

TRANSFORMING HEALTHCARE WITH AI POWERED

DISEASE PREDICTION BASED ON PATIENT DATA

STUDENTNAME:[K.SOPNA,J.SANDHIYA,K.SRIMATHI,K.PREETHA]

REGISTERNUMBER:[511923205043,511923205040,511923205036,511923205045]

INSTITUTION: PRIYADHARSHINI ENGINEERING COLLEGE

DEPARTMENT:B.TECH[INFORMATION TECHNOLOGY]

DATE OF SUBMISSION:[05\05\2025]

GITHUB REPOSITORY LINK:<https://github.com/sopna01/ai-prediction-in-healthcare.git>

1. PROBLEM STATEMENT

Healthcare systems are complex and challenging for all stakeholders, but artificial intelligence (AI) has transformed various fields, including healthcare, with the potential to improve patient care and quality of life. Rapid AI advancements can revolutionize healthcare by integrating it into clinical practice. Reporting AI's role in clinical practice is crucial for successful implementation by equipping healthcare providers with essential knowledge and tools.

This review article provides a comprehensive and up-to-date overview of the current state of AI in clinical practice, including its potential applications in disease diagnosis, treatment recommendations, and patient engagement. It also discusses the associated challenges, covering ethical and legal considerations and the need for human expertise. By doing so, it enhances understanding of AI's significance in healthcare and supports healthcare organizations in effectively adopting AI technologies.

The emergence of artificial intelligence (AI) in healthcare has been groundbreaking, reshaping the way we diagnose, treat and monitor patients. This technology is drastically improving healthcare research and outcomes by producing more accurate diagnoses and enabling more personalized treatments. AI in healthcare's ability to analyze vast amounts of clinical documentation quickly helps medical professionals identify disease markers and trends that would otherwise be overlooked. The potential applications of AI and healthcare are broad and far-reaching, from scanning radiological images for early detection to predicting outcomes from electronic health records.

2. PROJECT OBJECTIVES

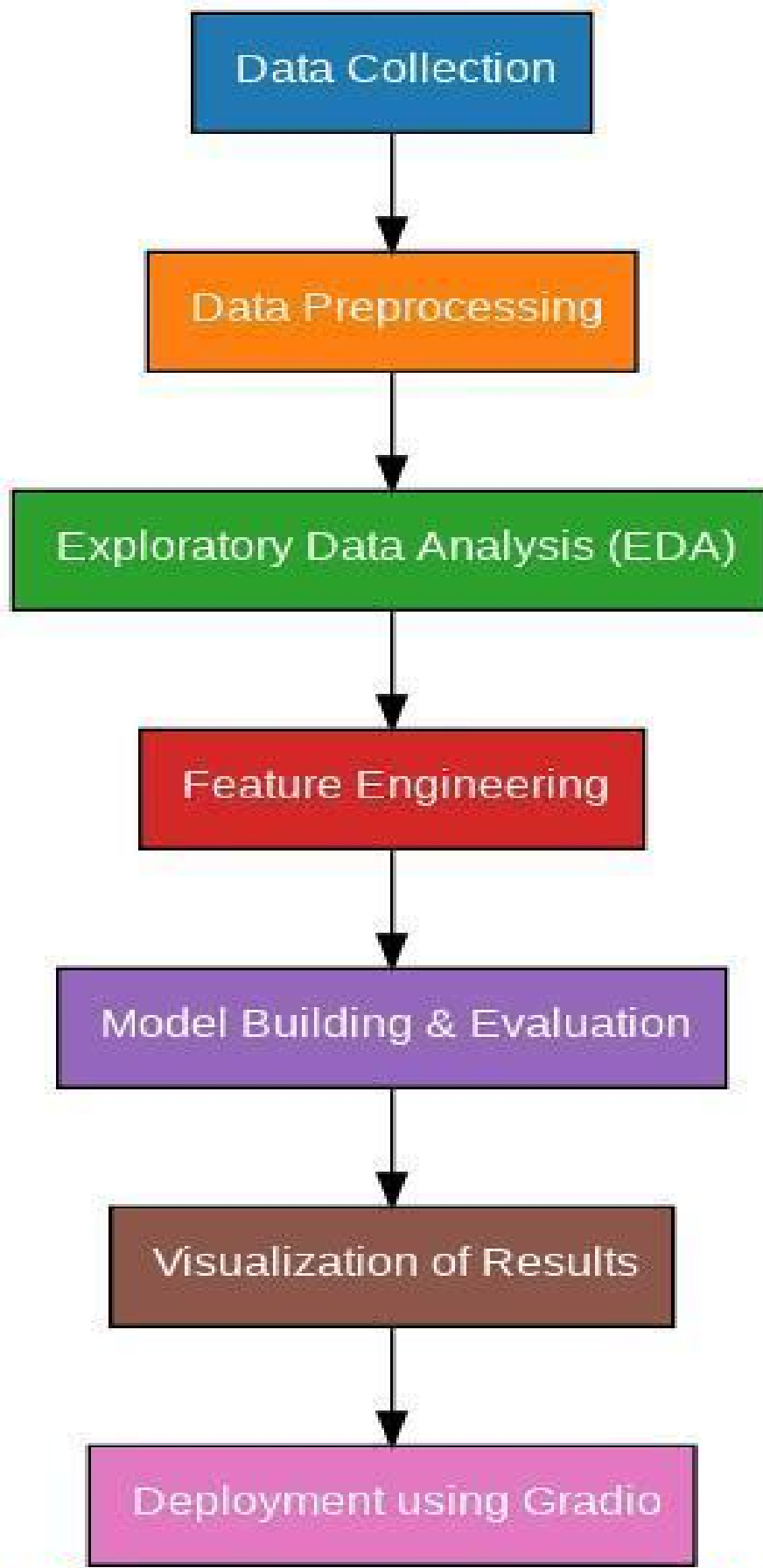
By incorporating AI into medical diagnostics, healthcare can be revolutionized, enabling the prediction and timely treatment of conditions that could otherwise lead to mortality. Early detection is particularly crucial in cases like cancer, where delayed diagnosis often results in fatal outcomes.

