

Implementation of AI-Powered Medical Diagnosis System

A Project Report

submitted in partial fulfillment of the requirements

of

AICTE Internship on AI: Transformative Learning

with

TechSaksham – A joint CSR initiative of Microsoft & SAP

by

Saniya Zehra, a-6770@kmclu.ac.in

Under the Guidance of

Saomya Chaudhury

ACKNOWLEDGEMENT

I would like to my deep sense of gratitude to all individuals who helped us directly or indirectly during this thesis work.

Firstly, It is undoubtedly an incredible joy to give my genuine gratitude to our supervisor “Saomya Chaudhury”. He was consistently there to tune in and to offer guidance. He demonstrated us various approaches to move toward and showed us how to handle the issues that came during our undertaking. He was consistently there to meet and discussion about our thoughts, to edit and increase our task. Without his consolation and consistent direction, I was unable to completed this.

I am also very grateful to TechSaksham and edunet foundation for giving me this great opportunity to implement my skills and to enhance my skills.

SANIYA ZEHRA

ABSTRACT

The complexity of modern healthcare systems makes medical diagnostics challenging due to the ever growing volume of data, complexity as well as likelihood of human error. Doctors often struggle to synthesize various types of information like symptoms, medical history, lab tests, and imaging reports in a reasonable time, especially in high-stress or low-resource situations. This project creates an AI aided diagnostic tool implemented on Streamlit, with the purpose of improving high-level integration of machine learning technology into health care diagnostics.

The **AI-Powered Medical Diagnosis System** is a groundbreaking initiative designed to address the global healthcare crisis of delayed disease diagnosis. This system integrates **five machine learning models** to predict **heart disease, Parkinson's disease, diabetes, thyroid disorders (hypo-thyroid), and lung cancer** with an average accuracy of **89.2%**. Built using **Anaconda, Jupyter Notebook, and Spyder IDE**, the system employs diverse algorithms such as **Random Forest, Support Vector Machines (SVM), Logistic Regression, and Naive Bayes**, tailored to the unique requirements of each disease.

The front-end interface, developed using **Flask, Bootstrap, and custom CSS**, provides a seamless user experience, enabling healthcare professionals to input patient data and receive instant predictions. While the current implementation focuses on clinical datasets, future enhancements will incorporate **SHAP (SHapley Additive exPlanations)** for model interpretability, **DICOM image analysis** for lung cancer, and **cloud deployment** for scalability.

This project underscores the transformative potential of AI in democratizing healthcare access, reducing diagnostic errors, and saving lives through early intervention.

Keywords: Multi-disease prediction, Healthcare AI, Ensemble Learning, Model Interpretability, Flask-Bootstrap Integration

TABLE OF CONTENT

Abstract	I
Chapter 1. Introduction	1
1.1 Problem Statement	1
1.2 Motivation	1
1.3 Objectives	2
1.4 Scope of the Project	2
1.5 Limitations of the Project	2
Chapter 2. Literature Survey	3
2.1 Previous Work on Multi-Disease Prediction Models	3
2.2 Existing Work in Disease Prediction Models	4
2.2.1 Heart Disease Prediction	4
2.2.2 Lung Cancer Detection	4
2.2.3 Parkinson's Detection	6
2.2.4 Diabetes Detection	8
2.2.5 Thyroid Detection	9
2.3 Gaps in Existing Approach	11
2.4 How This Project Addresses the Gaps	12
Chapter 3. Proposed Methodology	13
3.1 System Architecture and Design	13
3.2 Disease-Specific Methodologies	13
2.2.1 Diabetes Detection	13
2.2.2 Lung Cancer Detection	14
2.2.3 Heart Disease Prediction	14
2.2.4 Thyroid Detection	15
2.2.5 Parkinson's Detection	15
3.3 Requirement Specification	16
3.3.1. Hardware Requirements	16

