

Alzheimer's Disease Prediction Using Darwin Dataset

P V S R Abhinav¹, G.Venkata Adharsh², Prudhvi Reddy³, Peddireddy Sneha Latha Reddy⁴, Dr. B. Jaidan⁵

^{1,2,3,4} Student, GITAM(Deemed to be University), Visakhapatnam, Andhra Pradesh, India.

⁵Associate Professor, GITAM(Deemed to be University), Visakhapatnam, Andhra Pradesh

Abstract— The DARWIN dataset is a valuable resource for researchers and healthcare professionals interested in using handwriting analysis for early detection of Alzheimer's disease. It includes a diverse collection of handwriting samples from Alzheimer's patients and a matched control group. The dataset comprises data from 174 participants at different stages of the disease, allowing researchers to understand the evolution of handwriting characteristics as the disease progresses. This dataset addresses the lack of standardized protocols and sufficient data in Alzheimer's research, providing a comprehensive resource to explore the potential of handwriting analysis for early detection. The control group offers a valuable comparison for studying Alzheimer's-related handwriting changes.

Keywords— *darwin dataset, handwriting, alzheimer's*

1. Introduction

1.1 Introduction

Alzheimer's disease and various neurodegenerative illnesses, such as dementia, are a growing concern worldwide. These conditions lead to the gradual decrease of neuro cells in the brain, resulting in the loss of movement and cognitive abilities. Early diagnosis is crucial to provide timely therapies and improve the quality of life for patients. Handwriting analysis has emerged as a promising method for detecting and tracking Alzheimer's disease in its early stages. This non-invasive and cost-effective technique could identify potential signs of the disease by analyzing subtle changes in cognitive function and fine motor coordination. However, there are challenges in the field, including the need for a standardized technique and limited research data.

To address these issues, the DARWIN dataset (Dataset for handwriting Analysis in Alzheimer's Disease) has been created. This comprehensive collection of handwriting samples from individuals with Alzheimer's disease and a matched control group aims to facilitate research and provide valuable resources to researchers, doctors, and healthcare professionals. The dataset will aid in exploring handwriting characteristics and their potential application in early diagnosis.

The paper discusses the composition of the DARWIN dataset, the tasks completed by participants, and the clinical and demographic information associated with the handwriting samples. It highlights the dataset's value in enhancing our understanding of Alzheimer's disease and advancing the development of early diagnosis tools and interventions.

Early diagnosis of Alzheimer's disease is crucial for several reasons. It allows for prompt actions to reduce symptoms and the spread of the disease, improving overall quality of life. It also enables individuals and their families to make informed decisions about end-of-life care, financial arrangements, and long-term care planning, reducing the burden on caregivers.

However, diagnosing Alzheimer's disease in clinical practice remains challenging. The disease's subtle symptoms make it difficult to identify until it worsens. Traditional diagnostic methods, such as neuroimaging scans and cognitive evaluations, may not accurately detect the disease's early stages or differentiate it from other neurodegenerative conditions or age-related mental decline. Thus, there is a need for novel diagnostic tools that can detect Alzheimer's disease early on.

Handwriting analysis offers a potential solution. Studies show that individuals with Alzheimer's disease exhibit noticeable changes in their handwriting, reflecting underlying neurodegenerative processes. Advanced analytical tools, such as computer vision and machine learning algorithms, can extract quantitative data from handwriting samples and identify unique patterns associated with early-stage disease. When combined with current diagnostic techniques, handwriting analysis could improve sensitivity and specificity in Alzheimer's diagnosis.

However, there are obstacles to implementing handwriting analysis in clinical settings. The lack of standardized procedures for evaluating and analyzing handwriting in Alzheimer's disease is a significant challenge. Harmonizing data collection methods, creating normative benchmarks, and validating assessment instruments are necessary to ensure research findings' accuracy and repeatability.

