations) to gain a deeper understanding of the contributions made by various features to the anticipated outcomes.
- Make sure the models can be somewhat understood to encourage medical professionals to adopt and trust them.

### 7.2.6. Implementation and Practical Assessment

- Create a deployment pipeline that enables the model to be integrated into clinical workflows to detect liver disease in real time.
- Collaborate with healthcare providers to conduct practical testing to verify the efficacy and dependability of the model.

### 7.2.7. Frequent Monitoring and Updates

- Update the model often with fresh data to keep it accurate and relevant over time.
- Install monitoring systems to keep tabs on the model's performance and identify any forecast drifts or anomalies.

We can enhance the model's precision, resilience, and practicality by tackling these upcoming avenues, which will ultimately help with the early identification and management of liver disorders.

**Student Reflection :**

Using cutting-edge data mining and machine learning techniques in conjunction with medical knowledge, this project on liver disease detection was incredibly illuminating. Utilizing a dataset containing vital medical parameters like bilirubin levels, protein levels, enzyme levels, and demographic data allowed for a thorough grasp of the variables affecting liver health. I now have a deep understanding of the complexity of medical diagnostics and how machine learning can improve efficiency and accuracy in this field thanks to this project.

Being proficient in the preprocessing of medical data was one of the most important learning objectives. An important part of getting the dataset ready for analysis was handling missing values, normalizing the data, and choosing pertinent features. I was also able to comprehend the advantages and disadvantages of different machine learning algorithms in the context of medical data by applying them to logistic regression, decision trees, and support vector machines. My comprehension of their practical implications was further cemented by evaluating model performance using metrics such as accuracy, precision, recall, and F1-score.

The significance of interdisciplinary collaboration was also emphasized by this project. Medical terminology and diagnostic criteria understanding were necessary in addition to technical skills for the integration of medical knowledge with data science techniques. This encounter has motivated me to learn more about the nexus between technology and healthcare since I now see the revolutionary potential of data science to enhance patient outcomes.

To sum up, the project marked a critical turning point in my academic career and gave me invaluable experience in data mining, machine learning, and medical data analysis. It has encouraged me to look into this field more thoroughly and to develop applications, with the ultimate goal of using cuttingedge technology to improve healthcare.

# References

1.  Dhumane, A., Pawar, S., Aswale, R., Sawant, T., & Singh, S. (2023). Effective Detection of Liver Disease Using Machine Learning Algorithms. In *ICT Analysis and Applications* (Vol. 782, pp. 161–171). Springer Singapore Pte. Limited. https://doi.org/10.1007/978-981-99-65687_15

2.  Mohammad Reza Shahraki, & Mahboubeh Mesgar. (2019). Evaluation of Data Mining Algorithms for Detection of Liver Disease. Payavard salamat, 13(1), 81–90

3.  Md, A. Q., Kulkarni, S., Joshua, C. J., Vaichole, T., Mohan, S., & Iwendi, C. (2023). Enhanced Preprocessing Approach Using Ensemble Machine Learning Algorithms for Detecting Liver Disease. *Biomedicines*, *11*(2), 581. https://doi.org/10.3390/biomedicines11020581

4.  Singh, G., Agarwal, C., & Gupta, S. (2022). Detection of Liver Disease Using Machine Learning Techniques: A Systematic Survey. In *Communications in Computer and Information Science* (Vol. 1591, pp. 39–51). Springer International Publishing. https://doi.org/10.1007/9783-031-07012-9_4

5.  Singh, G., Agarwal, C., & Gupta, S. (2022). Detection of Liver Disease Using Machine Learning Techniques: A Systematic Survey. In Communications in Computer and Information Science (Vol. 1591, pp. 39–51). Springer International Publishing. https://doi.org/10.1007/9783-031-07012-9_4

6.  Kalejahi, B. K., Meshgini, S., Danishvar, S., & Khorram, S. (2022). Diagnosis of liver disease by computer- assisted imaging techniques: A literature review. Intelligent Data Analysis, 26(4), 1097–1114. https://doi.org/10.3233/IDA-216379

7.  Altaf, I., Butt, M. A., & Zaman, M. (2021). Disease Detection and Prediction Using the Liver Function Test Data: a Review of Machine Learning Algorithms. *Advances in Intelligent Systems and Computing*, 785–800. https://doi.org/10.1007/978-981-16-2597-8_68

