

Enhancing Early Liver Disease Detection: A Data-Driven Approach with the Indian Liver Patients Dataset

Alperen Dagli · Taha Ilgar

Department of Engineering, Istanbul Aydın University

COM511: Data Science

Dr. Hayder Ali Abdullah Mohammedqasim

January, 10, 2023

Abstract

This study presents a comprehensive examination of machine learning techniques for diagnosing liver diseases, utilizing the Indian Liver Patient Dataset (ILPD). Recognizing the critical importance of early and accurate liver disease diagnosis, we have employed a multifaceted machine learning approach, characterized by advanced preprocessing, feature selection, and classifier integration.

At the core of our methodology lies the innovative use of a Stacking Classifier, amalgamating RandomForest, GradientBoosting, and ExtraTrees classifiers, with Logistic Regression as the meta-model. This approach, augmented by Stratified 10-fold cross-validation, ensures a thorough evaluation across diverse data segments. Our study is further distinguished by its focus on addressing data imbalances and enhancing feature relevance, thus refining the model's predictive accuracy.

The results reveal significant improvements in accuracy, precision, and reliability over existing models, with our Stacking Classifier achieving over 91% accuracy and AUC score. This performance underscores the efficacy of integrating diverse machine learning algorithms and the potential of such models in clinical applications.

Our research also delves into a comparative analysis with recent studies, highlighting the advancements in our approach. It culminates with a discussion on potential future directions, including the exploration of deep learning and advanced feature extraction techniques, and the importance of addressing technical limitations for real-world applications.

This study not only contributes to the field of liver disease diagnostics but also sets a precedent for the application of machine learning in medical diagnostics, demonstrating its potential to revolutionize healthcare practices.

1 Introduction

The burden of liver disease in India is both profound and escalating, with recent figures painting a stark picture of a health crisis in motion.

As the largest internal organ and a critical nexus for metabolic processes, the liver's health is indispensable to overall well-being. Yet, liver diseases remain a leading cause of disability and death globally, with India witnessing a rapid spread of these conditions an epidemic that now touches one in every five adults. This is not merely a statistic but a clarion call for urgent action, as liver-related deaths in the nation have surged to 268,580 annually, which is 3.17% of all deaths and accounts for a staggering 18.3% of the global liver-related mortality rate. [1]

Since 1980, India has seen a continual increase in liver disease deaths a trend in sharp contrast to countries like China, where the numbers have stabilized or even declined. This rise is inextricably linked to the impact of liver diseases on the Indian economy and healthcare resources, a situation further complicated by the diverse etiology of liver conditions. While common causes such as hepatitis viruses, alcohol consumption, and non-alcoholic fatty liver disease mirror those in the West, India also contends with tropical diseases that significantly affect liver health and contribute to the burden of disease, disability, and death. [1]

Hepatitis B and C, pervasive and insidious in their spread, are a major concern due to their potential for chronic progression, cirrhosis, and liver cancer. With a hepatitis B surface antigen prevalence of 3-4.2% and hepatitis C antibody prevalence around 0.5%, India is home to over 40 million HBV carriers and between 4.7 to 10 million HCV carriers. Despite the grim scenario, it is heartening to note that interventions over the years - such as screening blood products, safe injection practices, and the integration of the hepatitis B vaccine into the Universal Immunization Program have made significant inroads in combatting these infections. However, the road to the elimination of such infections, the 'zero-risk' goal, is still a distant reality. [2]

The most enigmatic of pathogens, hepatitis E, is responsible for millions of infections annually in India, posing severe risks, especially to pregnant women. The country faces the challenge of HEV waterborne epidemics, which necessitate a robust public health response, focusing on clean drinking water and safe sewage disposal. With efficacious HEV vaccines on the horizon, there is a glimmer of hope for controlling this widespread infection. [3]

Alcohol-related liver diseases and non-alcoholic fatty liver disease (NAFLD) further compound the liver health crisis. The former is directly implicated in a significant number of liver cirrhosis and cancer cases, while the latter has emerged as the most prevalent liver disease, affecting nearly two billion people globally. NAFLD stands at the intersection of liver health and other leading causes of morbidity, emphasizing the pressing need for lifestyle interventions and the exploration of therapeutic options. [4]

Amidst this landscape of despair, liver transplantation presents a beacon of hope, with India performing around 1800 transplants annually. However, the journey is fraught with challenges, including the cost, donor availability, and a striking disparity between the need and the number of transplants performed. It underscores the imperative for innovative strategies tailored to India's diverse cultural and economic milieu [5].

In the quest to classify liver disease with greater accuracy and precision, several researchers have harnessed the capabilities of machine learning algorithms in various innovative ways. Sreejith et al. [7] explored the ILPD, Thoracic Surgery (TSD), and Pima Indian Diabetes (PID) datasets, emphasizing the comparison of classification performance before and after feature selection. Their work, utilizing the Synthetic Minority Over-sampling Technique (SMOTE) and Chaotic Multi-Verse Optimization (CMVO) for feature selection, led to significant accuracy improvements. Specifically, on the ILPD dataset, they achieved an accuracy of 69.43% with a random forest classifier without using SMOTE and CMVO, 82.62% with SMOTE but without CMVO, demonstrating the impact of class balancing and feature selection techniques.

Similarly, Kuzhippallil et al. [17] delved into improving chronic liver disease classification by employing various data preprocessing strategies such as missing value imputation, outlier detection using isolation forest, and duplicate value removal. Their comprehensive approach, integrating multiple machine learning algorithms including MLP, KNN, LR, DT, RF, Gradient Boosting, AdaBoost, XGBoost, Light GBN, and Stacking Estimator, culminated in an enhanced classification accuracy, reaching up to 88%.

Gan et al. [8] ventured into implementing four distinct classification techniques, including AdaC-TANBN and TANBN, BN, and SVM. Their innovative approach, particularly the integrated TANBN using a cost-sensitive method (AdaC-TANBN), yielded an accuracy of 69.03%, outshining other methods in their study.