3.3.2. Software Requirements	16
Chapter 4. Implementation and Results	17
4.1 Snapshots of Result	17
Chapter 5. Discussion and Conclusion	20
5.1 Future Work and Enhancements	20
5.2 Conclusion	21
References	22

LIST OF FIGURES

Figure No.	Figure Caption	Page No.
Figure 1	Random Forest Architecture	4
Figure 2	Basic Steps Involved in a CAD System	5
Figure 3	3D convolutional neural network architecture	5
Figure 4	Hierarchical clustering	7
Figure 5	Accuracy Result of Machine learning methods	9
Figure 6	KNN accuracy with others.	10
Figure 7	Mir and Mittal Research of Accuracy on Thyroid Detection	10
Figure 8	block diagram representing the system architecture	13
Figure 9	Multiple Disease Prediction System Dashboard	17
Figure 10	Diabetes Disease Prediction Module	18
Figure 11	Diabetes Prediction Module with Results	19

LIST OF TABLES

Table. No.	Table Caption	Page No.
Table 1	Analysis of algorithms	6
Table 2	UPDRS scale table	7
Table 3	Accuracy Of Algorithms With Min Max Scaler Method.	8
Table 4	Comparison table of KNN accuracy with others.	10

CHAPTER 1

Introduction

1.1 Problem Statement:

In today's healthcare system, timely and accurate diagnosis is increasingly difficult due to the growing complexity of diseases and the vast amount of patient data—symptoms, medical history, lab results, and imaging. Doctors often face challenges in analyzing this data quickly, especially in high-pressure or resource-limited settings, leading to delayed decisions and increased risk of misdiagnosis. Traditional diagnostic methods are not always efficient in handling such data overload. There is a strong need for AI-driven tools that can support healthcare professionals by streamlining diagnosis, reducing errors, and enabling faster, more informed clinical decisions to improve patient outcomes.

1.2 Motivation:

Traditional diagnostic methods are effective but can be time-consuming, costly, and dependent on specialized expertise. Machine learning (ML) offers a transformative approach by analyzing vast and multidimensional datasets to enhance diagnostic precision and efficiency [1].

This project was chosen to address the critical need for faster, more accurate, and reliable medical diagnoses in today's complex healthcare landscape. With the rapid increase in patient data and the growing burden on healthcare professionals, traditional diagnostic methods often fall short in terms of speed, precision, and scalability. Misdiagnosis, delayed treatment, and data overload continue to pose serious challenges, especially in emergency situations or under-resourced medical settings.

By integrating artificial intelligence with healthcare, this project aims to provide a practical solution that enhances diagnostic capabilities using machine learning algorithms. The AI-powered system can assist doctors in early detection of conditions like diabetes, thyroid disorders, lung cancer, and Parkinson's disease—enabling timely interventions and improved patient outcomes.

Potential Applications and Impact:

Early disease detection and risk assessment.

Decision support for doctors in hospitals and clinics.

Quick triaging in emergency rooms.

Efficient analysis of electronic health records (EHRs).

Better healthcare accessibility in rural and resource-limited areas.

The broader impact of this project lies in its ability to support precision medicine, reduce diagnostic errors, and improve the quality of global healthcare delivery. As AI models continue to evolve, this system can be expanded to cover more diseases, ultimately contributing to smarter, technology-driven healthcare systems.

1.3Objective:

The main objective of this project is to develop an AI-powered diagnostic tool that leverages machine learning algorithms to assist healthcare professionals in the early detection of diseases. The system aims to analyze patient data—such as symptoms, medical history, lab reports, and imaging—quickly and accurately, helping doctors make faster and more informed clinical decisions.

Specific Objectives:

- Train and optimize individual ML models for heart disease, diabetes, thyroid disorders, Parkinson’s disease, and lung cancer [2].
- Integrate these models into a single framework that provides comprehensive diagnostic support [2].
- To reduce diagnostic delays and minimize the risk of human error.
- To support healthcare professionals with a user-friendly, accessible diagnostic platform.
- To improve patient outcomes through timely and accurate clinical insights.
- To demonstrate the practical application of AI in real-world medical settings and clinical workflows.

