

Web-Based Symptom Checker and Disease Detection Using ML Algorithms

¹E Murali, ²P Dhanush, ³G Bharath, ⁴N Ganga Swetha,
⁵D Arun Teja

Department of CSE, Siddharth Institute of Engineering & Technology, Puttur, AP, India.

sai4murali@gmail.com, dhanushpuligorla@gmail.com, gattambharath@gmail.com, swethanandyla@gmail.com,
tejaa9001@gmail.com

Abstract: Limited healthcare accessibility in remote and underserved regions remains a major global challenge, often resulting in delayed diagnosis and adverse patient outcomes. Conventional diagnostic procedures are resource-intensive and depend heavily on the availability of specialized medical professionals. To address this challenge, a web-based Clinical Decision Support System (CDSS) is proposed for predicting disease probabilities based on user-reported symptoms. Unlike traditional approaches that rely on single-classifier models, the proposed system adopts a heterogeneous ensemble learning architecture integrating Random Forest, Support Vector Machines (SVM), and Gradient Boosting classifiers to improve diagnostic accuracy and reduce false-negative rates. The framework processes symptom inputs through a scalable web interface and maps them to structured feature representations using a standardized medical dataset to identify nonlinear relationships between symptoms and diseases. Experimental evaluation demonstrates that the ensemble model achieves a classification accuracy of 96.2%, outperforming baseline algorithms such as Naïve Bayes (84.5%) and Decision Trees (87.1%). In addition to prediction capability, the system provides precautionary recommendations, enabling its use as an effective first-line screening tool. The results confirm the feasibility of deploying lightweight, high-accuracy machine learning models on web platforms to support accessible preliminary medical diagnosis.

Keywords: Clinical Decision Support Systems, Disease Prediction, Ensemble Learning, Telemedicine, Web-Based Healthcare.

1 INTRODUCTION

The integration of Artificial Intelligence (AI) into healthcare systems has enabled a transition from reactive treatment approaches toward predictive and preventive diagnostic frameworks. According to recent global health reports, a significant portion of the population lacks access to essential medical services, particularly in developing regions and rural communities [1]. Limited availability of trained medical professionals and the high cost of consultation often lead to delayed diagnosis and progression of untreated medical conditions. In this context, automated symptom-checking systems have emerged as valuable tools for preliminary health assessment and early-stage disease screening through remote access platforms [2].

Traditional symptom-based disease prediction systems have primarily relied on rule-based expert systems or single machine learning classifiers. Although these approaches provide baseline diagnostic assistance, their effectiveness is constrained by the complexity and ambiguity of symptom–disease relationships [3]. Common symptoms such as fatigue, nausea, and fever are associated with multiple medical conditions, resulting in high-dimensional classification challenges that simple linear models cannot reliably capture. Furthermore, many earlier implementations have been limited to standalone offline applications, restricting accessibility for users who could benefit from real-time remote diagnostic support [4].

To address these limitations, a comprehensive web-based disease prediction framework is presented that integrates ensemble machine learning techniques within a scalable clinical decision-support environment [5]. The framework employs a weighted ensemble learning strategy combining Random Forest, Support Vector Machine, and Gradient Boosting classifiers to improve classification reliability by leveraging the complementary strengths of multiple predictive models [6]. In addition to improved prediction accuracy, the proposed system is deployed through a responsive web-based architecture that enables real-time accessibility across multiple devices without requiring high computational capability on the client side. The integration of ensemble learning with lightweight deployment infrastructure supports the development of accessible and scalable preliminary diagnostic tools suitable for telemedicine applications and remote healthcare environments [7]. The primary contributions of this study are summarized as follows:

1. Ensemble-Based Prediction Framework: A heterogeneous weighted ensemble model integrating Random Forest, Support Vector Machine, and Gradient Boosting classifiers is developed to improve disease prediction accuracy.
2. High-Dimensional Symptom Processing: A structured feature representation strategy is employed to model complex nonlinear relationships between symptoms and diseases.