AI is poised to revolutionize healthcare, with the global AI healthcare market expected to **reach \$102.7 billion by 2028**. This transformation will bring about **annual cost savings of \$150 billion** by optimizing operations, minimizing diagnostic errors, and enhancing treatment efficiency. AI-powered tools have the potential to improve disease diagnosis, boosting accuracy by 20-30%, while dramatically speeding up drug discovery processes, cutting development times by 50%. Wearables and remote monitoring systems could help reduce hospital readmissions by 38%, and precision medicine is forecasted to lower treatment costs by 20%. Furthermore, AI-driven robots are anticipated to perform 30% of surgeries by 2030, leading to better patient outcomes. In underserved areas, mobile AI solutions could dramatically increase healthcare access for billions, bridging gaps and improving healthcare delivery on a global scale.

- Develop a machine learning model that accurately predicts the presence of disease.
- Identify and rank the most influential features impacting disease risk.
- provide insights into how clinical and socio-economic variables affect health outcomes.
- Ensure model interpretability and clinical usability.
- Build a user-friendly Gradio interface for testing predictions.

The primary objectives of using AI for disease prediction in healthcare is to enhance diagnostics, treatment planning, and personalized care by identifying potential health risks early, improving diagnostic accuracy, and optimizing treatment strategies. This ultimately aims to improve patient outcomes, reduce healthcare costs, and enable more proactive, preventative care.

3. FLOWCHART OF THE PROJECT WORKFLOW



4. DATA DESCRIPTION

In oncology, AI can detect small tumors or abnormal growths in imaging scans, enabling earlier treatment and better outcomes. In cardiology, AI can monitor heart rhythms and predict the likelihood of heart attacks or strokes by analyzing patterns in heart rate data.

Based on the keywords identified for search, a set of searches was performed on two of the digital databases: ScienceDirect and PubMed.

- Dataset Name: e.g., Heart Disease UCI Dataset
- Source: UCI Machine Learning Repository/Kaggle
- Type of Data: Structured tabular data
- Records and Features: (e.g., 303 records, 14 features)
- Target Variable: Disease presence (0 or 1)
- Static or Dynamic: Static dataset
- Attributes Covered: Demographics (age, sex), clinical measures (blood pressure, cholesterol, ECG results), behavioral indicators (smoking, exercise habits).