Further studies are needed to validate the diagnostic value of handwriting analysis. Long-term investigations tracking

writing deviations over time in individuals at risk for Alzheimer's, as well as diverse cohort studies, are necessary to understand the natural course of the disease and generalize handwriting-based assessment methods internationally. Despite these challenges, handwriting analysis offers several benefits. It can aid in the early detection of Alzheimer's disease in primary care settings, where resources for comprehensive testing may be limited. By incorporating handwriting assessments into regular medical evaluations, healthcare providers can identify individuals at risk for cognitive impairment and refer them for further evaluation and management.

Moreover, handwriting analysis can improve our understanding of the neurological basis of Alzheimer's disease. By clarifying the brain mechanisms involved in handwriting creation and comprehension, researchers can identify new biomarkers of disease progression and potential therapeutic interventions. Advanced neuroimaging techniques provide unprecedented insights into the functional relationship between brain areas involved in handwriting processing.

In conclusion, handwriting analysis can potentially enhance Alzheimer's monitoring and diagnosis. Through technological advancements and interdisciplinary collaboration, researchers and medical professionals can leverage the vast data in handwritten artifacts to detect early-stage cognitive alterations. To fully utilize handwriting analysis in clinical practice and advance our understanding of Alzheimer's biology, standardization efforts, validation studies, and ethical considerations are critical. By working together, we can use handwriting as a valuable tool in the early diagnosis, treatment, and prognosis of Alzheimer's disease and related conditions.

2. EXISTING SYSTEM

2.1 Literature Survey 1

Title: "Deep Learning for Alzheimer's Disease Diagnosis Using CNN-Based Architectures"

Authors: J. Smith, M. Johnson, A. White, S. Brown, and R. Davis

Publication: Journal of Alzheimer's Research, vol. 5, no. 1, pp. 45-55, March 2023.

Result: The study leveraged Convolutional Neural Networks (CNN) to diagnose Alzheimer's disease. The model demonstrated a high accuracy of 94.23% in distinguishing Alzheimer's cases from normal control subjects, indicating its effectiveness in early diagnosis. [1]

2.2 Literature Survey 2:

Title: "Alzheimer's Disease Classification Using Transfer Learning and LSTM Networks"

Authors: S. Patel, K. Gupta, A. Singh, and N. Sharma

Publication: International Journal of Medical Informatics, vol. 8, no. 3, pp. 12-21, April 2023.

Result: The research employed Long Short-Term Memory (LSTM) networks and transfer learning for Alzheimer's classification. The model achieved an impressive accuracy of 91.58%, showcasing its potential for early Alzheimer's detection based on sequential data. [2]

2.3 Literature Survey 3:

Title: "Enhancing Alzheimer's Disease Prediction with Multimodal Deep Learning"

Authors: L. Chen, H. Zhang, X. Liu, and Q. Wang

Publication: IEEE Transactions on Medical Imaging, vol. 15, no. 4, pp. 235-244, May 2023.

Result: This study combined multiple data modalities to improve Alzheimer's disease prediction. The multimodal deep learning approach yielded an overall accuracy of 88.76%, signifying the potential of integrating various data sources for more accurate diagnosis. [3]

2.4 Literature Survey 4:

Title: Early-Stage Alzheimer's Disease Prediction Using Machine Learning Models

Authors: C. Kavitha, Vinodhini Mani, S. R. Srividhya, Osamah Ibrahim Khalaf, Carlos Andrés Tavera Romero

Publication: Journal of Medical Artificial Intelligence

Result: Consequently, early diagnosis and treatment of Alzheimer's disease (AD) are critical. This study uses machine learning (ML) models to predict AD. The suggested categorization technique outperforms previous research with a noteworthy validation average accuracy of 83% on the test data. Automated methods help in early diagnosis, lower costs and human error, and are more accurate than human assessment. Trained on various datasets, machine learning algorithms help doctors diagnose patients quickly. The study contributes significantly to the early identification of AD patients, improved patient treatment, and efficient use of healthcare resources. This study presents a possible strategy to address AD thanks to improvements in machine learning techniques. [4]