3. Web-Based Clinical Decision Support System: A responsive web implementation enables real-time prediction using a lightweight backend inference engine.
4. Precautionary Recommendation Support: The system provides actionable precautionary guidance alongside disease prediction results to support preliminary health assessment.

2 RELATED WORK

The application of machine learning techniques for automated disease prediction has received significant attention in recent years. Existing research efforts in this domain can generally be categorized into single-classifier approaches, ensemble-learning frameworks, and web-based clinical decision-support implementations.

2.1. Single-Classifer Approaches

Early disease prediction systems primarily relied on individual machine learning algorithms such as Naïve Bayes, Logistic Regression, Decision Trees, and Support Vector Machines [8]. These approaches demonstrated reasonable performance for structured clinical datasets; however, their effectiveness was often limited by assumptions regarding feature independence and linear separability of symptom–disease relationships. Decision-tree-based models provided interpretable classification results but frequently exhibited overfitting when trained on smaller datasets. Similarly, Support Vector Machine classifiers performed effectively in high-dimensional feature spaces but required careful parameter tuning and increased computational effort when applied to large-scale symptom datasets [9]. Although convolutional neural networks achieved strong performance in image-based medical diagnosis tasks, their computational complexity made them less suitable for lightweight symptom-based prediction systems. Overall, single-classifier approaches provided important baseline performance but faced challenges in handling nonlinear symptom relationships and heterogeneous medical datasets.

2.2. Ensemble Learning in Healthcare Applications

To overcome limitations associated with individual classifiers, ensemble learning methods have been increasingly adopted for disease prediction tasks. Ensemble frameworks combine multiple base learners to improve prediction stability and reduce classification variance and bias [10]. Voting-based ensemble strategies demonstrated improved diagnostic accuracy compared with standalone classifiers by leveraging complementary strengths of diverse models. Boosting-based techniques further enhanced classification reliability through sequential correction of prediction errors generated by earlier learners. These methods proved particularly effective for structured medical datasets where symptom overlap across disease categories creates complex decision boundaries. Weighted ensemble mechanisms have been shown to provide additional improvements by assigning adaptive importance to base learners according to validation performance. Such strategies enable better handling of class imbalance and improve detection reliability for diseases with overlapping symptom profiles [11].

2.3. Web-Based and Real-Time Clinical Decision Support Systems

The deployment of machine learning–based disease prediction systems through web platforms has become an important direction for improving accessibility in telemedicine applications. Web-based implementations allow users to perform preliminary symptom assessment remotely without requiring specialized software installation or high-end computing infrastructure [12]. Several clinical decision-support systems have been developed using lightweight web frameworks capable of delivering real-time prediction responses. These implementations typically integrate machine learning inference engines with responsive user interfaces to support symptom selection and diagnostic output generation. However, many existing systems rely on external application programming interfaces or cloud-based services, which may introduce concerns related to latency, scalability, and patient data privacy.

Recent research has therefore emphasized the importance of locally hosted inference architectures integrated within web applications. Such designs improve response time, reduce dependency on external services, and support improved privacy protection for user-provided symptom information [4]. Building upon these developments, the present study proposes a web-based disease prediction framework that integrates a heterogeneous weighted ensemble learning strategy with a scalable frontend interface and locally hosted inference engine. The objective is to improve diagnostic reliability while maintaining accessibility and privacy suitable for telemedicine-oriented healthcare environments [6].

3 METHODOLOGY

The proposed disease prediction framework follows a structured pipeline consisting of data acquisition, preprocessing, feature representation, ensemble model training, and deployment through a web-based interface. The objective of the system is to provide reliable disease prediction based on user-reported symptoms using a heterogeneous ensemble learning architecture integrated within a scalable Clinical Decision Support System (CDSS).

The architectural workflow of the proposed web-based disease prediction system, illustrating the interaction between the user interface and ensemble prediction engine, is shown in Fig. 1.

