

Prediction of Alzheimer's Disease Using Machine Learning

Dharshan M S¹, Dr. K. Santhi²

Department Of Information Technology, Dr. N.G.P Arts and Science College, Coimbatore, Tamil Nadu, India¹

Professor, Department of Information Technology, Dr. N.G.P Arts and Science College, Coimbatore, Tamil Nadu, India²

Abstract: Alzheimer's disease (AD) is a progressive neurodegenerative disorder and the leading cause of dementia worldwide, affecting approximately 50 million people with projections exceeding 150 million by 2050. Early detection of AD is critically important as interventions initiated during the prodromal phase—particularly mild cognitive impairment—have the greatest potential to slow disease progression and preserve cognitive function. Machine learning, particularly ensemble methods like Random Forest, has emerged as a powerful tool for early AD prediction by analysing complex, multimodal data including neuroimaging, genetic markers, cognitive assessments, and fluid biomarkers. This paper provides a comprehensive review of Random Forest applications in AD prediction, synthesizing findings from recent studies. The combination of Backward Elimination Feature Selection with Artificial Ant Colony Optimization has achieved 95% accuracy while reducing computation time by 81%. Key challenges including class imbalance, model interpretability, and cross-cohort generalizability are addressed through techniques such as SMOTE and SHAP analysis. This review provides researchers and clinicians with a comprehensive understanding of Random Forest's role in early AD prediction and identifies promising directions for future research.

Keywords: Alzheimer's disease, machine learning, Random Forest, early prediction, ensemble learning, feature selection, biomarker analysis, neuroimaging, cognitive assessment

I. INTRODUCTION

Alzheimer's disease (AD) represents the predominant cause of dementia worldwide, affecting approximately 1 in 9 people aged 65 and older in the United States alone, with prevalence increasing exponentially with age. The disease is characterized by progressive cognitive decline, memory loss, and functional impairment, with neuropathological features including amyloid-beta plaque accumulation, tau protein hyperphosphorylation, synaptic dysfunction, and neuronal loss. Despite decades of research, there is currently no cure for AD; however, prescribed medications can mitigate the progression of the condition, making early identification critically important for treatment and further research.

The pathological changes associated with AD begin 10-15 years before clinical symptoms manifest, yet current diagnostic methods typically identify the disease at moderate to advanced stages when significant neurodegeneration has already occurred. Traditional diagnostic approaches rely on clinical assessments, cognitive testing, and neuroimaging, but these methods often identify disease only after substantial brain damage has accumulated. The key obstacles to the early assessment of Alzheimer's disease include the limited availability of well-characterized training samples, the high dimensionality of feature descriptions, the heterogeneity of disease presentation, and the lack of trustworthy AI-based solutions for detecting this disease in clinical settings.

Machine learning (ML) techniques offer promising solutions to these challenges by identifying subtle patterns in high-dimensional data that may precede overt clinical symptoms by years. Among various ML algorithms, Random Forest (RF) has emerged as particularly well-suited for AD prediction due to its ability to handle high-dimensional data, manage non-linear relationships, provide intrinsic feature importance rankings, resist overfitting through ensemble learning, and maintain robustness to missing data and outliers.

II. LITERATURE REVIEW

Alzheimer's disease (AD), the primary cause of dementia, involves progressive cognitive decline and neuropathological features like amyloid-beta plaques and tau tangles, with changes often starting 10–15 years before symptoms. Early detection is vital for interventions but challenging due to limitations in traditional methods (clinical tests, neuroimaging, biomarkers, genetics). Random Forest (RF), an ensemble algorithm using bagging and random feature selection for