8.  Arbain, A. N., & Balakrishnan, B. Y. P. (2019). A Comparison of Data Mining Algorithms for Liver Disease Prediction on Imbalanced Data. *International Journal of Data Science and Advanced Analytics*, *1*(1), 1–11. https://www.ijdsaa.com/index.php/welcome/article/view/2

9.  Ayeldeen, H., Shaker, O., Ayeldeen, G., & Anwar, K. M. (2015, November 1). *Prediction of liver fibrosis stages by machine learning model: A decision tree approach*. IEEE Xplore. https://doi.org/10.1109/ICoCS.2015.7483212

10. Baitharu, T. R., & Pani, S. K. (2016). Analysis of Data Mining Techniques for Healthcare Decision Support System Using Liver Disorder Dataset. *Procedia Computer Science*, *85*, 862–870. https://doi.org/10.1016/j.procs.2016.05.276

11. *Data-Centric AI Solutions and Emerging Technologies in the Healthcare Ecosystem*. (n.d.).

12. Dritsas, E., & Trigka, M. (2023). Supervised Machine Learning Models for Liver Disease Risk Prediction. *Computers*, *12*(1), 19. https://doi.org/10.3390/computers12010019

13. Dutta, K., Chandra, S., & Gourisaria, M. K. (2022). Early-Stage Detection of Liver Disease Through Machine Learning Algorithms. *Advances in Data and Information Sciences*, 155–166. https://doi.org/10.1007/978-981-16-5689-7_14

14. Engy El-Shafeiy, El-Desouky, A. I., & Elghamrawy, S. M. (2018). Prediction of Liver Diseases Based on Machine Learning Technique for Big Data. *Advances in Intelligent Systems and Computing*, 362–374. https://doi.org/10.1007/978-3-319-74690-6_36

15. Ganie, S. M., Dutta Pramanik, P. K., & Zhao, Z. (2024). Improved liver disease prediction from clinical data through an evaluation of ensemble learning approaches. *BMC Medical Informatics and Decision Making*, *24*(1). https://doi.org/10.1186/s12911-024-02550-y

16. Ghosh, M., Mohsin Sarker Raihan, Md., Raihan, M., Akter, L., Kumar Bairagi, A., S. Alshamrani, S., & Masud, M. (2021). A Comparative Analysis of Machine Learning Algorithms to Predict Liver Disease. *Intelligent Automation & Soft Computing*, *30*(3), 917–928.https://doi.org/10.32604/iasc.2021.017989

17. Gogi, V. J., & Vijayalakshmi M.N. (2018). *Prognosis of Liver Disease: Using Machine Learning Algorithms*. https://doi.org/10.1109/icrieece44171.2018.9008482

18. Harshpreet Kaur, G. S. (2021). The Diagnosis of Chronic Liver Disease using Machine Learning Techniques. *INFORMATION TECHNOLOGY in INDUSTRY*, *9*(2), 554–564. https://doi.org/10.17762/itii.v9i2.382

19. Jain, D., & Singh, V. (2018). Feature selection and classification systems for chronic disease prediction: A review. *Egyptian Informatics Journal*, *19*(3), 179–189. https://doi.org/10.1016/j.eij.2018.03.002

20. Jesty, B., Jesty, B., & Niranjan, M. (2019, May 1). *Machine learning for liver disease classification*. Eprints.soton.ac.uk. https://eprints.soton.ac.uk/433924/

21. Kumar, Y., & Sahoo, G. (2013). Prediction of different types of liver diseases using rule based classification model. *Technology and Health Care*, *21*(5), 417–432. https://doi.org/10.3233/thc130742

22. Muthuselvan, S., Rajapraksh, S., Somasundaram, K., & Karthik, K. (2018). Classification of Liver Patient Dataset Using Machine Learning Algorithms. *International Journal of Engineering & Technology*, *7*(3.34), 323. https://doi.org/10.14419/ijet.v7i3.34.19217

23. Polat, K., & Sentürk, U. (2018, October 1). *A Novel ML Approach to Prediction of Breast Cancer: Combining of mad normalization, KMC based feature weighting and AdaBoostM1 classifier*. IEEE Xplore. https://doi.org/10.1109/ISMSIT.2018.8567245

24. Shaheamlung, G., Kaur, H., & Kaur, M. (2020). A Survey on machine learning techniques for the diagnosis of liver disease. *2020 International Conference on Intelligent Engineering and Management (ICIEM)*. https://doi.org/10.1109/iciem48762.2020.9160097