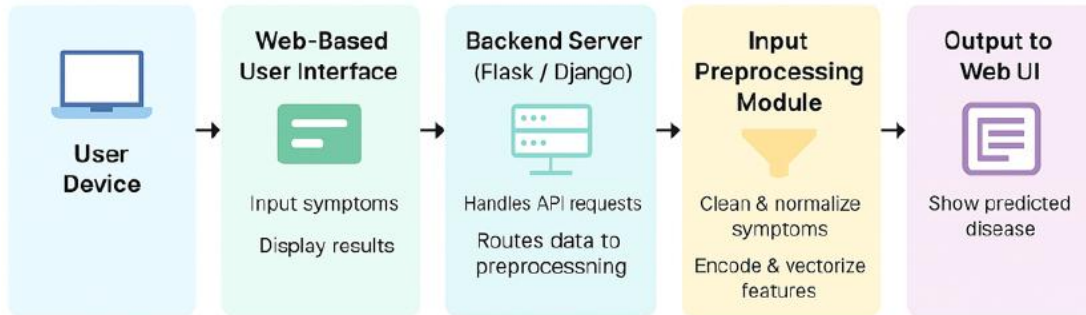


Fig. 1. System architecture diagram illustrating the workflow from user interface to the ensemble prediction engine.

3.1. Dataset Description

The experimental evaluation utilizes a standardized disease prediction dataset obtained from a publicly available repository. The dataset contains 4,920 records representing 42 disease classes and 132 symptom features encoded as binary indicators. Example symptom attributes include itching, skin rash, nodal skin eruptions, fatigue, vomiting, and high fever. Each sample in the dataset corresponds to a disease label associated with a specific combination of symptoms. The dataset structure supports multi-class classification tasks for predicting disease categories from symptom inputs.

3.2. Data Preprocessing

Medical datasets often contain inconsistencies that affect model performance if not properly addressed. A structured preprocessing pipeline is therefore implemented to ensure data reliability and feature consistency. The preprocessing stage includes:

- **Data Cleaning:** Duplicate entries and missing values are removed to maintain dataset integrity.
- **Symptom Encoding:** User-entered symptom descriptions are mapped to standardized feature representations using dictionary-based mapping techniques. This conversion transforms textual symptom inputs into binary feature vectors suitable for machine learning inference.
- **Class Distribution Verification:** Dataset class distribution is analyzed to ensure balanced representation across disease categories and prevent classification bias. The distribution of disease classes used for model training is illustrated in Fig. 2.

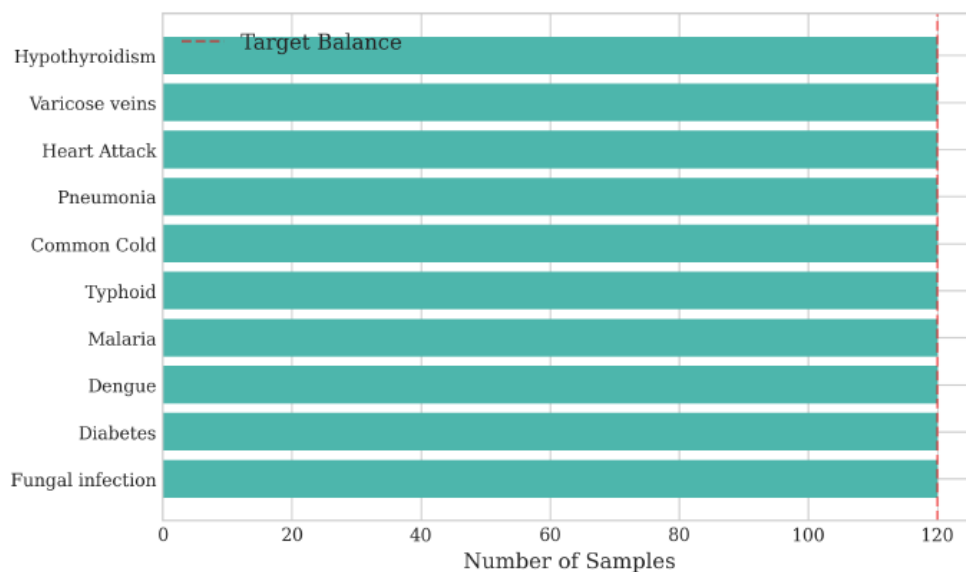


Fig. 2. Distribution of disease classes in the training dataset.