majority-vote predictions, excels in AD prediction by handling high-dimensional multimodal data, capturing non-linear patterns, ranking feature importance, resisting overfitting, and managing missing values. Key hyperparameters (e.g., number of trees, max_features) benefit from optimization via nature-inspired methods. Recent 2025–2026 studies show strong performance: Soladoye et al. (2025) combined Backward Elimination with Ant Colony Optimization on ~2,149 samples, achieving ~95% accuracy/precision/F1, 94% recall, and 98% AUC, identifying 26 key features (demographics, cognition, genetics) while cutting computation time by 81%. Other works report RF accuracies of 86–96% in multi-stage classification (e.g., NACC data, AUC ~0.95), survival analysis (concordance ~0.84), and longitudinal EHR keyword-based prediction (AUROC up to 0.86 near diagnosis, ~0.58–0.60 at 10+ years pre-onset). Top predictors often include hippocampus/entorhinal volumes, MMSE/CDR scores, with sex/age-specific patterns (e.g., right hippocampus in younger groups). Implementation uses datasets like ADNI/NACC/OASIS, preprocessing (missing data handling, encoding, stratified splits), cross-validation, and metrics (accuracy, F1, AUC-ROC), enhanced by interpretability tools (SHAP, PDPs). Challenges include data heterogeneity, missingness, bias risks, and the need for external validation and equitable performance across demographics. Optimized RF models offer interpretable, efficient early AD detection, often outperforming alternatives in tabular/multimodal settings, with promise for clinical tools pending prospective validation.

III. ALZHEIMER'S DISEASE: PATHOPHYSIOLOGY AND DIAGNOSTIC CHALLENGES

A. Disease Pathophysiology

Alzheimer’s disease is a progressive neurodegenerative disorder that affects neurotransmitters, neurons, and neural tissue, thereby impairing sensory perception, memory, and behavior. The neuropathological hallmarks of AD include extracellular amyloid-beta plaques, intracellular neurofibrillary tangles composed of hyperphosphorylated tau protein, synaptic loss, and neuronal degeneration. These pathological changes follow a characteristic spatial and temporal progression, typically beginning in the entorhinal cortex and hippocampus before spreading to neocortical association areas.

The amyloid cascade hypothesis posits that accumulation of amyloid-beta peptide triggers a cascade of events including tau hyperphosphorylation, oxidative stress, neuroinflammation, and ultimately neuronal death. However, the relationship between these pathological processes and clinical symptoms is complex, with substantial inter-individual variability in the rate of progression and the specific cognitive domains affected.

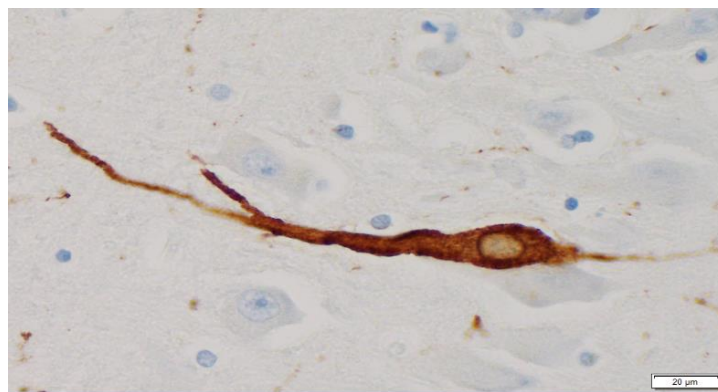


Fig 1: Tau Neurofibrillary Tangles

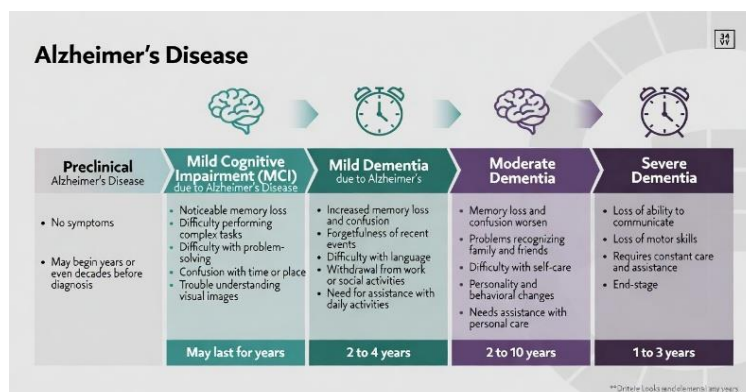


Fig 2: Alzheimer's progression chart

B. Clinical Stages and Early Detection Windows

The clinical progression of AD is typically conceptualized as a continuum from cognitively normal (CN) through mild cognitive impairment (MCI) to dementia. Mild cognitive impairment represents an intermediate state in which individuals have objective cognitive deficits but maintain functional independence. The annual conversion rate from MCI to AD dementia is approximately 10-15%, though this varies substantially based on the presence of biomarkers and other risk factors