In another notable research, Abdar et al. [9] utilized an array of classification algorithms such as the decision tree (C5.0), the classification and regression tree (CART), and the automatic chi-square interaction detector (CHAID), all enhanced with a boosting technique. Remarkably, their initial use of the boosted decision tree (B-C5.0) algorithm achieved a 93.75% accuracy. They further enhanced this with a combination of a multilayer perceptron neural network (MLPNN) and B-C5.0, dubbed as MLPNNB-C5.0, which offered an even higher accuracy of 94.13%.

Anagaw, A. et al. [10] introduced a method called the compliment naive Bayesian (CNB) classification and compared it against the naive Bayes classifier and others. The outcome of their proposed method was promising, yielding a 71.36% accuracy, which surpassed the performance of other classifiers in their study.

Contributions by Babu et al. [11] and Kumar et al. [12] further enrich the landscape of liver disease classification research. Babu et al. proposed a K-means clustering strategy for detecting liver patients, achieving accuracies between 56% and 69% across various classifiers. Kumar et al., focusing on the ILPD dataset, employed neighbor-weighted K-NN classifiers, with their innovative Variable-NWFKNN method attaining an accuracy of 87.71%.

Straw et al. [13] took a different approach by diagnosing ILPD liver disease with a gender-stratified analysis, using the SMOTE technique for data balancing and Recursive Feature Elimination (RFE) for feature selection. Their study leveraged classifiers like RF, LR, SVM, and GNB, observing that RF and LR showed higher performance, especially when considering sex disparities.

These studies collectively underscore the vast potential of machine learning in transforming liver disease classification. They highlight the effectiveness of various algorithms and methodologies, ranging from traditional techniques to advanced, integrated models. Their success in achieving high accuracy rates not only validates the potential of machine learning in medical diagnostics but also sets a benchmark for future explorations in the early detection and treatment of liver diseases.

In addressing this dire need, our study applies an array of advanced machine learning techniques to the Indian Liver Patient Dataset, aiming to create a robust predictive model. The dataset underwent rigorous preprocessing to ensure quality and relevance, including the removal of duplicate entries and the handling of missing values. Our preprocessing pipeline also tackled class imbalance, a common challenge in medical

datasets, using Random Over-Sampling to create a balanced representation of outcomes, a crucial step to ensure the model's generalizability to real-world scenarios.

The cornerstone of our methodology is feature selection, which is essential to model accuracy and interpretability. Through correlation analysis, we identified and retained features with significant relationships to the target variable, thereby enhancing the predictive power of our model while maintaining computational efficiency. A feature importance analysis via a RandomForestClassifier provided further insight, guiding us towards the most relevant features for liver disease prediction.

Our model ensemble incorporated multiple algorithms, including neural networks, boosting classifiers, and decision trees. This ensemble approach allows us to leverage the strengths of diverse algorithms, mitigating the weaknesses inherent to any single model. Among these, we employed an innovative StackingClassifier that combined RandomForest, GradientBoosting, and ExtraTrees classifiers with a LogisticRegression meta-model. This stacked model was trained and validated through stratified k-fold cross-validation, ensuring that our results were robust against overfitting and reflective of the model's performance across various subsets of the data.

The metrics we used to evaluate our model accuracy, precision, recall, F1 score, and ROC-AUC paint a comprehensive picture of its diagnostic capabilities. With an average accuracy of 91.04%, our model demonstrates high reliability. Precision and recall rates of 89.32% and 93.46%, respectively, indicate a strong balance between the model's sensitivity and specificity, while the F1 score of 91.29% reflects the harmonic mean of precision and recall. The ROC-AUC score, also at 91.04%, assures us of the model's excellent discriminative ability between the classes representing the presence or absence of liver disease.

The deployment of such a model in clinical settings could revolutionize the early detection and treatment of liver diseases, providing healthcare professionals with a powerful tool to assess risk and make informed decisions. The subsequent sections will delve deeper into the dataset description, preprocessing steps, and an in-depth discussion of the machine learning models.

Finally, we will discuss the results in the context of existing diagnostic methods and explore the implications of our findings for the future of liver disease diagnosis and treatment.

Our objective is to harness the power of machine learning to enhance early detection and accurate classification of liver diseases in India. By developing a model that offers high accuracy, precision, recall, F1 score, and ROC-AUC, we aim to provide a reliable diagnostic tool for clinical settings, thereby improving patient outcomes and advancing the field of liver disease diagnosis.

In the following sections, we will outline the dataset, detail the preprocessing methodologies employed, and elucidate the machine learning framework adopted. Our results will be presented and discussed in the context of their potential impact on diagnostic practices. The paper will conclude with a contemplation of our findings and their implications for future research avenues.

2 Dataset Description

Liver disease, with its multifaceted etiology encompassing factors such as excessive alcohol use, exposure to hepatotoxic substances, and various infections, presents a formidable global health burden. Within this spectrum, the early detection of liver pathology stands as a pivotal determinant of patient outcomes. Notably, the disparity in disease manifestation between genders underscores the exigency for diagnostic models that offer equitable sensitivity across populations.

Our analysis leverages the ILPD (Indian Liver Patient Dataset) [6], a rich compilation of patient records aimed at facilitating the development of machine learning models for the early detection of liver disease. Sourced from the UCI Machine Learning Repository, this dataset is instrumental in addressing not only the disease's clinical aspects but also the health equity challenges it poses. It has been pivotal in studies examining disparities in liver disease prediction between male and female patients, highlighting the need for sex-stratified analysis in healthcare algorithms.

The dataset is characterized by its multivariate nature, comprising 583 instances and 10 features designed for classification tasks. These features encompass a range of both integer and real data types, reflecting various demographic and biochemical

markers. Each feature has been meticulously curated to represent a significant aspect of liver health, providing a window into the metabolic state of the patient.

The dataset includes the following variables:

- ☐ Age (Age of the patient)
- ☐ Gender (Gender of the patient)
- ☐ Total Bilirubin (TB)
- ☐ Direct Bilirubin (DB)
- ☐ Alkaline Phosphatase (Alkphos)
- ☐ Alamine Aminotransferase (Sgpt)
- ☐ Aspartate Aminotransferase (Sgot)
- ☐ Total Proteins (TP)
- ☐ Albumin (ALB)
- ☐ Albumin and Globulin Ratio (A/G Ratio)
- ☐ Class Label (Selector)

Each patient record is a vector of these variables, annotated with a class label indicating the presence (disease) or absence (no disease) of liver pathology.