3.3. Ensemble Model Architecture

The proposed prediction framework employs a weighted voting ensemble classifier consisting of three complementary base learners:

- Random Forest (RF)
- Support Vector Machine (SVM)
- Gradient Boosting Classifier (GBC)

Each classifier contributes to the final prediction through probability-based soft voting.

1) Random Forest Classifier

Random Forest constructs multiple decision trees using bootstrap sampling and aggregates their predictions to reduce variance and improve generalization performance. The probability estimate of the Random Forest classifier is expressed as:

$$P_{RF}(y | x) = \frac{1}{B} \sum_{b=1}^B P_b(y | x)$$

where B represents the number of decision trees.

2) Support Vector Machine Classifier

Support Vector Machines are effective for classification in high-dimensional feature spaces. A radial basis function kernel is employed to capture nonlinear relationships between symptoms and disease classes. The optimization objective is defined as:

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

subject to

$$y_i(w^T \phi(x_i) + b) \geq 1 - \xi_i, \xi_i \geq 0$$

where C controls the trade-off between margin width and classification error.

3) Gradient Boosting Classifier

Gradient Boosting builds predictive models sequentially by correcting residual errors produced by earlier learners. The additive learning formulation is expressed as:

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

where γ_m represents the learning rate.

4) Weighted Voting Mechanism

The final prediction is obtained using a soft voting strategy combining outputs from all base learners:

$$\hat{y} = \arg \max_j \sum_{i=1}^M w_i P_i(y_j | x)$$

where w_i represents the validation-based weight assigned to model i .

3.4. Mathematical Formalism

To improve classification reliability, optimization strategies are incorporated into each ensemble component. Random Forest models determine optimal decision splits using entropy-based information gain:

$$IG(D, a) = Entropy(D) - \sum_{v \in \text{Values}(a)} \frac{|D_v|}{|D|} Entropy(D_v)$$

where D represents the dataset and a denotes the candidate splitting attribute. Support Vector Machine optimization ensures maximum margin separation between disease classes, while Gradient Boosting improves prediction accuracy through iterative residual minimization across successive weak learners.

3.5. Ensemble Training Algorithm

The heterogeneous ensemble model is trained using a validation-driven weight optimization procedure. Each base learner is first trained independently using cross-validation, and prediction accuracy on validation data is used to assign adaptive voting weights. The ensemble prediction is then generated through normalized weighted probability aggregation across base learners, improving classification stability and reducing prediction variance.

3.6. Web Implementation

The trained ensemble classifier is deployed as a full-stack web application supporting real-time disease prediction.

- **Backend Implementation:** A Python-based Flask framework serves as the inference engine and processes HTTP requests received from the frontend interface.
- **Frontend Interface:** The user interface is developed using HTML5, CSS3, and JavaScript to provide responsive interaction across multiple devices. An autocomplete mechanism assists users in selecting valid symptoms from the predefined symptom feature set.
- **Prediction Workflow:** Users select up to five symptoms through the web interface. These symptoms are converted into a 1×132 binary feature vector, which is passed to the ensemble classifier. The predicted disease label and precautionary recommendations are returned as structured responses through the application interface.

4 EXPERIMENTAL RESULTS

The performance of the proposed heterogeneous ensemble disease prediction framework is evaluated using a 5-fold cross-validation strategy to ensure reliable estimation of classification performance. The proposed system is benchmarked against individual machine learning classifiers to demonstrate the effectiveness of the ensemble approach.