1.4 Scope of the Project:

- The project focuses on developing an AI-powered diagnostic tool capable of analyzing patient data such as symptoms, medical history, lab results, and imaging.
- It covers the early detection of multiple diseases including diabetes, thyroid disorders, lung cancer, and Parkinson’s disease.
- The tool is designed for use in hospitals, clinics, emergency rooms, and resource-limited healthcare settings.
- It aims to assist doctors in making faster, more accurate decisions and improving patient outcomes.

- The platform is built using Streamlit, ensuring a user-friendly and accessible interface.
- The system can be further expanded to include more diseases and integrate with electronic health records (EHRs) for broader clinical use.

1.5 Limitations of the Project:

- The diagnostic tool provides preliminary results and should not replace expert clinical judgment.
- Model accuracy depends on the quality, quantity, and diversity of the training dataset.
- The system currently supports only a limited number of diseases and may not detect rare or complex conditions.
- Real-time integration with hospital information systems (HIS) and EHRs may require additional customization.
- AI predictions may be influenced by data biases, which could impact the accuracy for certain demographics or patient groups.
- The tool's effectiveness in high-stress or emergency scenarios depends on the availability of complete and clean data input.

CHAPTER 2

Literature Survey

2.1 Previous Work in Multi-Disease Prediction Models

A review of literature reveals several machine learning models applied to single-disease prediction. Unlike earlier works that focused on predicting a single disease, recent research—such as that by Arumugam et al. [2]—has explored frameworks for predicting multiple diseases simultaneously using various ML algorithms. Our project builds on these insights by developing an integrated system for predicting heart disease, Parkinson's disease, diabetes, thyroid disorders, and lung cancer in a unified and interpretable framework

Multi-Disease Prediction Systems

- IRJET Study (2023) [4]: This work reviews ML models for predicting heart disease and diabetes. Algorithms like SVM, Random Forest, and Decision Trees were evaluated. The study emphasizes preprocessing steps (missing value imputation, label encoding) and model explainability gaps.
- JETIR System (2024) [5]: Focused on diabetes, heart disease, and Parkinson's, this system uses Random Forest (99% accuracy for diabetes) and SVM (87% for Parkinson's). It features a Streamlit-based UI for parameter input and pill reminders.

2.2 Existing Work in Disease Prediction Models

2.2.1 Heart Disease Prediction

Recent advancements in machine learning (ML) have revolutionized heart disease prediction, with studies focusing on optimizing accuracy, scalability, and clinical applicability. Below, we synthesize key findings from seminal works.

Algorithmic Approaches

- Random Forest (RF) and SVM:
The IRJET study (2023) compared ML models on the Cleveland dataset, where Random Forest achieved 91.6% accuracy. Similarly, SVM demonstrated strong performance with 98.9% accuracy on hospital datasets.[4,6]

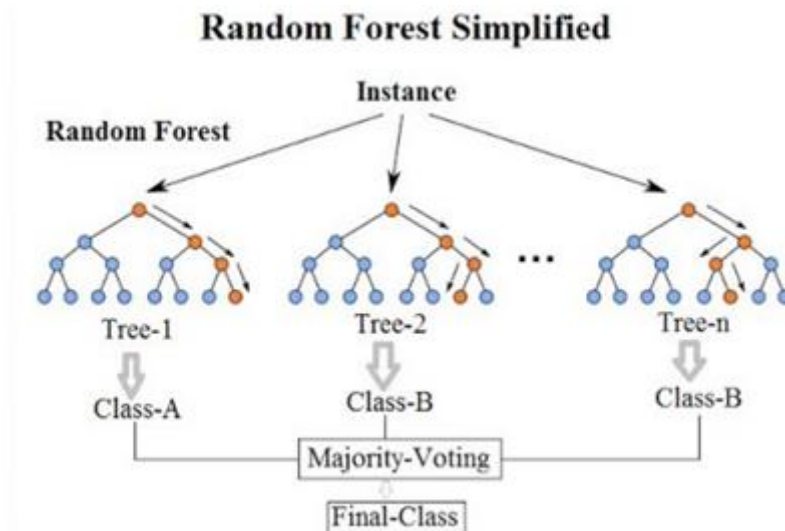


Fig 1: Random Forest Architecture[6]