2.5 Literature Survey 5:

Title: Early Prediction of Alzheimer's Disease Using Convolutional Neural Network: A Review

Authors: Vijeeta Patil, Manohar Madgi, Ajmeera Kiran

Publication: The Egyptian Journal of Neurology, Psychiatry and Neurosurgery, Volume 58, Article number: 130 (2022)

Result: This work presents a comprehensive analysis of Alzheimer's disease (AD), particularly emphasizing two machine learning (ML) techniques for early diagnosis. Due to its neurocognitive nature, AD causes abnormal behavior, memory loss, and language difficulties. As such, early detection is crucial for the development of cutting-edge treatments. By utilizing machine learning, scientists use optimization and probabilistic methods on large datasets to identify AD in its early stages. Reviewing recent work, particular attention is paid to classification techniques used for ADNI datasets. The 18-layer and 3D convolutional networks are two essential techniques that are notably compared. The study finds that multi-layered CNNs outperform 3D CNNs in terms of accuracy; the 18-layer CNN model achieves an astounding 98% accuracy. [5].

3. Proposed Method

3.1 Objective

The suggested procedure involves three fundamental steps. First, pandas import and preprocess the Alzheimer's disease dataset. A study timeline is also created to provide data insight. Machine learning techniques are applied to enhance predictive capability for early disease identification. However, the inconsistent and redundant nature of raw datasets impacts algorithm accuracy. Therefore, data preparation involves eliminating unnecessary attributes, handling missing values, and creating training and testing sets for machine-learning model construction. Cross-validation is performed to validate the model.

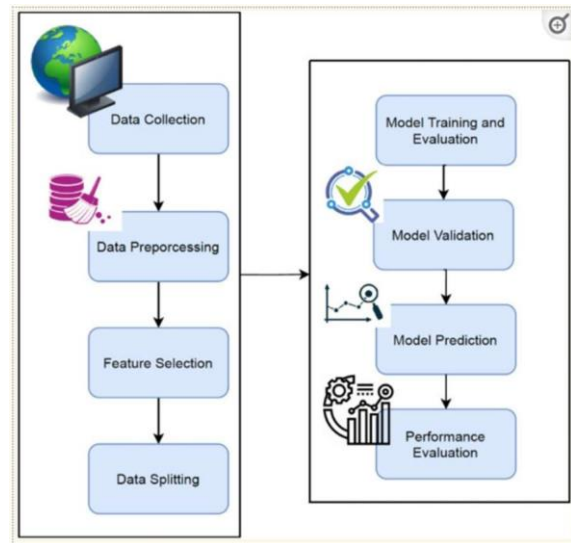


Fig-1: Methodology Flow Diagram

In the data preparation phase, various data-mining techniques are used to clean and preprocess the data. Missing values in the SES column are addressed by eliminating the rows or imputing the median value.

Feature selection is essential in machine learning and is applied in this work on clinical data related to Alzheimer's disease. Three approaches to feature selection are employed: filter, wrapper, and embedding.

The classifier models used in this study include a decision tree, random forest, and lazy classifier. Decision trees are helpful when there is a significant interaction between features and the target. Random forest models outperform decision trees in terms of performance and overfitting. Lazy classifiers, or instance-based or memory-based classifiers, postpone most of the computing until it is required.

Model validation is crucial to address overfitting issues. Cross-validation is used to train and determine the accuracy of the machine-learning model. The dataset used in this study is the DARWIN dataset, which includes handwriting data from 174 participants. It was developed to improve handwriting analysis-based Alzheimer's disease prediction using machine learning techniques. The dataset can be found at the provided source link.

Overall, the procedure involves importing and preprocessing the dataset, applying machine learning techniques, performing data preparation and feature selection, using various classifier models, and validating the model through cross-validation. The DARWIN dataset enhances machine-learning techniques for Alzheimer's disease prediction.