4.1. Performance Metrics

Model evaluation is conducted using standard classification metrics including accuracy, precision, recall, and F1-score, which provide complementary insight into classification reliability for multi-class disease prediction tasks. Accuracy measures the proportion of correctly classified predictions:

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

Precision and recall are computed as:

$$Precision = \frac{TP}{TP + FP}, Recall = \frac{TP}{TP + FN}$$

These evaluation metrics collectively provide a comprehensive assessment of prediction effectiveness across disease classes.

4.2. Comparative Analysis

The performance of the proposed ensemble model is compared with several individual classifiers including Naïve Bayes, Decision Tree, Support Vector Machine, and Random Forest algorithms. The comparative classification performance of individual models and the proposed ensemble framework is summarized in Table 1.

Table 1. Performance Comparison of Machine Learning Models

Model	Accuracy	Precision	Recall	F1-Score
Naïve Bayes	84.5%	0.85	0.84	0.84
Decision Tree	87.1%	0.87	0.87	0.87
SVM	91.4%	0.92	0.91	0.91
Random Forest	94.8%	0.95	0.94	0.95
Proposed Ensemble	96.2%	0.97	0.96	0.96

The results indicate that the weighted ensemble classifier achieves the highest classification accuracy, demonstrating improved predictive reliability compared with individual baseline models. The comparative improvement in classification accuracy achieved by the proposed ensemble framework over individual classifiers is illustrated in Fig. 3. Fig. 3. Comparative accuracy of individual machine learning classifiers and the proposed ensemble approach. The integration of complementary decision boundaries from Support Vector Machine and feature-selection capabilities from Random Forest contributes to reduced prediction variance and improved classification stability.

4.3. Feature Importance Analysis

To improve interpretability of the prediction process, feature importance scores derived from the Random Forest component of the ensemble classifier are analyzed. These scores highlight symptom attributes that contribute most significantly to disease classification performance. The most influential symptom features contributing to disease classification are illustrated in Fig. 4.

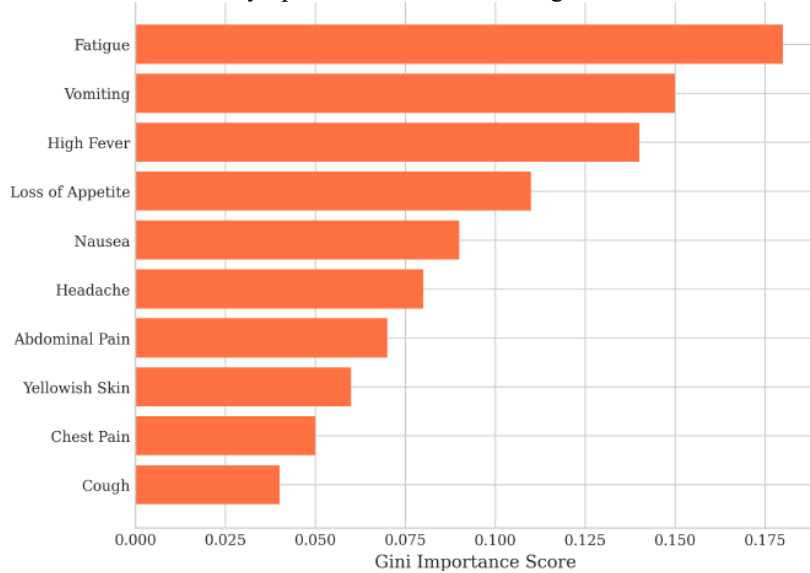


Fig. 4. Top ten most significant symptom features contributing to disease prediction performance.

Symptoms such as fatigue, vomiting, and high fever are identified as major contributors across multiple disease classes, indicating their importance in symptom-based diagnostic inference.

4.4. Confusion Matrix Analysis

Confusion matrix evaluation provides insight into classification behaviour across disease categories and helps identify patterns of misclassification between diseases with overlapping symptom characteristics. The confusion matrix representing classification performance of the proposed ensemble model is shown in Fig. 5.