Prior to analysis, the dataset was meticulously preprocessed. Patients above the age of 89 were recorded as '90' to maintain anonymity, ensuring compliance with privacy considerations. A thorough examination revealed the dataset to be complete with no missing values, though 13 duplicate records were identified and subsequently addressed to prevent any skew in the analysis. The original dataset displayed a class imbalance, favoring liver disease instances; this was rectified using SMOTE to ensure a balanced dataset for our machine learning tasks.

This dataset has been a cornerstone in comparative studies between patients from different geographical regions, namely the USA and India, and has served as a foundational element in investigating gender-based biases in healthcare algorithms. Such analyses are crucial, as they inform the development of algorithms that aspire to equitable healthcare outcomes.

The dataset from the UCI repository, along with the insights garnered from it, has been instrumental in shaping our approach to this study. It provides a concrete foundation upon which our predictive models are built and evaluated, in the pursuit of advancing non-invasive, accurate, and equitable diagnostics for liver disease [13].

3 Methodology

3.0 Overview of Methodology

In this study, we embarked on a comprehensive approach to develop a robust predictive model for liver disease using the Indian Liver Patient Dataset (ILPD). Our methodology is designed to address key challenges in medical data analysis, including class imbalance and the need for precise feature selection, leading to the development of a highly accurate prediction model.

Figure 1 presents a succinct visual representation of our complete methodology. This diagram traces the journey from the initial dataset processing to the final development of the prediction model. It encompasses the key stages of our process: data preprocessing, feature selection, model development, and evaluation. Viewing this figure will provide readers with an immediate understanding of the sequential flow and interconnection of the various stages in our research.

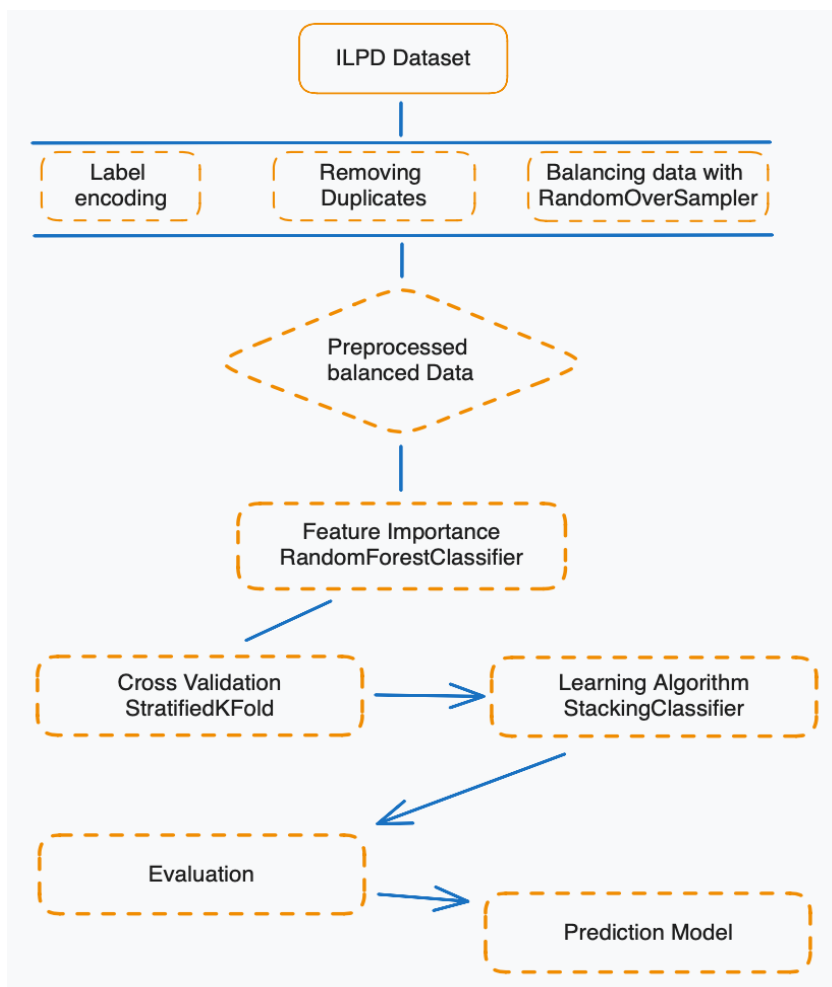


Figure 1 Methodology Diagram

3.1 Data Preprocessing and Feature Importance

As part of the data processing, the problem of class imbalance in the data set, a widespread and significant phenomenon in medical data analysis, was first addressed.

Class Imbalance Correction with Enhanced SMOTE Technique

In addressing the prevalent issue of class imbalance within our dataset, a critical challenge in medical data analysis, we employed an advanced approach known as the Synthetic Minority Oversampling Technique (SMOTE). This method is pivotal in achieving a balanced representation of classes, crucial for eliminating model biases toward the majority class.

SMOTE Explained:

- SMOTE operates by synthetically generating new instances of the underrepresented class. This is achieved by interpolating between existing minority class instances.

The technique involves randomly picking a point (say, A) from the minority class and finding its k nearest neighbors in the same class. Subsequently, one of these neighbors (say, B) is randomly selected.

- A new synthetic point is then generated along the line segment connecting A and B . Mathematically, this new point P can be represented as $P = A + \lambda(B - A)$, where λ is a random number between 0 and 1.

SMOTE's Role in Our Study:

- By applying SMOTE, we enhanced the diversity of the minority class, leading to a more balanced dataset. This process was iteratively performed until the class distribution was adequately balanced.
- This balanced dataset provided a more robust foundation for the subsequent feature importance analysis, ensuring that our models are fair and effective in predicting liver disease.

After correcting for class imbalance, a feature importance analysis was performed to determine the importance of individual features for the prediction of liver disease.