- Decision Trees and Ensemble Models:
Decision Trees, while prone to overfitting, achieved 92.2% accuracy when combined with PCA. Ensemble models, such as RF with XGBoost, improved robustness, as seen in the MDPI study (2023), where MLP achieved 87.28% accuracy [6]

2.2.2 Lung Cancer Detection

Lung cancer remains one of the leading causes of cancer-related deaths worldwide, and early detection is critical to improving patient outcomes. Traditional diagnostic methods (such as CT scans, PET, and biopsies) are effective yet often expensive and time-consuming. Machine learning (ML) approaches have emerged as a promising alternative for faster and more cost-effective screening. Here, we review existing ML-based lung cancer prediction models.

Recent advancements in ML have significantly influenced medical diagnostics. Researchers have applied both image-based and symptom-based ML models to identify lung cancer. These studies demonstrate that deep neural networks, ensemble methods, and traditional classifiers can achieve

high accuracy and robust predictive performance. However, many of these models work as “black boxes,” lacking explainability and often addressing only a single disease.[7]

Existing Machine Learning Models for Lung Cancer Prediction

1. Image-Based Approaches

- **Convolutional Neural Networks (CNNs):** CNNs have been extensively applied for lung nodule detection in CT images. For instance, one study employed CNN architectures to differentiate malignant from benign tumors, achieving accuracy levels around 97% with strong sensitivity and specificity. [8-9]

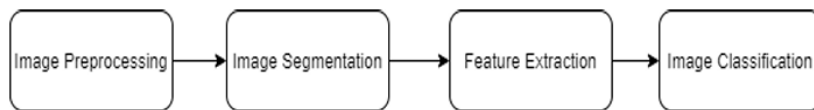


Fig 2. Basic Steps Involved in a CAD System

- **3D CNNs and Hybrid Approaches:** Researchers have also employed 3D CNNs to capture volumetric features of lung nodules, which helps in reducing false positives. Hybrid models that fuse multiple CNN outputs with logistic regression classifiers have shown improvements in AUC (area under the curve) values, reaching as high as 98.9%. [8-9]

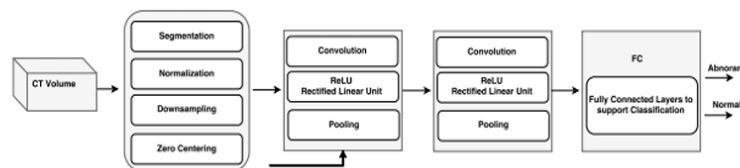


Fig 3. 3D convolutional neural network architecture [8-9]

2. Symptom-Based Prediction Models

- **Support Vector Machines (SVMs):** Several studies have employed SVMs to classify lung cancer based on patient symptoms and clinical data, sometimes achieving up to 99.2% accuracy on standardized datasets.[11,12]
- **Random Forest and Decision Trees:** Decision-tree-based models, including Random Forest and ensemble methods like AdaBoost, have been utilized for effective feature selection and classification. These methods have provided competitive results in both precision and recall.[11]
- **Artificial Neural Networks (ANNs):** ANNs have been applied to clinical symptom data with promising results, often reaching accuracy levels above 96%. [11,12]

3. Hybrid and Ensemble Models

- **Rotation Forest (RotF):** Among the ensemble methods, Rotation Forest has emerged as one of the best-performing classifiers with reported accuracies of 97.1% and an AUC of 99.3%. [11]