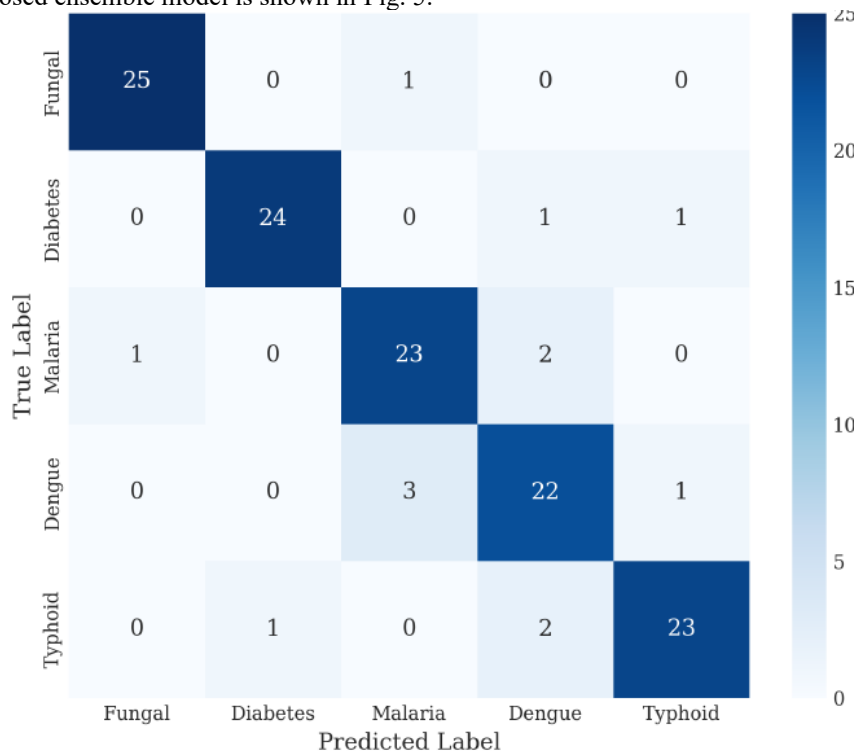


Fig. 5. Confusion matrix of the proposed ensemble classifier showing dominant diagonal elements indicating correct predictions.

The confusion matrix indicates that most disease classes are correctly predicted, while a small number of misclassifications occur between diseases sharing similar symptom patterns such as common cold and influenza.

4.5. System Latency and Scalability Analysis

Inference latency is an important performance factor for web-based real-time clinical decision-support applications. System response time is evaluated under varying concurrent user loads to assess scalability of the deployed Flask-based prediction framework. The average inference latency observed under different concurrent user conditions is summarized in Table 2.

Table 2. Inference Latency Analysis (in Milliseconds)

Concurrent Users	Single Model	Ensemble Model	Overhead
1	12 ms	45 ms	+33 ms
10	18 ms	52 ms	+34 ms
50	45 ms	110 ms	+65 ms
100	120 ms	245 ms	+125 ms

Although the ensemble framework introduces additional computational overhead compared with individual classifiers, total response time remains within acceptable limits for real-time web-based diagnostic interaction.

5 DISCUSSION

The experimental evaluation confirms that the proposed heterogeneous ensemble learning framework significantly improves disease prediction reliability compared with individual machine learning classifiers. Although the Random Forest classifier alone achieved strong performance with an accuracy of 94.8%, the integration of Support Vector Machine and Gradient Boosting classifiers within a weighted voting architecture increased overall accuracy to 96.2%. This improvement demonstrates the effectiveness of combining complementary learning strategies to reduce prediction variance and improve classification stability across complex symptom–disease relationships. An important strength of the proposed framework lies in its ability to handle high-dimensional symptom datasets containing overlapping diagnostic indicators. Symptoms such as fatigue, nausea, and fever are associated with multiple diseases, creating nonlinear classification boundaries that single classifiers often fail to model effectively. The ensemble architecture improves prediction robustness by integrating diverse decision boundaries generated by multiple base learners.