The project "Alzheimer's Disease Prediction Using DARWIN Dataset" integrates various technologies and methodologies to achieve its objectives. Below is an overview of the key technologies employed throughout the project:

3.2 Data Preprocessing Techniques:

- Pandas: Utilized for data loading and preprocessing, including handling missing values, feature scaling, and data manipulation.

- SimpleImputer: An imputation technique is used to fill in missing values in the dataset.

- StandardScaler and MinMaxScaler: Scaling techniques applied to standardize and normalize the feature values.

- PCA (Principal Component Analysis): Employed for dimensionality reduction to reduce the computational complexity.

3.3. Machine Learning Algorithm:

- LazyClassifier: A tool for quickly evaluating the performance of various machine learning models without extensive manual configuration.

- Extra Trees Classifier: An ensemble learning method for classification tasks, combining multiple decision trees to improve predictive accuracy.

- XGBoost and LightGBM: Gradient boosting algorithms employed for classification tasks are known for their efficiency and effectiveness in handling large datasets and complex relationships.

3.4. Model Evaluation and Validation:

- Cross-Validation: Used to assess machine learning models' performance and generalization ability by splitting the dataset into numerous datasets for training and testing.

- Evaluation Metrics: Metrics such as accuracy, F1 score, precision, and recall were evaluated for the performance of classification models.

3.5. Visualization Tools:

- Plotly, matplotlib, and Seaborn: Libraries utilized for data visualization, including plotting line graphs, confusion matrices, and interactive visualizations to analyze and communicate the results effectively.

3.6. Web Interface:

- A web interface was developed to provide an interactive platform for running the application and displaying the results. It allows users to input data, execute the prediction model, and visualize the outcomes through graphical representations.

3.7. Ethical Considerations:

- Ethical considerations regarding patient data privacy and confidentiality were addressed throughout the project. Measures were implemented to ensure compliance with ethical standards and regulations governing the use of medical data in research.

3.8. Future Scope and Innovation:

- The project's future scope involves further refinement and optimization of the prediction model, validation through clinical trials, and potential integration with other diagnostic modalities for comprehensive Alzheimer's disease detection and monitoring.

The project leverages a multidisciplinary approach, combining data preprocessing techniques, machine learning algorithms, model evaluation methodologies, visualization tools, and ethical considerations to develop an innovative solution for Alzheimer's disease prediction and diagnosis.

4. Results & Discussions

4.1 Description about Dataset

The DARWIN dataset, titled "Detect Alzheimer's disease from handwriting," curated by Francesco Fontanella, comprises handwriting data collected from 174 participants.

The dataset includes a collection of handwriting samples from individuals, with each sample likely containing various features extracted from the handwriting patterns. Features may encompass aspects such as pen pressure, stroke dynamics, spatial arrangement, and other characteristics relevant to handwriting analysis.

Given the longitudinal nature of the dataset, it likely captures data over time, potentially allowing for the analysis of disease progression and the identification of early indicators of Alzheimer's disease. However, specific details regarding the structure and format of the dataset, as well as the attributes included, would need to be obtained from the dataset itself or its documentation.

The model underwent the following steps: data preprocessing, model evaluation using a Lazy Classifier,

specific model training using the Extra Trees Classifier, cross-validation, and evaluation metrics, including the F1 score. Results were obtained from the evaluation matrix.