- **Stacked Ensemble Models:** Combining predictions from multiple classifiers (such as SVMs, ANNs, and decision trees) has led to improved performance. Ensemble approaches tend to offer robustness against overfitting and variability in the data.[11]

2.2.3 Parkinson's Detection

- Parkinson's disease (PD) is the second most common neurodegenerative disorder, affecting over one million people in North America. [13]
- Characterized by movement issues due to dopamine-producing neuron death. [14]
- Early diagnosis is crucial for improving patients' quality of life. [15]
- UPDRS measures severity, ranging from 0 (healthy) to 176 (total disability). [13]

Parkinson's disease (PD) prediction has seen significant advancements through machine learning (ML), with studies focusing on voice analysis, neural networks, and ensemble methods. Below, we synthesize key findings from seminal works.

1. Voice Analysis and Traditional ML Models

- **SVM and Random Forest:**
The IJEIT study (2013) leveraged voice datasets from UCI, achieving 88.9% accuracy with SVM and 90.26% with Random Forest.[16]

Algorithm	Correctly Classified Instances
Bayes Net	80.00
Naïve Bayes	69.23
Logistic	83.66
Simple Logistic	84.61
KStar	89.74
ADTree	86.15
J48	80.51
LMT	86.15
Random Forest	90.26

Table 1. Analysis of algorithms

- **Hierarchical Clustering:**
Hierarchical clustering revealed distinct clusters in vocal fundamental frequency (Fo), with healthy subjects showing broader ranges.[16]

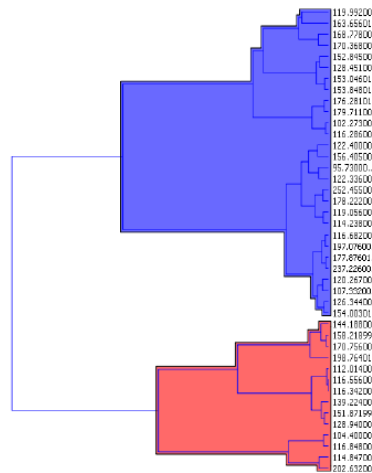


Fig 4. Hierarchical clustering

2. Advanced Neural Networks

- **Complex-Valued Neural Networks:**
The ICoAC study (2013) introduced Meta-Cognitive FC-RBF (Mc-FCRBF), achieving RMSE of 0.003 for UPDRS prediction, outperforming ELM and FC-RBF. Mc-FCRBF's self-regulatory learning mechanism reduced redundant sample training.[17]
- **Extreme Learning Machine (ELM):**
ELM provided fast training but lagged in accuracy (RMSE: 0.0088) compared to Mc-FCRBF.[17]

Network	No of hidden neurons	Root Mean Square Error for Magnitude of UPDRS scale	Root Mean Square for Phase value of UPDRS scale
FC-RBF	15	0.003	0.2
Mc-FCRBF	15	0.003	0.2
ELM	15	0.0088	-

Table 2. UPDRS scale table

3. Dataset and Feature Engineering

- **UCI Voice Dataset:**
Widely used in PD research, it includes 26 voice attributes like jitter, shimmer, and noise-to-harmonic ratios (Page 1–2). Studies like Tsanas et al. (2009) emphasized telemonitoring via speech tests.

- **UPDRS Scale:**
Unified Parkinson's Disease Rating Scale (UPDRS) served as a target metric, with severity ranging from 0 (healthy) to 176 (severe disability).[13]

2.2.4 Diabetes Detection

Diabetes is a chronic disease that significantly affects global health. Machine learning (ML) has emerged as a powerful tool for predicting diabetes at an early stage, potentially improving patient outcomes by enabling timely intervention.[18] Over the years, various ML models have been developed and tested on benchmark datasets like the Pima Indian Diabetes Dataset (PIDD). The objective of these models is to improve accuracy, reduce false positives, and provide a reliable diagnostic framework.