Early detection is critically important for several reasons. First, interventions initiated during the MCI stage have the greatest potential to slow disease progression. Second, accurate identification of individuals at high risk for progression enables more efficient design of clinical trials for disease-modifying therapies. Third, early diagnosis allows patients and families to plan for future care needs and participate in decisions about treatment options.

C. Current Diagnostic Modalities

Current diagnostic approaches for AD integrate multiple data sources:

Neuroimaging: Structural magnetic resonance imaging (MRI) detects brain atrophy patterns characteristic of AD, particularly in medial temporal lobe structures including the hippocampus and entorhinal cortex. Functional imaging modalities such as fluorodeoxyglucose positron emission tomography (FDG-PET) reveal metabolic deficits, while amyloid PET (AV45) and tau PET provide direct visualization of pathological protein deposition

Fluid Biomarkers: Cerebrospinal fluid (CSF) measurements of amyloid-beta 42 (A β 42), total tau (t-tau), and phosphorylated tau (p-tau) have high diagnostic accuracy for AD pathology. Recent advances in ultra-sensitive assays have enabled measurement of these biomarkers in blood, offering a less invasive approach suitable for screening.

Cognitive Assessment: Standardized neuropsychological tests including the Mini-Mental State Examination (MMSE), Montreal Cognitive Assessment (MoCA), and Clinical Dementia Rating (CDR) provide quantitative measures of cognitive function across multiple domains. These assessments are inexpensive and widely available but may lack sensitivity to early pathological changes.

Genetic Markers: The apolipoprotein E (APOE) ϵ 4 allele is the strongest genetic risk factor for late-onset AD, with each copy increasing risk approximately threefold. Genome-wide association studies have identified numerous additional risk variants with smaller effect sizes

The integration of these diverse data modalities presents both opportunities and challenges for machine learning approaches, as the complementary information they provide must be effectively combined while managing missing data, different scales, and complex interactions.

IV. RANDOM FOREST ALGORITHM: THEORETICAL FOUNDATIONS**A. Ensemble Learning Principles**

Random Forest, introduced by Leo Breiman in 2001, is an ensemble learning method that operates by constructing multiple decision trees during training and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. The algorithm is built upon two fundamental principles of ensemble learning: bootstrap aggregating (bagging) and random feature selection.

Prediction Process:

For classification, the final prediction for a new instance x is the class that receives the majority vote among all trees:

$$\hat{h}(x) = \text{majority vote } \{h_b(x)\}_{b=1}^B$$

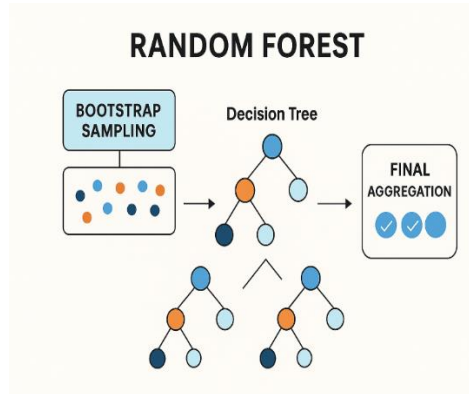


Fig 3 :Random forest Flow

B. Key Characteristics and Advantages for Medical Applications

Random Forest possesses several characteristics that make it particularly well-suited for medical prediction tasks such as AD detection:

High-Dimensional Data Handling:

Random Forest can effectively handle datasets with many features (p) relative to samples (n), a common situation in medical research where hundreds or thousands of potential predictors may be measured on a relatively small number of patients (Kuhn & Johnson, 2013).

Non-Linear Relationships: Decision trees can capture complex non-linear relationships and interactions between features without requiring explicit specification of functional forms, making them well-suited for modeling the complex pathophysiology of AD.