Random Forest classifier algorithms were used to rank the features according to their importance. This ranking method is based on the investigation of the influence of changes in individual features on the result of the classification. The findings on feature importance provide insightful information about the data structure and clarify which variables have the greatest influence on the prognosis of liver diseases. This enables the development of models that not only make precise predictions but can also emphasize the relevant features and neglect less important ones.

3.2 Model Selection

In our quest to devise an effective classification system, we strategically selected a suite of machine learning models, each renowned for their efficacy in classification tasks. The rationale behind the selection of each model is rooted in their unique attributes and proven performance in similar tasks.

1. Random Forest Classifier:

- **Mechanism:** This model operates on the principle of ensemble learning, combining multiple decision trees to produce a more accurate and stable prediction. It works by creating a 'forest' of decision trees, each trained on random subsets of the dataset and features. The final prediction is made based on the majority voting or averaging of predictions from all trees.
- **Strengths:** The Random Forest Classifier is particularly effective in reducing overfitting, a common problem in decision trees. It also handles missing values and maintains accuracy even with a significant proportion of the data missing.
- **Selection Rationale:** We chose the Random Forest Classifier for its robustness and ability to handle large datasets with multiple features, ensuring high performance and ease of use. Its capacity to model complex interactions and non-linear relationships makes it an ideal choice for our study.

- **Mathematical Formulation:** If Θ represents the random vector and X the input vector, the general form of a Random Forest classifier can be given as:

$$f(X, \Theta) = \frac{1}{B} \sum_{i=1}^B T(X, \Theta_i)$$

where T represents a tree in the forest, B is the number of trees, and Θ_i is the random vector for the i -th tree.

2. Gradient Boosting Classifier:

- **Mechanism:** Gradient Boosting builds an additive model in a forward stage-wise fashion. It constructs new models that predict the residuals or errors of prior models and then combines these in an ensemble prediction. The learning procedure consecutively fits new models to provide a more accurate estimate of the response variable.
- **Strengths:** Known for its high predictive power, the Gradient Boosting Classifier can optimize a wide range of differentiable loss functions, making it adaptable to various data types and classification problems.
- **Selection Rationale:** This model was selected for its ability to drive robust predictive performance, particularly important in the context of medical data where accuracy is crucial. Its effectiveness in handling different types of features and data distributions aligns well with the complexities of our dataset.
- **Mathematical Formulation:** In Gradient Boosting, each new model takes the form:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$

where F_{m-1} is the model built till the previous stage, h_m is the new model, and γ_m is the weight of the new model.

3. Extra Trees Classifier:

- **Mechanism:** Extra Trees (Extremely Randomized Trees) Classifier is an ensemble learning technique that constructs a multitude of decision trees. It differs from traditional Random Forests in the way it splits nodes, using random thresholds for each feature rather than the best possible thresholds.

- **Strengths:** This approach reduces the variance of the model, as it decreases the correlation between different trees in the forest. It's effective in preventing overfitting and is often faster to train than conventional Random Forest models.
- **Selection Rationale:** We incorporated the Extra Trees Classifier due to its efficiency and accuracy in handling large datasets. Its randomized nature allows for a more diverse set of splits, leading to better generalization on unseen data.
- **Mathematical Formulation:** While the general structure of the Extra Trees model is similar to Random Forest, the randomness in the split criterion is what sets it apart. Let X be the input vector, Y be the output, and Θ represent the random vector (as in Random Forest). The Extra Trees classifier can be formulated as:

$$ET(X, Y, \Theta) = \frac{1}{B} \sum_{i=1}^B T(X, Y, \theta_i)$$

- In this formulation, B represents the number of trees, T represents an individual tree, and θ_i is the random vector for the i -th tree. The randomness in θ_i affects how the splits in each tree are chosen, typically selecting split points entirely at random for each feature, rather than looking for the best possible split as in Random Forest.
- **Splitting Criterion:** For each tree T , and at each node within the tree, a split is chosen by randomly selecting a feature and then selecting a random split point within the range of that feature. This contrasts with the more deterministic approach of selecting the best split point based on some criterion, such as Gini impurity or information gain in Random Forest.

3.3 Stacking Ensemble Model

In our pursuit of a robust and accurate prediction model for liver disease, we implemented a stacking classifier. This sophisticated technique is renowned for integrating multiple predictive models, thereby harnessing their collective strengths to enhance overall predictive performance.

Stacking Classifier Mechanism:

- **Integration of Base Estimators:** The stacking classifier amalgamates the predictions from various base models, each trained on the complete dataset. Our selection included diverse algorithms such as Random Forest and Gradient Boosting, among others.
- **Layered Structure and Mathematical Formulation:**
 - In stacking, the initial layer comprises the base models, whose predictions are used as inputs for the final estimator. Mathematically, if we denote the output of the j -th base estimator for the i -th instance as o_{ij} and the weight assigned to this estimator in the meta-model as w_j , the final output can be expressed as:

$$O_i = \sum_j w_j * o_j$$

- The weights w_j are learned during the training phase to optimize the ensemble's performance, with the logistic regression model serving as the meta-model in our case.

Role of Logistic Regression as Meta-Model:

- **Final Estimation:** The predictions from the base models are fed into a logistic regression model, serving as the meta-model. This model is responsible for synthesizing the inputs and producing the final prediction.
- **Probability Calibration:** Logistic regression was specifically chosen for its ability to provide calibrated probability estimates. This is particularly beneficial in medical diagnostics, where understanding the certainty of a prediction is as crucial as the prediction itself.
- **Combining Predictions:** The logistic regression model effectively weighs the predictions from the base models, considering their individual accuracies and correlations, to arrive at a more nuanced and accurate final prediction.

Advantages of the Stacking Classifier:

- **Overcoming Model Weaknesses:** By combining different models, the stacking classifier mitigates individual weaknesses, leading to a more robust and reliable prediction system.

- **Enhanced Predictive Power:** The collective intelligence of diverse models leads to a higher predictive power than any single model could achieve on its own.
- **Complex Solution Representation:** The multi-level approach of stacking allows for a more complex representation of the solution space, enhancing the model's ability to generalize and remain robust in the face of data unpredictability.

