# Liver Cirrhosis detection using Random forest

**Deepak Kumar**
M.Tech Scholar
School of Computer Technology, Sanjeev Agrawal Global Educational (SAGE) University, Bhopal, India,
deepak25042000kumar@gmail.com

**Dr. Prashant Kumar Shrivastava**
Associate professor
School of Computer Technology, Sanjeev Agrawal Global Educational (SAGE) University, Bhopal, India,
Prashants@sageuniversity.edu.in

*Abstract*—**Tuberculosis (TB) remains a significant global health challenge, particularly in low- and middle-income countries, where timely diagnosis and treatment are critical to controlling the disease. Traditional diagnostic methods, such as sputum microscopy and chest radiography, are often time-consuming, require trained personnel, and may lack sensitivity, especially in resource-limited settings. Recent advancements in machine learning (ML) have opened new avenues for enhancing TB diagnosis by leveraging clinical, demographic, and radiological data for early prediction. Among various ML algorithms, Random Forest (RF)—a robust ensemble learning technique based on decision trees—has demonstrated high accuracy, efficiency, and scalability in medical diagnosis applications. This review paper provides a comprehensive analysis of the use of Random Forest in the prediction and diagnosis of tuberculosis. It explores various studies that employ RF models, discusses the types of datasets used, preprocessing techniques, feature selection methods, and performance evaluation metrics. The paper also compares Random Forest with other machine learning models in terms of diagnostic accuracy and interpretability. Furthermore, it highlights the current limitations in existing research, such as data imbalance and generalizability, and outlines future directions to enhance the application of RF in clinical settings. Through this review, we aim to demonstrate the potential of Random Forest as a reliable decision support tool for TB prediction and contribute to the ongoing efforts in adopting intelligent healthcare solutions.**

*Keywords*—**Tuberculosis (TB), Random Forest, Machine Learning, Medical Diagnosis, Predictive Modeling, Classification**

## I INTRODUCTION

Tuberculosis (TB) is one of the deadliest infectious diseases, responsible for over 1.5 million deaths annually, according to the World Health Organization. Despite being preventable and curable, TB remains a persistent challenge, particularly in developing regions where diagnostic infrastructure is often limited. The conventional methods of TB diagnosis—such as sputum smear microscopy, culture testing, and chest radiography—are either time-consuming, less sensitive in early-stage detection, or require trained personnel and laboratory facilities. These limitations call for the development of more efficient, automated, and cost-effective diagnostic tools.

In recent years, machine learning (ML) has emerged as a transformative tool in healthcare, enabling predictive modeling and pattern recognition in large and complex datasets. Within the scope of TB detection, ML algorithms have shown promise in identifying critical patterns in clinical, demographic, and radiographic data that may not be evident through traditional diagnostic procedures.

Among these algorithms, the Random Forest (RF) classifier has gained considerable traction due to its high accuracy, stability, and interpretability. RF operates by constructing multiple decision trees during training and aggregating their outputs to make a final prediction. This ensemble learning approach not only reduces the risk of overfitting but also handles imbalanced data and missing values effectively—common challenges in medical datasets.

What distinguishes RF in the context of TB prediction is its ability to process heterogeneous data types, such as symptoms, blood markers, and image features, while providing insights into feature importance
[1] [2] have reported that RF outperforms traditional statistical methods and even other ML models in TB classification tasks. For instance, Li et al. integrated RF with artificial neural networks for pulmonary TB diagnosis and achieved significantly improved sensitivity and specificity compared to standalone models.

Additionally, Patel et al. [3] conducted a meta-analysis of ML-based TB diagnostic systems and noted that RF consistently yielded high performance metrics across multiple datasets, with reported accuracies exceeding 95% in some cases. This consistency makes RF a strong candidate for deployment in clinical decision support systems, especially in low-resource environments where rapid and reliable diagnosis is critical.

## II RELATED WORK

The application of machine learning (ML) algorithms, particularly Random Forest (RF), in the prediction and diagnosis of tuberculosis (TB) has garnered significant attention in recent years. This section reviews key studies that have explored RF's efficacy in TB-related tasks, highlighting methodologies, datasets, and performance metrics.

### A. RF-Based Models for TB Diagnosis

Differentiation Between TB and Sarcoidosis[4]

RF model to distinguish between lung sarcoidosis and TB using clinical and laboratory data from 252 patients. The model achieved an area under the receiver operating characteristic (ROC) curve of 0.915, with a classification error rate of 24.9%. Key contributing factors included angiotensin-converting enzyme and prothrombin time, demonstrating RF's capability in handling complex diagnostic tasks.