25. Shahid Mohammad Ganie, Dutta, K., & Zhao, Z. (2024). Improved liver disease prediction from clinical data through an evaluation of ensemble learning approaches. *BMC Medical Informatics and Decision Making*, *24*(1). https://doi.org/10.1186/s12911-024-02550-y

26. Singh, A., & Pandey, B. (2014). *Intelligent techniques and applications in liver disorders: a survey*. *16*(1), 27–27. https://doi.org/10.1504/ijbet.2014.065638

27. Singh, G., Agarwal, C., & Gupta, S. (2022). Detection of Liver Disease Using Machine Learning Techniques: A Systematic Survey. *Communications in Computer and Information Science*, 39–51. https://doi.org/10.1007/978-3-031-07012-9_4

28. Siuly, S., & Zhang, Y. (2016). Medical Big Data: Neurological Diseases Diagnosis Through Medical Data Analysis. *Data Science and Engineering*, *1*(2), 54–64. https://doi.org/10.1007/s41019-016-0011-3

29. Smith, A., Baumgartner, K., & Bositis, C. (2019). Cirrhosis: Diagnosis and Management. *American Family Physician*, *100*(12), 759–770. https://www.aafp.org/pubs/afp/issues/2019/1215/p759.html/amp/

30. Takkar, S., Singh, A., & Pandey, B. (2017). Application of Machine Learning Algorithms to a Well Defined Clinical Problem: Liver Disease. *International Journal of E-Health and Medical Communications*, *8*(4), 38–60. https://doi.org/10.4018/ijehmc.2017100103

31. Velu, S. R., Ravi, V., & Tabianan, K. (2022). Data mining in predicting liver patients using classification model. *Health and Technology*, *12*(6), 1211–1235. https://doi.org/10.1007/s12553-022-00713-3

32. Venkata Ramana, B., Babu, M. Surendra. P., & Venkateswarlu, N. B. (2011). A Critical Study of Selected Classification Algorithms for Liver Disease Diagnosis. *International Journal of Database Management Systems*, *3*(2), 101–114. https://doi.org/10.5121/ijdms.2011.3207

33. Zhang, Y. (2010). New Advances in Machine Learning. In *Google Books*. BoD – Books on Demand. https://books.google.co.uk/books?hl=en&lr=&id=XAqhDwAAQBAJ&oi=fnd&pg=PA19&dq=Ayodele

34. Zhao, X., Wang, R., & Cao, Z. (n.d.). A Data Mining Approach to Prediction of Liver Diseases You may also like Application Progress of Bear Bile Powder and Ursodeoxycholic Acid in Liver Disease and its Mechanism of Action. *Journal of Physics*. https://doi.org/10.1088/17426596/1529/3/032002

35. S., Dr. Aneeshkumar, & Jothi, V. C. (2012, November 10). *Estimating the Surveillance of Liver Disorder using Classification Algorithms*. Ssrn.com. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3509372