Feature Importance Ranking: The algorithm provides intrinsic measures of feature importance, either through mean decrease in impurity (how much each feature reduces impurity across all splits) or mean decrease in accuracy (how much prediction accuracy decreases when a feature is permuted). This property is particularly valuable for identifying the most informative biomarkers and understanding disease mechanisms.

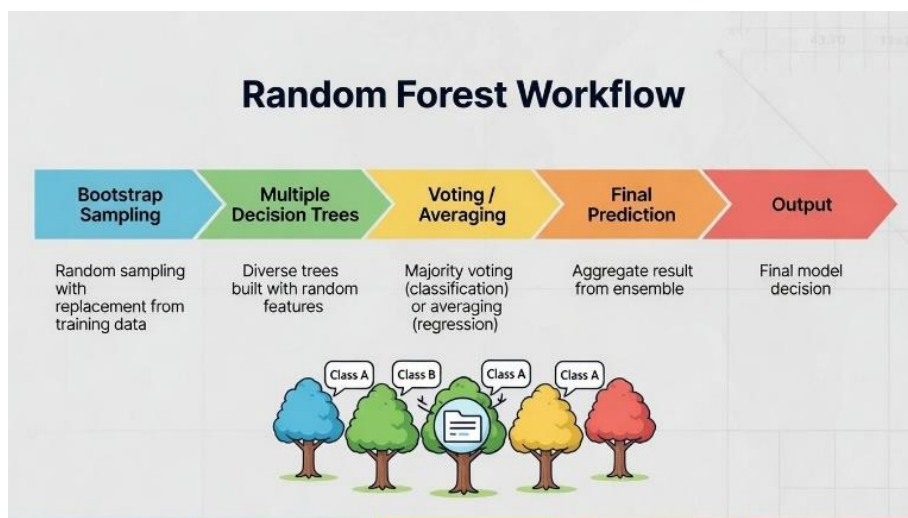


Fig 4: Random forest workflow

Robustness to Overfitting: The ensemble nature of Random Forest, combined with bootstrap sampling and random feature selection, makes it resistant to overfitting, even with many features and complex relationships.

Handling Missing Data: Random Forest can accommodate missing data through proximity-based imputation or by using surrogate splits, reducing the need for complex imputation procedures.

V. RECENT ADVANCES IN RANDOM FOREST FOR AD PREDICTION

A. State-of-the-Art Performance

Recent research has significantly advanced RF-based AD prediction through sophisticated feature selection and hyperparameter optimization techniques. A landmark study by Soladoye et al. (2025) proposed a novel framework combining Backward Elimination Feature Selection with Artificial Ant Colony Optimization for Random Forest hyperparameter tuning.

The study utilized a dataset of 2,149 instances with 34 features representing demographic information, cognitive assessments, genetic markers, and clinical variables. Preprocessing included Min Max normalization and SMOTE (Synthetic Minority Oversampling Technique) applied to training data only to prevent data leakage. Multiple feature selection techniques were evaluated, including the Whale Optimization Algorithm, Artificial Bee Colony, and Backward Elimination Feature Selection. Hyperparameter optimization algorithms compared included Artificial Ant Colony Optimization and Bald Eagle Search.

The combination of Backward Elimination Feature Selection with Artificial Ant Colony Optimization achieved the highest performance metrics:

- Accuracy: **95% ± 1.2%**
- Precision: **95% ± 1.1%**
- Recall: **94% ± 1.3%**
- F1-score: **95% ± 1.0%**
- Area Under the Curve (AUC): **98% ± 0.8%**

This approach identified **26 significant features** associated with Alzheimer's disease and demonstrated substantial computational efficiency advantages over empirical grid search approaches (18 minutes versus 133 minutes, representing an **81% reduction** in computation time). The identified features spanned multiple domains, confirming the importance of multimodal assessment for accurate AD prediction.