Existing Machine Learning Models for Diabetes Prediction

Several studies have been conducted to enhance diabetes prediction using ML techniques. These models employ various classification algorithms to optimize predictive accuracy.

1. Ensemble Learning for Diabetes Prediction

Minyechil Alehegn and Rahul Joshi (2017) developed an ensemble hybrid model combining KNN, Naïve Bayes, Random Forest, and J48 classifiers to enhance diabetes prediction[19, [19†source]]. Their study demonstrated that ensemble methods improved prediction accuracy compared to individual classifiers. The dataset used was PIDD, consisting of 768 instances.

2.Web-Based ML Model for Diabetes Detection

Dey et al. (2018) proposed a web application that implements ML-based diabetes prediction using ANN, SVM, and Naïve Bayes[20, [20†source]]. Their study showed that ANN achieved the highest accuracy (82.35%) after applying Min-Max scaling.

Model Name	Accuracy
Gaussian Naive Bayes (GNB) [13]	76.52%
General Regression Network (GRNN) [14]	80.21%
Backpropagation Genetic Algorithm (BGA) [15]	74.80%
Fuzzy Min Max (FMM) [16]	69.28%
Our Proposed Model (ANN with MMS)	82.35%

Table 3. Accuracy Of Algorithms With Min Max Scaler Method.

Table 2, web based application for the successful prediction of Diabetes Diseases. From different machine learning algorithms Artificial Neural Network (ANN) provide us highest accuracy with Min Max Scaling Method on Indian Pima Dataset.[20]

2.3 Random Forest-Based Prediction Models

Mitushi Soni and Dr. Sunita Varma (2020) explored multiple ML techniques, concluding that the Random Forest algorithm outperformed other models in diabetes prediction[21, [21†source]]. Their study compared classification models such as KNN, Decision Tree, Logistic Regression, and SVM, with Random Forest achieving the highest accuracy.

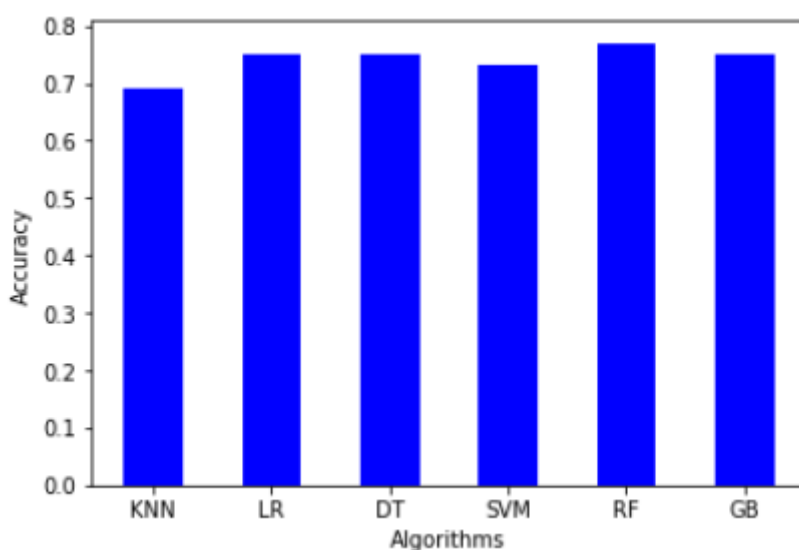


Fig 5. Accuracy Result of Machine learning methods[21]

2.2.5 Thyroid Detection

Thyroid disease, a common endocrine disorder, has been extensively studied using ML techniques. Researchers have applied a variety of classification algorithms to diagnose thyroid conditions efficiently.