Active vs. Latent TB Infection [5]

A conditional random forest (c forest) model to differentiate between active TB (ATB) and latent TB infection (LTBI). Utilizing laboratory data, the c forest model achieved an ROC-AUC of 0.995 in the training set and 0.978 in the test set, with sensitivities of 98.85% and 93.39%, respectively. The model highlighted the importance of biomarkers such as ESAT-6 and CFP-10 in distinguishing between ATB and LTBI .

Integration of RF and Artificial Neural Networks [6]

A hybrid model combining RF and artificial neural networks (ANN) for diagnosing pulmonary TB. The RF algorithm was used to select significant biomarkers, which were then evaluated by the ANN. This approach leveraged the strengths of both models, enhancing diagnostic accuracy and robustness .

### B. Systematic Reviews and Comparative Analyses

Systematic Review of ML Approaches in TB Diagnosis [7] Analyzed 19 studies focusing on biomarker-based TB diagnosis using ML. The review found that RF, along with Support Vector

Machine (SVM), were among the top-performing algorithms, achieving accuracies up to 97%, sensitivities up to 99.2%, and specificities up to 98%. The study emphasized the potential of ML in improving TB diagnosis, especially in low-resource settings .

[8] Predicting Treatment Success in TB Patients the use of various ML models, including RF, to predict sputum culture conversion in TB patients during treatment. The study found that RF, along with Decision Tree (DT) and SVM, achieved high performance with area under the curve (AUC) values greater than 80%. The findings suggest that ML models can be instrumental in monitoring treatment progress and predicting outcomes .

### C. Challenges and Future Directions

While RF has demonstrated promising results in TB prediction, several challenges remain:

Data Quality and Availability: The performance of RF models heavily depends on the quality and quantity of data. Incomplete or biased datasets can lead to overfitting and reduced generalizability.

Interpretability: Although RF provides feature importance scores, the overall model can be perceived as a "black box," making it challenging to interpret and trust in clinical settings. Integration with Clinical Workflows: Implementing RF models in real-world clinical environments requires seamless integration with existing systems and workflows, which can be complex and resource-intensive.

Future research should focus on:

Enhancing Data Collection: Collaborating with healthcare institutions to collect diverse and comprehensive datasets Improving Model Interpretability: Developing techniques to make RF models more transparent and interpretable for clinicians.

Real-World Validation: Conducting prospective studies to validate the effectiveness of RF models in actual clinical settings. techniques like one-hot encoding or label encoding.

**Normalization and Scaling**: Feature values are scaled using Min- Max or Z-score normalization to ensure uniformity across inputs.

Balancing the Dataset: Since TB-positive cases are often underrepresented, techniques like SMOTE (Synthetic Minority Oversampling Technique) or random undersampling are used to balance the dataset.

### C. Feature Selection

Identifying relevant features to improving model accuracy and reducing overfitting. Feature selection in TB prediction involves: **Correlation Analysis**: Removing highly correlated or irrelevant features.

**Random Forest Feature Importance:** Using RF's built-in feature importance scores to retain only top-ranking predictors.

**Domain Knowledge**: Consulting clinical experts or literature to retain biologically or clinically relevant variables.

Example: Patel et al. [7] found that features like chest pain, ESR level, and lymphocyte count were among the top contributors to TB prediction in RF models.

### D. Random Forest Model Training

The RF algorithm will form a large number of decision trees during training. Each tree is trained on a bootstrap sample (randomly selected subset) of the dataset, and splits are made based on the best feature selected using a criterion like Gini Index or Information Gain.

Key hyperparameters include:

**Number of Trees (n_estimators):** More trees generally improve accuracy but increase computational cost.

**Maximum Tree Depth (max_depth):** Controls the complexity of the model.

**Minimum Samples Split (min_samples_split):** Determines the minimum number of samples required to split node.

The RF model outputs a class label based on a majority vote among individual trees.

### III       METHODOLOGY

This section outlines the typical methodological pipeline used in various studies for the prediction of Tuberculosis (TB) using the Random Forest (RF) algorithm. The methodology generally consists of multiple key stages: data collection, preprocessing, feature selection, model training, evaluation, and validation. Each component plays a crucial role in the performance and reliability of the prediction model.