5. DATA PREPROCESSING

We define data curation as the entirety of procedures and actions after data gathering that refer to data

management, creation, modification, verification, extraction, integration, standardisation, conversion, maintenance, quality assurance, integrity, validation, traceability and reproducibility. According to this broad definition, multiple tools and applications could arguably be serving the goal of data curation. To provide a concise yet comprehensive list of curation tools, we, hence, focus on tools with broad application that cover the

most frequent use-cases of curation in medical imaging such as DICOM conversion, modification and validation.

Medical data typically contain inconsistencies and missing information. The preprocessing

step includes:

- ✓ Data Cleaning: Removal of duplicate records and filling up missing values using imputation methods.
- ✓ Feature Engineering: Feature selection of relevant features like biomarkers, symptoms, and laboratory findings.
- ✓ Normalization and Standardization: Scaling the features in order to achieve a standard data distribution.
- Verified dataset integrity: handled missing/null values.
- Removed irrelevant or constant features.
- Handled categorical variables via one-hot encoding.
- Applied StandardScaler to normalize numerical features.
- Detected and treated outliers via boxplots and z-score analysis.
- Cleanse the data to remove noise, errors, and inconsistencies. Handle missing values, outliers, and redundant features. Standardize or normalize the data to ensure uniformity across variables. Perform feature engineering to extract relevant features and enhance predictive performance.

6. EXPLORATORY DATA ANALYSIS(EDA)

Exploratory Data Analysis(EDA) is a crucial step in AI-powered disease prediction using patient data, helping to understand the data, identify patterns, and prepare it for modeling. It involves techniques like statistical summaries, visualizations, and feature engineering to uncover relationships and insights that can improve prediction accuracy.

Data Understanding and Cleaning

Data Overview:

- ✓ EDA begins with understanding the structure and nature of the patient data. This includes identifying data types, missing values, and potential errors.

Data Cleaning:

- ✓ EDA helps in identifying and addressing issues like outliers, inconsistent data entries, and inconsistencies that can negatively impact model performance.

Data Transformation:

- ✓ EDA can reveal the need to transform data (e.g., scaling, normalization) to make it more suitable for AI algorithms.

Univariate Analysis:

- focuses on examining individual variables one at a time.
- This involves calculating descriptive statistics (e.g., mean, median, standard deviation, percentiles) and visualizing data distribution (e.g., histograms, box plots).
- Univariate analysis helps identify potential outliers, data inconsistencies, and understand the range and distribution of each variable.
- Histograms for variables like age, cholesterol.
- Count plots for disease presence by gender, chest pain type.
- Bivariate/Multivariate Analysis: Correlation matrix to study clinical variable relations.
- Scatter plots of age vs cholesterol vs disease status.
- AI-powered disease prediction in healthcare, particularly utilizing exploratory data analysis (EDA), involves leveraging machine learning algorithms to analyze patient data and predict disease risk or onset. This

approach enhances early detection, enables personalized treatment plans, and optimizes resource allocation.

Key Insights:

Early Detection:

- ✓ AI and EDA can help identify patients at risk of developing diseases before symptoms appear, enabling timely intervention and potentially preventing severe outcomes.

Personalized Medicine:

- ✓ AI-driven predictive models can tailor treatment plans to individual patient needs, leading to more effective and efficient care.

Improved Diagnostic Accuracy:

- ✓ AI can assist clinicians in making more accurate diagnoses, reducing errors and improving patient outcomes.

Optimized Resource Allocation:

- ✓ Predictive analytics can help healthcare providers prioritize resources and ensure that patients at highest risk receive the care they need.

Proactive Disease Management:

- ✓ AI-driven models can predict disease progression and help clinicians anticipate potential complications, allowing for proactive management and intervention strategies.
- ✓ High cholesterol, older age, and abnormal ECGs correlated with higher disease risk.
- ✓ Males show a slightly higher incidence in some datasets.