B. Multi-Class Classification of AD Stages

While much research focuses on binary classification (AD versus non-AD), clinical decision-making requires differentiation among multiple disease stages. Bulut et al. (2025) addressed this need by focusing on determining four different ordinal stages of Alzheimer's disease using multiple machine learning algorithms.

Their study utilized data from the National Alzheimer's Coordinating Center (NACC) database, incorporating demographic information, genetic markers, neurocognitive inventory data, and brain volume/thickness measurements from MRI scans. Models compared included Deep Neural Networks, Ordinal Logistic Regression, Random Forest, Gaussian Naive Bayes, XG Boost, and Light GBM.

The highest classification rate of AD stages was achieved by the Random Forest model:

- Accuracy: **0.86**
- F1 score: **0.86**
- AUC: **0.95**

Importantly, the study employed **SHapley Additive ex Planations (SHAP)** to explain model outputs, indicating that non-invasive markers and machine learning models can be used effectively in early diagnosis and decision support systems. This focus on interpretability is crucial for clinical adoption, as physicians must understand and trust model recommendations.

C. Feature Selection and Dimensionality Reduction

Effective feature selection is essential for Alzheimer's disease (AD) prediction because medical datasets are often high dimensional and clinical adoption requires interpretability. The **Applied Predictive Modeling** textbook by **Max Kuhn** and **Kjell Johnson** provides detailed demonstrations of feature selection strategies using the Alzheimer Disease dataset from the Applied Predictive Modeling package. A central approach is Recursive Feature Elimination (RFE), which iteratively removes less informative predictors based on cross-validated model performance to identify an optimal subset of features. To ensure stability and reduce variance in the selection process, repeated k-fold cross-validation—commonly five repeats of 10-fold CV—is employed. The framework also emphasizes comparing multiple classifiers, including Random Forest, Linear Discriminant Analysis, Support Vector Machine, Naive Bayes, Logistic Regression, and K-Nearest Neighbor, to determine which algorithm benefits most from the selected feature set. Model evaluation relies on

robust performance metrics such as ROC, sensitivity, and specificity, enabling balanced assessment of discrimination ability and clinical relevance.

D. Sex-Specific and Age-Specific Patterns

A nuanced understanding of AD requires consideration of how prediction models perform across demographic subgroups. Employed the Random Forest algorithm on numerical data derived from anatomical MRI scans, achieving 92.87% accuracy in detecting AD from MCI and cognitively normal individuals.

Critically, their subgroup analyses across nine sex- and age-based cohorts revealed important differences in predictive patterns. The hippocampus, amygdala, and entorhinal cortex emerged as consistent top-rank predictors across all groups. However, younger males and females (aged 69-76) exhibited volume decreases in the right hippocampus, suggesting its importance in early AD stages. Older males (77-84) showed substantial volume decreases in the left inferior temporal cortex, while the left middle temporal cortex demonstrated decreased volume specifically in females.

These findings indicate potential sex-specific neuroanatomical patterns in AD progression and highlight the importance of developing and validating prediction models that perform equitably across demographic subgroups. Failure to account for such differences could lead to biased predictions and disparities in early detection.

VI. METHODOLOGY FOR RANDOM FOREST IMPLEMENTATION IN AD PREDICTION

A. Data Sources and Preparation

Successful implementation of Random Forest for Alzheimer's disease (AD) prediction begins with careful data selection and preparation. Commonly used datasets include the **Alzheimer's Disease Neuroimaging Initiative (ADNI)**, a longitudinal multicenter study that provides rich multimodal data such as MRI, PET, cerebrospinal fluid biomarkers, genetics, and cognitive assessments, making it the most widely used resource in AD research. The **National Alzheimer's Coordinating Center (NACC)** maintains a large standardized database of clinical and neuropathological data collected from Alzheimer's Disease Centers funded by the U.S. National Institute on Aging and is frequently used for multi-stage classification studies. The **Open Access Series of Imaging Studies (OASIS)** offers cross-sectional and longitudinal MRI datasets for normal aging and AD, widely applied in neuroimaging-based prediction work. In addition, the **Kaggle Alzheimer's Disease dataset**, containing 35 features and 2,149 samples, is often used for benchmarking machine learning algorithms. Data preparation typically involves several critical steps. Handling missing values is essential; although Random Forest can tolerate some missing data, systematic gaps should be examined and treated using approaches such as casewise deletion, median or mode imputation, multiple imputation, or model-based methods. Feature encoding must be applied appropriately—label encoding for ordinal variables and one-hot encoding for nominal categorical variables such as genotype. Finally, the dataset should be divided into training (about 70–80%) and testing (20–30%) subsets using stratified sampling to maintain class balance, with cross-validation on the training portion employed for robust model development and hyperparameter tuning.