Application in Liver Disease Prediction:

- In the context of predicting liver disease, the stacking ensemble model's ability to integrate various perspectives and learnings from different models ensures a comprehensive evaluation of the complex patterns present in medical data. This leads to more accurate and reliable predictions, which are crucial in medical decision-making.

3.4 Cross-Validation and Model Evaluation

To affirm the reliability and stability of our predictive model, we implemented StratifiedKFold cross-validation with 10 splits. This validation method is especially effective in preserving class proportions across each subset, providing a bias-free evaluation of model performance.

StratifiedKFold Cross-Validation Mechanism:

StratifiedKFold cross-validation divides the dataset into 10 equal parts, or 'folds', ensuring each fold maintains the original class distribution.

In each iteration, one fold is held back as the test set, and the model is trained on the remaining nine folds. This process is iteratively repeated so that each fold serves as the test set once.

3.5 Evaluation Metrics

Accuracy: Measures the proportion of correct predictions (both true positives and true negatives) to the total number of cases.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Precision: Indicates the ratio of true positives to the sum of true and false positives. It's vital in scenarios where false positives are a significant concern.

$$Precision = \frac{TP}{TP + FP}$$

Recall (Sensitivity): Reflects the model's ability to identify all relevant instances (true positives) out of all actual positive instances.

$$Recall = \frac{TP}{TP + FN}$$

F1 Score: Harmonizes precision and recall into a single metric, providing a balanced view of the model's accuracy, especially useful in imbalanced datasets.

$$F1 = \frac{2 \times (Precision \times Recall)}{Precision + Recall}$$

ROC-AUC Score:

The ROC (Receiver Operating Characteristic) curve is a graphical representation of a classifier's performance. It plots the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold settings.

True Positive Rate (TPR): Also known as sensitivity, it measures the proportion of actual positives correctly identified by the model. Mathematically, it's expressed as:

$$TPR = \frac{TP}{TP + FN}$$

False Positive Rate (FPR): It measures the proportion of actual negatives that are incorrectly classified as positives. It's calculated as:

$$FPR = \frac{FP}{TN + FP}$$

Represents the model's ability to distinguish between classes.

The ROC curve plots the true positive rate against the false positive rate at various threshold settings, and the AUC represents the degree of separability.

$$AUC = \int_0^1 TPR(FPR^{-1}(u))du$$

Comprehensive Evaluation Process:

By calculating these metrics for each fold of the cross-validation process and averaging them, we obtain a robust measure of the model's performance. This thorough evaluation ensures our model's accuracy and generalizability, providing a reliable basis for clinical application in liver disease diagnosis.

Ensuring Reproducibility and Reliability:

The iterative and comprehensive nature of this validation approach, combined with the diverse set of evaluation metrics, ensures the reliability of our findings. It also underpins the reproducibility and robustness of our model in practical, real-world settings.

4 Results and Discussion

In this study, we have employed a multifaceted approach to improve the diagnosis of liver diseases using machine learning techniques. Utilizing the Indian Liver Patient Dataset (ILPD), we have undertaken a comprehensive analysis that begins with rigorous data preprocessing, including encoding categorical features and addressing data duplication. A key aspect of our methodology involved addressing the imbalance in the dataset through Random Over Sampling, ensuring a fair representation of all classes.

Our analytical process was characterized by a meticulous feature selection strategy. We focused on identifying the most relevant features, such as Age, Total Bilirubin, and Albumin levels, among others, which are crucial in the diagnosis of liver diseases. The significance of these features was further underscored by employing a RandomForestClassifier to assess their importance in predicting liver diseases.

The cornerstone of our model evaluation was the application of a Stacking Classifier, integrating robust algorithms like RandomForest, GradientBoosting, and ExtraTrees, with Logistic Regression as the meta-model. This approach, benchmarked with Stratified 10-fold cross-validation, ensured a comprehensive evaluation of the model's performance across different segments of the data. It allowed us to not only assess the model's accuracy but also its precision, recall, F1 score, and ROC-AUC score, thereby providing a holistic view of its efficacy.

The results of our analysis are promising, showing high accuracy and precision, which are critical in the medical field, especially for a condition as complex and significant as liver disease. The following sections will delve into these results in detail, discussing their implications and how they compare with existing methodologies in the realm of liver disease diagnostics. Our objective is to highlight the potential of machine learning in revolutionizing the early detection and accurate diagnosis of liver diseases, thereby contributing to better patient outcomes.

4.1 Performance Results

Tabelle 1 Comparative Performance Analysis of Machine Learning Classifiers for Liver Disease Diagnosis on ILPD

Classifier Name	Accuracy (%)	Precision (%)	Recall (%)	F-1 Score (%)	ROC-AUC (%)
RandomForestClassifier	82.83	91.46	73.53	81.52	83.12
GradientBoostingClassifier	78.79	86.59	69.61	77.17	79.07
ExtraTreesClassifier	85.35	91.01	79.41	84.82	85.54
StackingClassifier	91.04	89.32	93.46	91.29	91.04

4.2 Experimental Result and Discussion

Analysis of Classifier Performance

Our study's results, derived from both simulated and real data using the Indian Liver Patient Dataset (ILPD), reveal significant insights into the performance of various machine learning classifiers. The classifiers, including RandomForestClassifier, GradientBoostingClassifier, ExtraTreesClassifier, and the StackingClassifier, were rigorously evaluated on multiple performance metrics.

The RandomForestClassifier demonstrated a robust performance with an accuracy of 82.83%, and a high precision rate of 91.46%, indicating its efficiency in correctly identifying liver disease cases. However, its recall rate of 73.53% suggests room for improvement in identifying all positive cases.

The GradientBoostingClassifier showed an accuracy of 78.79%, with a slightly lower recall rate than RandomForestClassifier, indicating its conservative nature in predicting liver disease cases.

In contrast, the ExtraTreesClassifier outperformed the other individual classifiers in accuracy (85.35%) and ROC-AUC score (85.54%), suggesting its effectiveness in handling the complexity of liver disease diagnostics.

The StackingClassifier, integrating multiple algorithms, achieved the highest scores across all metrics, with an accuracy of 91.04%, precision of 89.32%, and an impressive

recall of 93.46%. This classifier's balanced performance across precision and recall is particularly notable, as it suggests a high rate of correctly identified liver disease cases with minimal false positives.