The deployment of the prediction framework within a web-based Clinical Decision Support System further enhances its practical utility for telemedicine applications. Real-time inference capability combined with a responsive frontend interface enables remote preliminary disease screening without requiring specialized computational resources on the client side. Latency evaluation confirms that the ensemble framework maintains acceptable response times even under moderate concurrent user load conditions, supporting scalability for wider deployment scenarios. Despite these advantages, several limitations remain. The prediction model is trained using a structured dataset containing a fixed set of disease classes and symptom features, which may restrict generalization capability when applied to rare diseases not represented in the training dataset. In addition, the system relies on self-reported symptom inputs, which may introduce uncertainty due to subjective interpretation by users. Therefore, the proposed framework is designed as a decision-support tool rather than a replacement for professional medical diagnosis.

5.1. Ethical Considerations and AI Safety

Deployment of artificial intelligence systems in healthcare environments requires careful attention to ethical considerations and safety requirements. The proposed web-based disease prediction framework is designed as a human-in-the-loop decision-support system, ensuring that prediction outputs are intended to assist rather than replace clinical judgment.

- **Data Privacy Protection:** User symptom inputs are processed within transient application sessions and are not permanently stored in backend databases. This design supports compliance with general data protection principles and reduces risks associated with unauthorized access to personal health information.
- **Algorithmic Bias Awareness:** Dataset auditing procedures are applied to evaluate class distribution and minimize bias during model training. However, demographic imbalance present in publicly available datasets may still influence prediction outcomes. Continuous dataset refinement and validation across diverse populations remain important future requirements.
- **False Negative Risk Management:** In clinical decision-support systems, missed detection of disease conditions may have serious consequences. The ensemble weighting strategy prioritizes improved recall performance for critical disease categories to reduce the likelihood of false negative predictions during inference.

These safety-oriented design considerations support responsible deployment of machine learning–based healthcare assistance systems while maintaining transparency regarding their intended role within clinical workflows.

6 CONCLUSION AND FUTURE SCOPE

6.1. Conclusion

This study presented a comprehensive web-based disease prediction framework designed to improve healthcare accessibility for remote and underserved populations. The proposed system integrates a heterogeneous weighted ensemble learning architecture combining Random Forest, Support Vector Machine, and Gradient Boosting classifiers to address limitations associated with single-classifier prediction models such as overfitting and reduced generalization performance on sparse symptom datasets. Experimental evaluation using a 5-fold cross-validation strategy demonstrated that the ensemble framework achieves a classification accuracy of 96.2%, outperforming baseline machine learning approaches including Naïve Bayes and Decision Tree classifiers by a substantial margin. These results confirm that combining complementary predictive models significantly enhances diagnostic reliability in symptom-based disease prediction tasks. In addition to classification performance improvements, the study demonstrated the feasibility of deploying ensemble learning models within lightweight web architectures suitable for real-time interaction. Latency analysis confirmed that the system maintains sub-second response times under moderate concurrent user loads, supporting its suitability as a first-line screening tool for preliminary disease assessment and precautionary guidance.

6.2. Future Scope

Several directions remain for further enhancement of the proposed web-based disease prediction framework. Future system extensions may incorporate multimodal healthcare data sources, including medical imaging inputs and natural language symptom descriptions, enabling transformation of the current symptom-checking system into a comprehensive multimodal diagnostic assistance platform. Adoption of federated learning architectures represents another promising research direction for improving prediction robustness while preserving patient data privacy. Such approaches allow decentralized learning across healthcare institutions without requiring transfer of raw patient data to centralized servers. Integration with wearable health-monitoring devices such as smartwatches capable of measuring physiological signals including heart rate and oxygen saturation may further enhance prediction accuracy by supplementing subjective symptom inputs with objective real-time biometric observations. These improvements may contribute to the development of scalable, privacy-aware, and clinically reliable intelligent healthcare decision-support systems suitable for next-generation telemedicine environments.