**Appendix A : Code**

```python
# 1 Logistic Model
logistic_mdl=LogisticRegression(random_state=42)
logistic_mdl.fit(X_train,y_train)

# Make predictions and evaluate
y_pred_log=logistic_mdl.predict(X_test)
print("Logistic Regression Results:")
print(f'Accuracy:{accuracy_score(y_test,y_pred_log)*100:.2f}%')
print(f'Precision:{precision_score(y_test,y_pred_log,zero_division=0)*100:.2f}%')
print(f'Recall: {recall_score(y_test, y_pred_log)*100:.2f}%')
print(f'F1-Score: {f1_score(y_test, y_pred_log)*100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_log)*100:.2f}%\n')

import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

# Calculate the confusion matrix
confu_matrix=confusion_matrix(y_test, logistic_mdl.predict(X_test), labels=[0, 1])

# Display the confusion matrix
disp=ConfusionMatrixDisplay(confusion_matrix=confu_matrix,display_labels=['No Liver Disease', 'Liver Disease'])
disp.plot(cmap=plt.cm.Blues)
plt.title('Confusion Matrix for Logistic Regression')
plt.show()
```

```python
#2.Decision Tree Classifier

from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

decision_tree_model=DecisionTreeClassifier(random_state=42)
decision_tree_model.fit(X_train, y_train)

# Make predictions and evaluate
y_pred_tree=decision_tree_model.predict(X_test)
print("Decision Tree Results:")
print(f' Accuracy: {accuracy_score(y_test, y_pred_tree) * 100:.2f}%')
print(f'  Precision: {precision_score(y_test, y_pred_tree, zero_division=0) * 100:.2f}%')
print(f'  Recall: {recall_score(y_test, y_pred_tree) * 100:.2f}%')
print(f'  F1-Score: {f1_score(y_test, y_pred_tree) * 100:.2f}%')
print(f'  ROC-AUC: {roc_auc_score(y_test, y_pred_tree) * 100:.2f}%\n')

import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

# Calculate the confusion matrix
confu_matrix=confusion_matrix(y_test,decision_tree_model.predict(X_test),labels=[0,1])

# Display the confusion matrix
display=ConfusionMatrixDisplay(confusion_matrix=confu_matrix,display_labels=['No Liver Disease', 'Liver Disease'])
disp.plot(cmap=plt.cm.Blues)
plt.title('Confusion Matrix for Decision Tree')
plt.show()
```

```python
#3 Random Forest Classifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
forest_model = RandomForestClassifier(random_state=42)
forest_model.fit(X_train, y_train)

# Make predictions and evaluate
y_pred_forest=forest_model.predict(X_test)
print("Random Forest Results:")
print(f'Accuracy: {accuracy_score(y_test, y_pred_forest)*100:.2f}%')
print(f'Precision: {precision_score(y_test, y_pred_forest, zero_division=0)*00:.2f}%')
print(f'Recall:{recall_score(y_test, y_pred_forest)*100:.2f}%')
print(f'F1-Score:{f1_score(y_test, y_pred_forest)*100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_forest) * 100:.2f}%\n')
confu_matrix_forest=confusion_matrix(y_test,forest_model.predict(X_test),labels=[0, 1])
disp_forest = ConfusionMatrixDisplay(confusion_matrix=confu_matrix_forest, display_labels=['No Liver Disease', 'Liver Disease'])
disp_forest.plot(cmap=plt.cm.Greens)
plt.title('Confusion Matrix for Random Forest')
plt.show()

from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
```

```python
#4 KNN
knn_model=KNeighborsClassifier(n_neighbors=5)
knn_model.fit(X_train,y_train)

# Make predictions and evaluate
y_pred_knn=knn_model.predict(X_test)
print("K-Nearest Neighbors Results:")
print(f'Accuracy: {accuracy_score(y_test, y_pred_knn)*100:.2f}%')
print(f'Precision: {precision_score(y_test, y_pred_knn, zero_division=0)*100:.2f}%')
print(f'Recall: {recall_score(y_test, y_pred_knn)*100:.2f}%')
print(f'F1-Score: {f1_score(y_test, y_pred_knn)*100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_knn)*100:.2f}%\n')
confu_matrix_knn=confusion_matrix(y_test,knn_model.predict(X_test),labels=[0, 1])
disp_knn=ConfusionMatrixDisplay(confusion_matrix=confu_matrix_knn,display_labels=['No Liver Disease', 'Liver Disease'])
disp_knn.plot(cmap=plt.cm.Oranges)
plt.title('Confusion Matrix for K-Nearest Neighbors')
plt.show()
```

```python
#5 Gradient Boosting Classifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

# Initialize and train the model
gb_model=GradientBoostingClassifier(random_state=42)
gb_model.fit(X_train,y_train)

# Make predictions and evaluate
y_pred_gb = gb_model.predict(X_test)
print("Gradient Boosting Results:")
print(f'Accuracy: {accuracy_score(y_test, y_pred_gb) * 100:.2f}%')
print(f' Precision: {precision_score(y_test, y_pred_gb, zero_division=0) * 100:.2f}%')
print(f'Recall: {recall_score(y_test, y_pred_gb) * 100:.2f}%')
print(f'F1-Score: {f1_score(y_test, y_pred_gb) * 100:.2f}%')
print(f'ROC-AUC: {roc_auc_score(y_test, y_pred_gb) * 100:.2f}%\n')
confu_matrix_gb= confusion_matrix(y_test,gb_model.predict(X_test),labels=[0, 1])
disp_gb=ConfusionMatrixDisplay(confusion_matrix=confu_matrix_gb, display_labels=['No Liver Disease', 'Liver Disease'])
disp_gb.plot(cmap=plt.cm.Purples)
plt.title('Confusion Matrix for Gradient Boosting')
plt.show()
```