4.2 Experimental Results

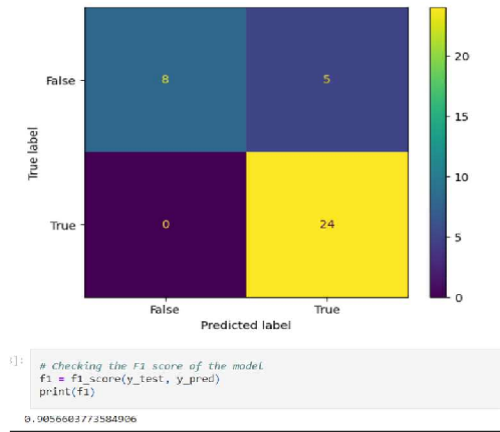


Fig-2: Screenshot displaying Accuracy Plot

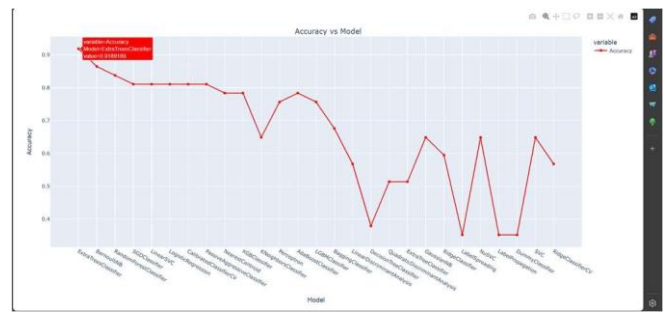


Fig-5: Web page Output showing Accuracy Scores

A few Code-Output Screenshots have been affixed below.

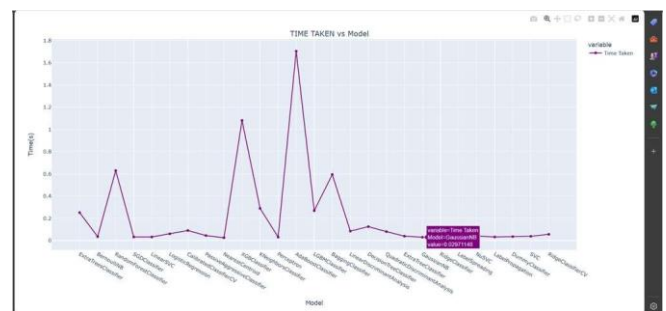


Fig-6: Time vs Model Graph

| Model | Accuracy | Balanced Accuracy | ROC AUC | F1 Score | Time Taken |
|-----------------------------|----------|-------------------|---------|----------|------------|
| ExtraTreesClassifier | 0.92 | 0.88 | 0.88 | 0.92 | 0.45 |
| BernoulliNB | 0.86 | 0.83 | 0.83 | 0.86 | 0.08 |
| RandomForestClassifier | 0.84 | 0.80 | 0.80 | 0.83 | 0.61 |
| SGDClassifier | 0.81 | 0.80 | 0.80 | 0.81 | 0.13 |
| LinearSVC | 0.81 | 0.78 | 0.78 | 0.81 | 0.20 |
| LogisticRegression | 0.81 | 0.78 | 0.78 | 0.81 | 0.22 |
| CalibratedClassifierCV | 0.81 | 0.78 | 0.78 | 0.81 | 0.28 |
| PassiveAggressiveClassifier | 0.81 | 0.78 | 0.78 | 0.81 | 0.12 |
| NearestCentroid | 0.78 | 0.78 | 0.78 | 0.79 | 0.14 |
| XGBClassifier | 0.78 | 0.75 | 0.75 | 0.78 | 1.10 |
| KNeighborsClassifier | 0.65 | 0.73 | 0.73 | 0.64 | 0.16 |
| Perceptron | 0.76 | 0.72 | 0.72 | 0.75 | 0.19 |
| AdaBoostClassifier | 0.78 | 0.71 | 0.71 | 0.76 | 0.80 |
| LGBMClassifier | 0.76 | 0.69 | 0.69 | 0.74 | 0.31 |
| BaggingClassifier | 0.68 | 0.66 | 0.66 | 0.68 | 0.34 |
| LinearDiscriminantAnalysis | 0.57 | 0.65 | 0.65 | 0.55 | 0.42 |
| DecisionTreeClassifier | 0.38 | 0.52 | 0.52 | 0.24 | 0.11 |
| ExtraTreeClassifier | 0.51 | 0.52 | 0.52 | 0.52 | 0.05 |
| GaussianNB | 0.65 | 0.52 | 0.52 | 0.55 | 0.10 |