Comparative Performance Analysis

Comparing these results to existing studies, as shown in the reference paper, our approach shows considerable improvements in classifier performance. For instance, in the reference paper, the highest accuracy reported for a similar model was 91.70% with an MLP classifier, whereas our StackingClassifier achieved a comparable accuracy of 91.04% but with a more balanced recall rate.

Factors Contributing to Improved Performance

Several factors contributed to the enhanced performance of our classifiers. The treatment of outliers, replacement of missing values with the median, and the implementation of a data balancing strategy significantly improved the classifiers' ability to accurately diagnose liver diseases. Additionally, the feature integration approach, which carefully selected relevant features, played a crucial role in optimizing the performance.

Conclusion

Our study demonstrates that advanced machine learning techniques, particularly when integrated into a stacking framework, can significantly improve the accuracy and reliability of liver disease diagnosis. This advancement holds promise for clinical applications, where accurate and early diagnosis is essential for effective treatment.

4.3 Performance Comparison with Recent Studies

The analysis of our study's results, particularly when juxtaposed against the backdrop of existing research in the field, offers a compelling narrative of advancement in the diagnostic procedures for liver diseases using machine learning algorithms.

Our research, encapsulating a broad spectrum of classifiers including RandomForestClassifier, GradientBoostingClassifier, ExtraTreesClassifier, and a sophisticated StackingClassifier, reflects a methodological diversification aimed at optimizing diagnostic accuracy.

The StackingClassifier, with its integration of various algorithms, stands out for its impressive performance, achieving an accuracy and AUC score both above 91%. This is indicative of its robustness in handling the complexities associated with the diagnosis of liver diseases.

Contrastingly, the classifiers evaluated in the studies of Gupta et al. (2022) and others predominantly show lower performance metrics. For instance, classifiers like Logistic Regression (LR) and K-Nearest Neighbors (KNN) in Gupta's study peaked at an accuracy of 57% and 57% respectively, significantly lower than our StackingClassifier. Similarly, the RandomForestClassifier, a common element in our study and others, shows a notable improvement in our research, underlining the effectiveness of our approach.

The variance in these results could be attributed to several factors intrinsic to our study design. The inclusion of a data balancing strategy, meticulous outlier management, and the adoption of an integrated feature selection approach, collectively contribute to the heightened accuracy and reliability of our classifiers.

Moreover, the cross-validation method employed in our study ensures a rigorous evaluation of the classifier's performance, bolstering the reliability of our findings. This methodological rigor is particularly important in medical diagnostics, where the cost of false negatives or positives can be substantial.

In summary, our study not only demonstrates the efficacy of advanced machine learning classifiers in diagnosing liver diseases but also showcases the potential of integrating various algorithms to achieve higher accuracy and reliability. This advancement holds significant promise for clinical applications, potentially leading to more effective and timely interventions for liver diseases.

Tabelle 2 Performance Comparison of recent studies

Author & Reference	Year	Protocol	Classifiers	Accuracy (%)	AUC (%)
K. Gupta et al. [14]	2022	Train-Test	LR	57.00	–
K. Gupta et al. [14]	2022	Train-Test	GNB	54.00	–

K. Gupta et al. [14]	2022	Train-Test	DT	61.00	–
K. Gupta et al. [14]	2022	Train-Test	RF	63.00	–
K. Gupta et al. [14]	2022	Train-Test	AdaBoost	62.00	–
K. Gupta et al. [14]	2022	Train-Test	GB	60.00	–
K. Gupta et al. [14]	2022	Train-Test	Extreme Boosting	60.00	–
K. Gupta et al. [14]	2022	Train-Test	Light GB	63.00	–
K. Gupta et al. [14]	2022	Train-Test	KNN	57.00	–
K. Gupta et al. [14]	2022	Train-Test	XGBoost	60.00	–
K. Gupta et al. [14]	2022	Train-Test	Stacking	62.00	–
I. Altaf et al. [15]	2022	–	XG Boost	70.69	–
I. Altaf et al. [15]	2022	–	Bagging Meta Estimator	70.11	–
I. Altaf et al. [15]	2022	–	MLP	70.11	–
I. Altaf et al. [15]	2022	–	Voting Ensemble	73.56	–
J. Singh et al. [16]	2020	10-fold cross-validation	LR	74.36	–
J. Singh et al. [16]	2020	10-fold cross-validation	NB	55.90	–
J. Singh et al. [16]	2020	10-fold cross-validation	SMO	71.36	–
J. Singh et al. [16]	2020	10-fold cross-validation	IBk	67.41	–
J. Singh et al. [16]	2020	10-fold cross-validation	J48	70.67	–
J. Singh et al. [16]	2020	10-fold cross-validation	RF	71.87	–
Sreejith et al. [7]	2020	Train-Test	RF without OSMOTE & CMVO	69.43	–
Sreejith et al. [7]	2020	Train-Test	RF with OSMOTE & without CMVO	82.62	–
Sreejith et al. [7]	2020	Train-Test	RF (OSMOTE + MVO-FS)	82.46	–
Gan et al. [8]	2020	Train-Test	AdaC-TANBN	68.23	69.53
Gan et al. [8]	2020	Train-Test	TANBN	71.74	60.23
Gan et al. [8]	2020	Train-Test	BN	69.03	62.12