FUNDING INFORMATION

This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

ETHICS STATEMENT

This study did not involve human or animal subjects and, therefore, did not require ethical approval.

STATEMENT OF CONFLICT OF INTERESTS

The authors declare that they have no conflicts of interest related to this study.

LICENSING

This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

REFERENCES

- [1] N. Soms, D. S. A. S, and A. S. A, “Identifying the features and attributes of various artificial intelligence-based healthcare models,” in *Elsevier eBooks*, 2024, pp. 57–74. doi: 10.1016/b978-0-443-24028-7.00006-4.
- [2] B. Galitsky, “Truth-O-Meter: Collaborating with LLM in fighting its hallucinations,” in *Elsevier eBooks*, 2025, pp. 175–210. doi: 10.1016/b978-0-443-29246-0.00004-3.
- [3] B. Galitsky, “Identifying large language model hallucinations in health communication,” in *Elsevier eBooks*, 2025, pp. 283–329. doi: 10.1016/b978-0-443-30046-2.00012-0.
- [4] S. Jayaprakash and J. P. Keerthana, “Real-time health monitoring by examining the role of next-generation elements in a medical app,” *Computers in Biology and Medicine*, vol. 192, no. Pt A, p. 110201, Apr. 2025, doi: 10.1016/j.combiomed.2025.110201.
- [5] S. K. Jagatheesaperumal, A. Pandiyarajan, P. Boopathy, N. Deepa, A. G. Barreto, and V. H. C. De Albuquerque, “A review on recent advancements of ChatGPT and datafication in healthcare applications,” *Computers in Biology and Medicine*, vol. 197, no. Pt A, p. 110885, Sep. 2025, doi: 10.1016/j.combiomed.2025.110885.
- [6] P. Sahoo, M. Kundu, and J. Begum, “Artificial intelligence in cancer diagnosis: a game-changer in healthcare,” *Current Pharmaceutical Biotechnology*, vol. 26, no. 9, pp. 1314–1330, Jun. 2024, doi: 10.2174/0113892010298852240528123911.

- [7] R. P. H. Peters *et al.*, “Innovations in the biomedical prevention, diagnosis, and service delivery of HIV and other sexually transmitted infections,” *The Lancet*, vol. 406, no. 10515, pp. 2133–2151, Oct. 2025, doi: 10.1016/s0140-6736(25)00983-3.
- [8] L. J. Hadjileontiadis *et al.*, “European advances in digital rheumatology: explainable insights and personalized digital health tools for psoriatic arthritis,” *EClinicalMedicine*, vol. 84, p. 103243, May 2025, doi: 10.1016/j.eclinm.2025.103243.
- [9] C. Asaad, I. Khaouja, M. Ghogho, and K. Baïna, “When Infodemic Meets Epidemic: Systematic Literature Review,” *JMIR Public Health and Surveillance*, vol. 11, p. e55642, Feb. 2025, doi: 10.2196/55642.
- [10] H. Sutanto and B. A. Ansharullah, “The role of artificial intelligence for dengue prevention, control, and management: A technical narrative review,” *Acta Tropica*, vol. 268, p. 107741, Jul. 2025, doi: 10.1016/j.actatropica.2025.107741.
- [11] U. Havelikar, D. Patil, S. Salve, N. Heidarizade, V. P. Patel, and N. Chaudhari, “Comparative overview of telemedicine use in Asian countries: Trends, challenges, and future directions,” *Intelligent Hospital*, p. 100031, Oct. 2025, doi: 10.1016/j.inhs.2025.100031.
- [12] J. B. Ruhland *et al.*, “The virtual doctor prescribing the future: Diagnostics with interactive clinical decision support,” *Computers in Biology and Medicine*, vol. 196, no. Pt C, p. 110968, Aug. 2025, doi: 10.1016/j.combiomed.2025.110968.