B. Feature Selection Techniques

Feature selection serves multiple purposes in AD prediction: improving model performance, reducing computational requirements, enhancing interpretability, and identifying potentially novel biomarkers. Approaches include:

Filter Methods: Evaluate feature relevance independently of the model using statistical measures such as correlation, mutual information, or chi-square tests. These methods are computationally efficient but may miss feature interactions.

Wrapper Methods: Evaluate feature subsets using the model's performance. Recursive Feature Elimination (RFE) iteratively removes the least important features and assesses model performance, providing a robust but computationally intensive approach.

Embedded Methods: Leverage the model's intrinsic feature importance rankings. Random Forest provides importance measures based on mean decrease in impurity or mean decrease in accuracy, enabling feature selection within the modeling process.

Hybrid Approaches: Recent work has combined multiple techniques, such as using Backward Elimination (a wrapper method) followed by Ant Colony Optimization for simultaneous feature and hyperparameter optimization

C. Model Evaluation and Validation

A comprehensive evaluation of Alzheimer’s disease (AD) prediction models requires the use of multiple performance metrics and rigorous validation strategies to ensure reliability and clinical usefulness. Accuracy offers a general measure of correct predictions but may be misleading in the presence of class imbalance; therefore, precision, recall (sensitivity), and F1-score provide a more balanced assessment of classification performance. The Area Under the Receiver Operating Characteristic Curve (AUC-ROC) evaluates the model’s discrimination capability across all decision thresholds. To obtain robust and unbiased performance estimates, K-fold cross-validation—commonly using 5 or 10 folds—is employed, while repeated cross-validation (e.g., five repeats of 10-fold CV) further reduces variance in the results. External validation using independent datasets from different populations is critical for assessing model generalizability, as performance often declines when models are applied to unseen data, highlighting the importance of diverse training samples and domain adaptation techniques. In addition, calibration assessment ensures that predicted probabilities reflect true outcome frequencies; tools such as calibration plots and the Brier score are commonly used for this purpose. For successful clinical adoption, interpretability of the Random Forest model is essential. Feature importance plots help identify key predictors influencing the model, using measures such as mean decrease impurity and permutation importance. Partial Dependence Plots (PDPs) illustrate how predictions change with variations in individual features while averaging out others, revealing patterns such as monotonic trends or threshold effects. Individual Conditional Expectation (ICE) plots further enhance interpretability by displaying prediction curves for individual instances, thereby exposing heterogeneity in feature effects across patients.

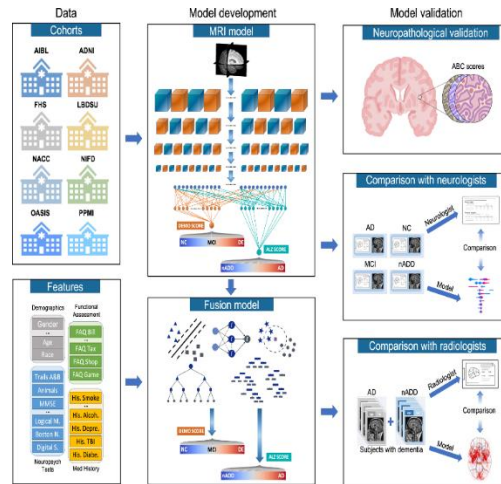


Fig 5: ML Pipeline Diagram for AD Prediction

VII. CHALLENGES AND LIMITATIONS

A.Data Heterogeneity and Missing Data

Despite significant advances, several challenges impede the clinical translation of ML models for AD prediction. Data heterogeneity across acquisition sites, scanner types, and protocols remains a major obstacle, limiting model generalizability. A model trained primarily on data from academic medical centers may perform poorly in community hospital settings or different geographic regions.