### A. Data Collection

The first and most critical step is acquiring relevant and high- quality datasets. TB prediction studies have used a variety of data sources:

Clinical Data: Includes patient symptoms (e.g., cough, weight loss, fever), demographic information (e.g., age, gender, location), and medical history.

Laboratory Data: Includes results of blood tests, sputum microscopy, chest X-ray scores, and biomarker levels.

Public Datasets: Studies often use open-source repositories like the NIH Chest X-ray Dataset, or datasets from WHO or national health departments.

Example: Yao et al. [4] collected clinical and lab data from 252 patients, while Li et al. [2] used biomarker profiles for prediction.

**B.   Data Preprocessing**

Preprocessing is essential for transforming raw data into a suitable format for training. Common preprocessing steps include:

**Handling Missing Values**: Missing data is imputed using statistical methods (mean, median) or using model-based imputers.

**Encoding Categorical Variables**: Categorical data (e.g., gender, region) is converted into numerical form using

**E.   Model Evaluation**

After training, the model is tested on a separate validation or test set to assess its predictive performance. Common evaluation metrics include:

**Accuracy:** Overall correctness of predictions.

**Precision:** Ratio of true positive predictions to all predicted positives.

**Recall:** Ratio of true positives to all actual positives. Specificity: Ability to identify true negatives.

**F1-Score:** Harmonic mean of precision and recall.

**AUC-ROC:** providing a measure of model discrimination.

Example: Zhang et al. [5] reported AUC values of 0.978 on test data using RF to distinguish active TB from latent TB.

**F.   Cross-Validation and Hyper parameter Tuning**

To avoid over fitting and validate model generalizability, K-fold cross-validation is applied. The dataset is divided into K equal parts; the model is trained on K-1 parts and validated on the remaining part. This process is repeated K times.

Hyper parameter tuning is performed using:

**Grid Search:** Exhaustively testing all possible combinations of hyperparameters.

**Random Search:** Randomly selecting combinations to find optimal settings faster.

**G.   Model Interpretation**

Although RF is considered a "black-box" model, its interpretability can be enhanced through:

**Feature Importance Ranking:** Helps understand which variables most influence predictions.

**Partial Dependence Plots (PDPs):** Illustrate the effect of a feature on the predicted outcome.

**SHAP (SHapley Additive exPlanations):** Advanced technique for detailed model explanation.

Data Collection → Preprocessing → Feature Selection → RF Training → Evaluation → Validation → Interpretation

IV                RESULTS AND DISCUSSION

A.   Performance Metrics of Random Forest in TB Prediction Recent studies have demonstrated the efficacy of Random Forest (RF) in predicting tuberculosis (TB) outcomes, showcasing high accuracy and robustness across various datasets.

Diagnostic Accuracy: In a study by Li et al., the RF model achieved an accuracy of 94.68% in predicting TB diagnosis, with a precision of 96.78% and a recall of 97.04% on the training dataset. The testing dataset yielded an accuracy of 92.18%, precision of 95.30%, and recall of 95.60%, indicating the model's generalizability .

Distinguishing Active from Latent TB: Zhang et al. developed a conditional RF model that distinguished between active TB (ATB) and latent TB infection (LTBI) using laboratory data. The model achieved an area under the receiver operating characteristic curve (AUC) of 0.995 in the training set and 0.978 in the test set, with sensitivities of 98.85% and 93.39%, respectively .

Sarcoidosis vs. TB Classification: Yao et al. applied RF to classify sarcoidosis and TB using clinical and laboratory data. The model achieved an AUC of 0.915, demonstrating RF's capability in handling complex diagnostic tasks .

B.   Comparative Analysis with Other Machine Learning Models RF has been compared with other machine learning algorithms in TB prediction, highlighting its strengths and limitations. Comparison with Decision Trees and Support Vector Machines: In a study by Ahamed Fayaz et al., RF outperformed Support Vector Machine (SVM) and Naïve Bayes (NB) models, achieving higher accuracy and F1 scores. However, the Decision Tree (DT) model achieved the highest accuracy of 92.72% and an AUC of 0.909, surpassing RF's AUC of 0.896 . Integration with Artificial Neural Networks: Li et al. proposed a hybrid model combining RF and artificial neural networks (ANN) for diagnosing pulmonary TB. The RF algorithm was used to select significant biomarkers, which were then evaluated by the ANN. This approach leveraged the strengths of both models, enhancing diagnostic accuracy and robustness .

C.   Feature Importance and Interpretability

RF provides insights into feature importance, aiding in the understanding of key predictors for TB diagnosis.