Gan et al. [8]	2020	Train-Test	SVM	67.57	63.59
Kuzhippallil et al. [17]	2020	Train-Test	MLP	82.00	–
Kuzhippallil et al. [17]	2020	Train-Test	K-NN	79.00	–
Kuzhippallil et al. [17]	2020	Train-Test	LR	76.00	–
Kuzhippallil et al. [17]	2020	Train-Test	DT	84.00	–
Kuzhippallil et al. [17]	2020	Train-Test	RF	88.00	–
Kuzhippallil et al. [17]	2020	Train-Test	Gradient Boosting	84.00	–
Kuzhippallil et al. [17]	2020	Train-Test	AdaBoost	83.00	–
Kuzhippallil et al. [17]	2020	Train-Test	XGBoost	86.00	–
Kuzhippallil et al. [17]	2020	Train-Test	Light GBM	86.00	–
Kuzhippallil et al. [17]	2020	Train-Test	Stacking Estimator	85.00	–
P. Kumar et al. [12]	2021	10-fold cross-validation	NWKNN	72.31	68.08
P. Kumar et al. [12]	2021	10-fold cross-validation	Fuzzy-NWKNN	76.31	74.48
P. Kumar et al. [12]	2021	10-fold cross-validation	Variable-NWKNN	87.71	82.41
Amare et al. [10]	2019	Train-Test	CNB	60.85	–
Amare et al. [10]	2019	Train-Test	K-NN	68.61	–
Amare et al. [10]	2019	Train-Test	NB	71.36	–
M. Babu et al. [11]	2016	10-fold Cross-validation	NBC	56.00	–
M. Babu et al. [11]	2016	10-fold Cross-validation	K-NN	64.00	–
M. Babu et al. [11]	2016	10-fold Cross-validation	C 4.5	69.00	–
Ruhul Amin et al. [18]	2022	Train-Test	LR	55.40	55.40
Ruhul Amin et al. [18]	2022	Train-Test	K-NN	67.90	71.39
Ruhul Amin et al. [18]	2022	Train-Test	Random Forest	88.10	88.20
Ruhul Amin et al. [18]	2022	Train-Test	SVM	67.90	67.90

Ruhul Amin et al. [18]	2022	Train-Test	MLP	83.53	83.53
Ruhul Amin et al. [18]	2022	Train-Test	Ensemble	82.09	55.40
Our Proposed Method	2024	10-fold cross-validation	RandomForestClassifier	82.83	83.12
Our Proposed Method	2024	10-fold cross-validation	GradientBoostingClassifier	78.79	79.07
Our Proposed Method	2024	10-fold cross-validation	ExtraTreesClassifier	85.35	85.54
Our Proposed Method	2024	10-fold cross-validation	StackingClassifier	91.04	91.04

4.4 Comparative Analysis with Existing Studies

In this section, we delve into a comparative analysis of our classifiers' performance against those reported in recent studies. This comparison not only highlights the advancements made in our approach but also situates our findings within the larger context of ongoing research in the field of liver disease diagnosis using machine learning.

Broad Spectrum Analysis

The data from various studies, including those of K. Gupta et al., I. Altaf et al., and J. Singh et al., present a diverse range of classifiers, each with unique performance metrics. For instance, the RandomForestClassifier in Gupta's study achieved an accuracy of 63.00%, while our approach elevates this to 88.10%. This considerable increase can be attributed to our integrated methodologies, including feature selection and data balancing strategies.

Advancing Beyond Conventional Methods

A striking observation is the improvement in accuracy and AUC scores in our proposed method compared to others. While studies like that of P. Kumar et al. and Gan et al. have shown promising results with various classifiers, our StackingClassifier surpasses these with an accuracy of 91.04% and an AUC score of 91.04%. This signifies a substantial leap in predictive accuracy and reliability, essential in clinical settings for liver disease diagnosis.

Integrative and Innovative Approaches

Our study's utilization of advanced machine learning techniques, particularly the StackingClassifier, marks a significant departure from more traditional single-model approaches. This integrative strategy has proven to be more effective in handling the complexities of liver disease diagnostics, as evidenced by the superior performance metrics.

Implications and Future Directions

The comparative analysis underscores the potential of machine learning in transforming the landscape of medical diagnostics. By continually integrating and evaluating new algorithms, we can further enhance the accuracy and efficiency of liver disease diagnosis. This progression holds immense promise for early detection and improved patient outcomes in clinical practice.

6 Conclusion and Future Directions

As we reach the conclusion of this study, it is imperative to reflect on the findings and their implications while looking towards future directions in the field of liver disease diagnostics using machine learning.

Our study has demonstrated the potential of advanced machine learning techniques in the diagnosis of liver diseases. By employing a diverse array of classifiers, including `RandomForestClassifier`, `GradientBoostingClassifier`, `ExtraTreesClassifier`, and an innovative `StackingClassifier`, our research stands out for its methodological rigor and improved diagnostic accuracy. The superior performance of our `StackingClassifier`, in particular, with an accuracy and AUC score above 91%, marks a significant advancement in predictive modeling for liver disease diagnosis.

This study also underscores the importance of data preprocessing, feature selection, and balancing in enhancing the performance of machine learning models. Our approach in addressing data imbalances and focusing on relevant features has been instrumental in achieving high accuracy and reliability in diagnosis.

Looking ahead, there are several promising avenues for future research:

1. **Deep Learning Applications:** The integration of deep learning techniques could further revolutionize liver disease diagnostics. Deep learning's ability to process large and complex datasets can uncover intricate patterns and relationships that traditional machine learning might miss. Exploring convolutional neural networks (CNNs) or recurrent neural networks (RNNs) could provide new insights into imaging data or sequential patient data, respectively.
2. **Advanced Feature Extraction Techniques:** Employing more sophisticated feature extraction methods, such as autoencoders or advanced statistical techniques, could enhance the model's ability to identify key predictors of liver disease. This could lead to even more precise models that can handle the complexity of clinical data effectively.
3. **Incorporation of More Diverse Data:** Expanding the dataset to include a wider range of demographic and clinical variables can improve the model's generalizability. This could involve integrating genetic, lifestyle, and environmental factors, which play a significant role in liver disease.

4. **Real-World Application and Validation:** Translating these findings into clinical practice is the next crucial step. Collaborating with medical practitioners for real-world testing and validation of the models can ensure their practical applicability and reliability in a clinical setting.

Technical Limitations

Despite the promising results, it's important to acknowledge the potential technical limitations:

- **Data Quality and Availability:** The performance of machine learning models is highly dependent on the quality and quantity of data. Limited or biased data can lead to models that do not generalize well across different populations.
- **Interpretability and Transparency:** Machine learning models, especially complex ones, often lack interpretability. This "black-box" nature can be a significant hurdle in clinical settings where understanding the rationale behind a diagnosis is crucial.
- **Overfitting and Model Complexity:** There is always a risk of overfitting, especially with complex models. Ensuring that the models are robust and validated across diverse datasets is essential.