Missing data is particularly problematic in longitudinal studies due to patient dropouts, missed visits, or incomplete assessments. While Random Forest can accommodate missing data to some extent through surrogate splits, systematic missingness patterns related to disease severity or patient characteristics can introduce bias if not properly addressed.

B. Interpretability and Explainability

For clinicians to trust and adopt ML models, they must understand the reasoning behind predictions. "Black box" models, including deep neural networks and complex ensembles, provide limited insight into which features drive their decisions. Recent work has increasingly incorporated explainable AI (XAI) techniques including SHAP, LIME, and visualization methods to address this limitation.

However, interpretability methods themselves have limitations and may produce inconsistent results across different techniques. Establishing standardized protocols for model interpretation and validation is essential for clinical

acceptance. Furthermore, explanations must be communicated effectively to clinicians with varying levels of technical expertise.

VIII. CONCLUSION

Machine learning, particularly the Random Forest algorithm, has transformed early Alzheimer's disease prediction, enabling identification of at-risk individuals years before clinical diagnosis through analysis of multimodal data. This comprehensive review has synthesized findings from recent research demonstrating that optimized Random Forest models achieve accuracy rates of 86-95% for AD stage classification and 91% for early detection.

The success of Random Forest in this domain can be attributed to several factors: its ability to handle high-dimensional data with complex non-linear relationships, its robustness to overfitting through ensemble learning, its provision of intrinsic feature importance rankings, and its relative ease of implementation compared to deep learning alternatives. Recent advances in feature selection and hyperparameter optimization—particularly the combination of Backward Elimination with Artificial Ant Colony Optimization—have pushed performance to 95% accuracy while reducing computation time by 81%.

REFERENCES

- [1]. Biswas, J., Hasan, M.N., Islam, M.M., Rahman, M.M., Torabi, A., & Saha, S. (2026). Performance-optimized Alzheimer's detection using machine learning with SMOTE and randomized hyperparameter tuning. *Discover Artificial Intelligence*, 6(1).
- [2]. Bulut, N., Çakar, T., Arslan, İ., Akıncı, Z.K., & Oner, K.S. (2025). Determination of Alzheimer's Disease Stages by Artificial Learning Algorithms. *OpenAIRE*.
- [3]. Chen, L., et al. (2026). Diagnosing Alzheimer's Disease using Hypergraph Neural Networks with Prompt Tuning. *Pattern Recognition*, 113290.
- [4]. García-Gutiérrez, F., et al. (2025). Deep multimodal learning for domain-level cognitive decline prediction in Alzheimer's disease. *Frontiers in Artificial Intelligence*, 8, 1731062.
- [5]. Hancerliogullari Koksalmis, G., et al. (2025). Artificial Intelligence for Personalized Prediction of Alzheimer's Disease Progression: A Survey of Methods, Data Challenges, and Future Directions. *arXiv:2504.21189*.
- [6]. Jena, K.K., & Prasad, K. (2022). Alzheimer Disease MRI Preprocessed Images: A Machine Intelligent Based Approach for Classification and Analysis. *International Journal of Case Studies in Business, IT, and Education*, 6(2), 174-189.
- [7]. Jogeshwar, B.K., et al. (2025). Neuroanatomical-based machine learning prediction of Alzheimer's Disease across sex and age. *Neuroscience*.
- [8]. Kuhn, M., & Johnson, K. (2013). *Applied Predictive Modeling*. Springer.
- [9]. Liu, Y., et al. (2026). Utilizing multimodal models to forecast Alzheimer's disease progression and clinical subtypes. *Health Information Science and Systems*, 14(1), 10.
- [10]. 10.Malik, S., Kumari, T., Bamnawat, S., & Sonali. (2025). Generalizability, Interpretability, and Clinical Readiness of Deep Learning Methods for Alzheimer's Disease: A Systematic Literature Review. *eLife*