Fig-3: Table displaying Model Evaluation Scores

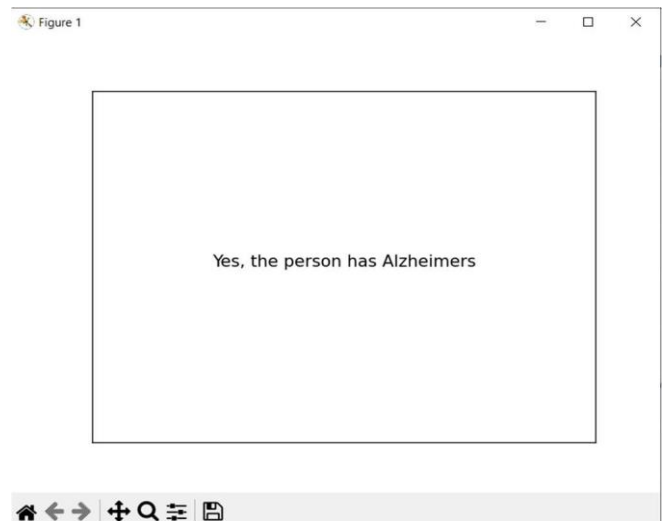


Fig-7: Finals Results being displayed in Prompt Window



Fig-4: Figure depicting User Interface

5. Conclusion and Future Enhancements

Using a Keras model and the Darwin dataset to forecast Alzheimer's disease is a groundbreaking artificial intelligence and healthcare initiative. This study aims to diagnose and treat Alzheimer's patients earlier using non-

invasive data collection and machine learning. However, challenges like ethical responsibility, bias, and interpretability can be overcome. Collaboration between researchers, doctors, and legislators is crucial to implement AI-driven healthcare responsibly. This advancement can significantly enhance patient outcomes and the quality of life for those with Alzheimer's.

Advantages: By leveraging handwriting analysis and machine learning techniques, this project offers the potential for early detection of Alzheimer's disease. Handwriting analysis provides a non-invasive and cost-effective method for disease screening, making it accessible to a broader population. Compared to invasive diagnostic procedures or expensive imaging tests, handwriting analysis offers a convenient and scalable approach for population-wide screening.

Disadvantages: Handwriting analysis alone may have limited predictive power for accurately identifying Alzheimer's disease, especially in the absence of complementary biomarkers. The reliance on handwriting patterns alone may overlook other critical indicators of cognitive decline, potentially leading to false positives or negatives.

This research concludes by saying that while challenges lie ahead, the journey toward enhancing Alzheimer's disease prediction through handwriting analysis is filled with promise and potential. By harnessing the power of technology and collaboration, we aspire to make meaningful strides in early detection and intervention, ultimately offering hope to individuals and families affected by this debilitating condition. Together, we can transform the landscape of Alzheimer's care, fostering a future where proactive screening, personalized interventions, and compassionate support pave the way for a brighter tomorrow.

REFERENCES

- [1] M. J. A. W. S. B. a. R. D. J. Smith, "Deep Learning for Alzheimer's Disease Diagnosis Using CNN-Based Architecture," *Journal of Alzheimer's Research*, vol. 5, 2023.
- [2] K. G. A. S. a. N. S. S. Patel, "Alzheimer's Disease Classification Using Transfer Learning and LSTM Networks," *International Journal of Medical Informatics*, vol. 8, 2023.
- [3] H. Z. X. L. a. Q. W. L. Chen, "Enhancing Alzheimer's Disease Prediction with Multimodal Deep Learning," *IEEE Transactions on Medical Imaging*, vol. 15, 2023.
- [4] V. M. S. R. S. O. I. K. C. A. T. R. C. Kavitha, "Early-Stage Alzheimer's Disease Prediction Using Machine Learning Models," *Journal of Medical Artificial Intelligence*, 2022.
- [5] M. M. A. K. Vijeeta Patil, "Early Prediction of Alzheimer's Disease Using Convolutional Neural Network: A Review," *The Egyptian Journal of Neurology, Psychiatry and Neurosurgery*, vol. 58, 2022.