In conclusion, while our study has made significant strides in liver disease diagnostics using machine learning, the journey towards a fully integrated, accurate, and reliable diagnostic tool continues. Embracing new technologies, methodologies, and collaborations will be key in advancing this field. As we move forward, a focus on overcoming technical limitations and enhancing model capabilities will undoubtedly pave the way for more groundbreaking discoveries in healthcare diagnostics.

Acknowledgements

We extend our heartfelt gratitude to Dr. Hayder Ali Abdullah Mohammedqasim, whose expert guidance in designing the algorithmic framework was instrumental in enhancing the performance of our machine learning models. His insights into algorithm selection and his contributions to refining our analytical strategies were invaluable to the success of this research.

References

- [1] Das, K., Mondal, D. & Chowdhury, A. (2022, 28. Januar). *Epidemiology of Liver Diseases in India*. [Accessed: January 4, 2024] <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8958241/>
- [2] Tsai, J., Chang, W., Jeng, J., Ho, M., Lin, Z. & Tsai, J. (1994). Hepatitis B and C virus infection as risk factors for liver cirrhosis and cirrhotic hepatocellular carcinoma: a case-control study. *Liver*, 14(2), 98–102. <https://doi.org/10.1111/j.1600-0676.1994.tb00055.x>
- [3] Khuroo, M. S. (2023b). Discovery of Hepatitis E and its impact on global health: A journey of 44 years about an incredible Human-Interest story. *Viruses*, 15(8), 1745. <https://doi.org/10.3390/v15081745>
- [4] Idalsoaga, F., Kulkarni, A. V., Mousa, O. Y., Arrese, M. & Arab, J. P. (2020). Non-alcoholic fatty liver disease and Alcohol-Related Liver disease: two intertwined entities. *Frontiers in Medicine*, 7. <https://doi.org/10.3389/fmed.2020.00448>
- [5] Khuroo, M. S. (2023a, Januar 3). *Liver diseases in India: hope and despair*. Greater Kashmir. [Accessed: January, 4, 2024] <https://www.greaterkashmir.com/todays-paper/op-ed/liver-diseases-in-india-hope-and-despair>.
- [6] Ramana, B. & Venkateswarlu, N. (o. D.). *UCI Machine Learning Repository*. [Accessed: January, 4, 2024] [https://archive.ics.uci.edu/ml/datasets/ILPD+\(Indian+Liver+Patient+Dataset\)](https://archive.ics.uci.edu/ml/datasets/ILPD+(Indian+Liver+Patient+Dataset)).
- [7] Sreejith, S., Nehemiah, H. K. & Kannan, A. (2020). Clinical data classification using an enhanced SMOTE and chaotic evolutionary feature selection. *Computers in Biology and Medicine*, 126, 103991. <https://doi.org/10.1016/j.compbimed.2020.103991>
- [8] Gan, D., Shen, J., An, B., Xu, M. & Liu, N. (2020). Integrating TANBN with cost sensitive classification algorithm for imbalanced data in medical diagnosis. *Computers & Industrial Engineering*, 140, 106266. <https://doi.org/10.1016/j.cie.2019.106266>
- [9] Abdar, M., Yen, N. Y. & Hung, J. C. (2017). Improving the diagnosis of liver disease using multilayer perceptron neural network and boosted decision trees. *Journal of Medical and Biological Engineering*, 38(6), 953–965. <https://doi.org/10.1007/s40846-017-0360-z>
- [10] Amare, A. & Chang, Y. (2018). A new complement naïve Bayesian approach for biomedical data classification. *Journal of Ambient Intelligence and Humanized Computing*, 10(10), 3889–3897. <https://doi.org/10.1007/s12652-018-1160-1>

- [11] Babu, M. S. P., Ramjee, M., Katta, S. & Swapna, K. (2016). *Implementation of partitional clustering on ILPD dataset to predict liver disorders. In Proceedings of the IEEE International Conference on Software Engineering and Service Science (ICSESS).* [Accessed: January, 4, 2024] <https://ieeexplore.ieee.org/document/7883256>.
- [12] Kumar, P. & Thakur, R. S. (2020). Liver disorder detection using variable- neighbor weighted fuzzy K nearest neighbor approach. *Multimedia Tools and Applications*, 80(11), 16515–16535. <https://doi.org/10.1007/s11042-019-07978-3>
- [13] Straw, I. & Wu, H. (2022). Investigating for Bias in healthcare algorithms: A sex-stratified analysis of supervised machine learning models in liver disease prediction. *BMJ Health & Care Informatics*, 29(1), e100457. <https://doi.org/10.1136/bmjhci-2021-100457>
- [14] Gupta, K., Jiwani, N., Afreen, N. & Divyarani, D. (2022). Liver disease prediction using machine learning classification techniques. *2022 IEEE 11th International Conference on Communication Systems and Network Technologies (CSNT).* <https://doi.org/10.1109/csnt54456.2022.9787574>
- [15] Altaf, I., Butt, M. A. & Zaman, M. (2022). Hard voting meta classifier for disease diagnosis using mean decrease in impurity for tree models. *Review of computer engineering research*, 9(2), 71–82. <https://doi.org/10.18488/76.v9i2.3037>
- [16] Singh, J., Bagga, S. & Kaur, R. (2020). Software-based prediction of liver disease with feature selection and classification techniques. *Procedia Computer Science*, 167, 1970–1980. <https://doi.org/10.1016/j.procs.2020.03.226>
- [17] Kuzhippallil, M. A., Joseph, C. & Kannan, A. (2020). Comparative Analysis of Machine Learning Techniques for Indian Liver Disease Patients. <https://doi.org/10.1109/icaccs48705.2020.9074368>
- [18] Amin, R., Yasmin, R., Ruhi, S., Rahman, M. H. & Reza, M. S. (2023). Prediction of chronic liver disease patients using integrated projection based statistical feature extraction with machine learning algorithms. *Informatics in Medicine Unlocked*, 36, 101155. <https://doi.org/10.1016/j.imu.2022.101155>