Identification of Critical Features: In the study by Yao et al., features such as age of onset, drug regimen, lung cavity size, and number of daily contacts were identified as significant predictors for TB classification .

Enhancing Model Interpretability: While RF is considered a "black-box" model, techniques like feature importance ranking and partial dependence plots have been employed to enhance interpretability, making the model more transparent for clinicians .

D.   Limitations and Challenges

Despite its promising performance, RF models face several challenges in TB prediction.

Data Quality and Availability: The performance of RF models heavily depends on the quality and quantity of data. Incomplete or biased datasets can lead to over fitting and reduced generalizability .

Interpretability: Although RF provides feature importance scores, the overall model can be perceived as a "black box," making it challenging to interpret and trust in clinical settings . Integration with Clinical Workflows: Implementing RF models in real-world clinical environments requires seamless integration with existing systems and workflows, which can be complex and resource-intensive .

Frontiers

E.   Future Directions

To address the limitations and enhance the applicability of RF

models in TB prediction, future research should focus on:

Enhancing Data Collection: Collaborating with healthcare institutions to collect diverse and comprehensive datasets, including clinical, laboratory, and imaging data.

Improving Model Interpretability: Developing techniques to make RF models more transparent and interpretable for clinicians, facilitating their adoption

in clinical practice.

Real-World Validation: Conducting prospective studies to validate the effectiveness of RF models in actual clinical settings, ensuring their reliability and robustness.

## V CONCLUSION

The application of Random Forest (RF) in the prediction and diagnosis of Tuberculosis (TB) represents a significant advancement in leveraging machine learning (ML) to improve healthcare outcomes. This review highlights the effectiveness of RF in handling diverse datasets, including clinical, laboratory, and imaging data, for diagnosing TB and differentiating between various TB-related conditions such as active and latent TB. Studies demonstrate that RF models can achieve high accuracy, with sensitivities and specificities often surpassing traditional diagnostic methods. Furthermore, RF's ability to handle missing data, identify important features, and generalize across various datasets makes it a powerful tool in medical diagnostics.

Despite these successes, several challenges remain, including the need for high-quality, balanced datasets, the model's interpretability, and the integration of RF models into existing clinical workflows. These hurdles can limit the full potential of RF-based systems in real-world applications, particularly in resource-limited settings. Future research should address these limitations by focusing on improving data quality, enhancing the transparency of model decisions, and validating RF models in diverse clinical environments.

In conclusion, while Random Forest demonstrates substantial promise as a tool for TB prediction, further work is required to optimize its performance, interpretability, and integration into clinical practice. With ongoing advancements in machine learning and the continuous availability of more diverse data, RF has the potential to revolutionize TB diagnosis, contributing to earlier detection, better treatment outcomes, and more effective control of this global health challenge.

## VI REFERENCE

[1]    Q. H. Li et al., "A united model for diagnosing pulmonary tuberculosis with random forest and artificial neural network," *Frontiers in Genetics*, vol. 14, p. 1094099, 2023.

[2]    Y. Q. Yao et al., "Development of a random forest model to classify sarcoidosis and tuberculosis," *American Journal of Translational Research*, vol. 13, no. 6, pp. 6166–6174, 2021.

[3]    P. Patel et al., "Machine learning approaches in diagnosing tuberculosis through biomarkers - A systematic review," *Journal of Clinical Microbiology*, vol. 61, no. 5, e02725-22, 2023.

[4] Y. Q. Yao et al., "Development of a random forest model to classify sarcoidosis and tuberculosis," *Am. J. Transl. Res.*, vol. 13, no. 6, pp. 6166–6174, Jun. 2021.

[5]    X. Zhang et al., "Development of diagnostic algorithm using machine learning for distinguishing between active tuberculosis and latent tuberculosis infection," *BMC Infect. Dis.*, vol. 22, no. 1, p. 1, 2022.

[6]    Q. H. Li et al., "A united model for diagnosing pulmonary tuberculosis with random forest and artificial neural network," *Front. Genet.*, vol. 14, p. 1094099, 2023.

[7]    P. Patel et al., "Machine learning approaches in diagnosing tuberculosis through biomarkers - A systematic review," *J. Clin. Microbiol.*, vol. 61, no. 5, e02725-22, 2023.

[8]    A. Fayaz et al., "Machine learning algorithms to predict treatment success for patients with pulmonary tuberculosis," *J. Clin. Microbiol.*, vol. 62, no. 4, e02725-23, 2024.