7. FEATURE ENGINEERING

AI can enable healthcare systems to achieve their 'quadruple aim' by democratising and standardising a future of connected and AI augmented care, precision diagnostics, precision therapeutics and, ultimately, precision medicine (Table 1).³⁰ Research in the application of AI healthcare continues to accelerate rapidly, with potential use cases being demonstrated across the healthcare sector (both physical and mental health) including drug discovery, virtual clinical consultation, disease diagnosis, prognosis, medication management and health monitoring.

AI today (and in the near future)

Currently, AI systems are not reasoning engines ie cannot reason the same way as human physicians, who can draw upon 'common sense' or 'clinical intuition and experience'.¹² Instead, AI resembles a signal translator, translating patterns from datasets. AI systems today are beginning to be adopted by healthcare organisations to automate time consuming, high volume repetitive tasks.

Moreover, there is considerable progress in demonstrating the use of AI in precision diagnostics (eg diabetic retinopathy and radiotherapy planning).

- Created interaction features (e.g., cholesterol/age ratio).
- Derived binary features like "high_bp" (1 if BP > 140 mmHg).
- Removed redundant features to reduce multicollinearity.
- Encoded binary categorical features.
- Applied StandardScaler for uniform scaling.

- Utilize techniques such as correlation analysis, feature importance ranking, and dimensionality reduction to select informative features for model training. Prioritize features that are clinically relevant and contribute to predictive accuracy
- Feature engineering is the process of creating and selecting relevant features from the raw data that can be used to train predictive models.
- This involves identifying patterns and relationships within the data to create features that are informative for predicting diseases.
- Examples of features include age, gender, medical history, lab results, imaging data, and lifestyle factors.

8. MODEL BUILDING

Predictive modeling, in particular, allows forecasting the course of the disease, enabling medical professionals to avoid health risks like adverse reactions to medicine, genetically determined resistance to treatment, and failure to adhere to the regimen.

An AI model is a program that has been trained on a set of data to recognize certain patterns or make certain decisions without further human intervention. Artificial intelligence models apply different algorithms to relevant data inputs to achieve the tasks, or output, they've been programmed for.

Algorithms Used:

- Logistic Regression (for baseline)
- Random Forest Classifier (for non-linear patterns)
- Model Selection Rationale:
- Logistic Regression: interpretable baseline.
- Random Forest: robust, handles complex feature interactions.
- Train-Test Split: 80% training, 20% testing using `train_test_split`.
- Evaluation Metrics: Accuracy, Precision, Recall, F1-Score, ROC-AUC Curve
- Choose appropriate machine learning or deep learning algorithms based on the nature of the problem and the characteristics of the data. Common algorithms include logistic regression, random forests, support vector machines (SVM), gradient boosting, convolutional neural networks (CNNs), and recurrent neural networks (RNNs).

- Train the models using labeled data, optimizing hyperparameters and model architectures through cross-validation or grid search. Implement ensemble methods or transfer learning to improve model robustness and generalization.
- Fine-tune models using techniques like regularization, dropout, and batch normalization to prevent overfitting and improve performance on unseen data.

9. VISUALIZATION OF RESULTS & MODEL INSIGHTS

Feature Importance:

- Visualized Random Forest feature importances.
- Model Comparison: Compared models based on F1 and AUC scores.
- ROC Curves: Plotted ROC curves to visualize model performance.
- Residual Analysis: Analyzed false positives and false negatives.
- User Testing: Deployed Gradio app where users can input patient features and get risk predictions.

10. TOOLS AND TECHNOLOGIES USED

- Programming Language: Python 3
- Notebook Environment: Google Colab/Jupyter Notebook
- Key Libraries:
 - ❖ pandas, numpy (data handling)
 - ❖ matplotlib, seaborn, plotly (visualizations)
 - ❖ scikit-learn (modeling and preprocessing)
 - ❖ Gradio (deployment interface)

11. Team Members and Contributions

- Data Collection: [K. Sopna]
- Data Cleaning and Preprocessing: [J. Sandhiya]
- EDA and Feature Engineering: [K. Preetha]

- Model Development: [K.Srimathi]
- Deployment and Documentation: [J.Sandhiya]