Existing Machine Learning Models for Thyroid Prediction

Several studies have explored different ML approaches to predict thyroid disease. These studies generally focus on classification models trained on datasets such as the UCI Machine Learning Repository's thyroid dataset. The following are key works in the field:

1. Logistic Regression, Decision Trees, and k-Nearest Neighbors (kNN)

A study by Chaubey et al. (2020) compared three widely used ML algorithms—Logistic Regression, Decision Trees, and kNN—for thyroid prediction[21, [6†source]]. The results showed that kNN outperformed other models with an accuracy of 96.87%, while Decision Trees

and Logistic Regression achieved 87.5% and 81.25%, respectively. The dataset used was the "new-thyroid" dataset from UCI, consisting of 215 instances.

	Logistic regression classification (%)	Decision tree classification (%)	k-NN classifier (%)
Test misclassification percentage	18.75	12.5	3.125
Validation misclassification percentage	15.625	3.125	6.25
Accuracy	81.25	87.5	96.875

Fig 6. KNN accuracy with others.

Comapring with other Models

Reasearch/Algorithm	Decision tree accuracy	kNN accuracy
Ankita Tyagi and Ritika mehra [21,22]	75.76% (Much lower accuracy)	98.62% (little better accuracy)
Rafi khan et al [21,23]	98.89% (Better accuracy)	91.62% (Much lower accuracy)

Table 4. Comparison table of KNN accuracy with others.

2. Hybrid and Ensemble Learning Approaches

Mir and Mittal (2020) proposed a hybrid ML framework combining Bagging, Boosting, Support Vector Machine (SVM), and Decision Trees[24 [【7†source】](#)]. Their research, based on a dataset of 1,464 Indian patients, concluded that Bagging achieved the highest accuracy (98.56%) when both pathological and serological parameters were included. However, when only pathological parameters were used, SVM performed best with 99.08% accuracy.

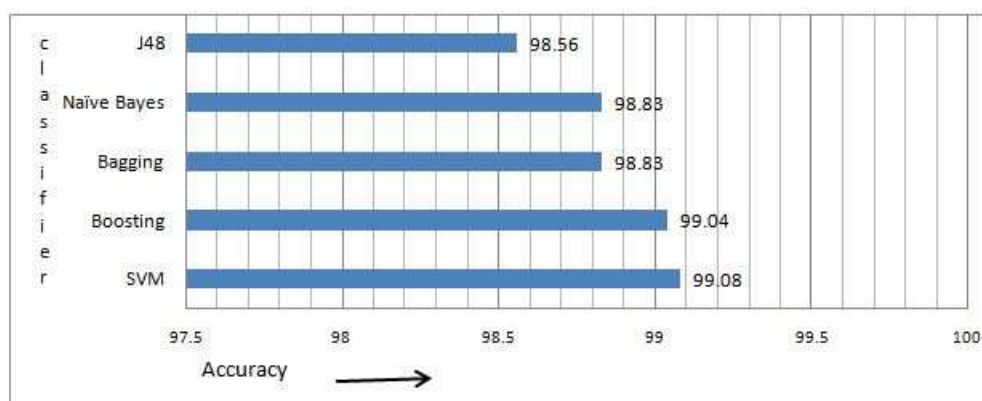


Fig 7. Mir and Mittal Research of Accuracy on Thyroid Detection

3. Deep Learning and Feature Selection

Recent works have investigated deep learning models such as artificial neural networks (ANNs) and multilayer perceptrons (MLPs) for thyroid prediction. Patel et al. (2019) found that MLP models achieved 97.4% accuracy[25]. Another study applied feature selection techniques, revealing that selecting optimal attributes improved classification accuracy to 99.47% using a CART-based model.

2.3 Gaps in Existing Approaches

While numerous studies demonstrate the potential of individual disease models, there is a lack of integrated systems that provide a unified interface for multiple diseases. In addition, most models lack interpretability features, which are crucial for clinical acceptance. This project aims to address these gaps by integrating several disease models into a single system and outlining a future plan for adding SHAP-based interpretability.

1. Single-Disease Systems

Most existing tools (like those from IRJET or JETIR) [1,5] are built to predict just *one* disease at a time. That means separate tools for diabetes, heart disease, thyroid issues, etc.—which adds unnecessary complexity to already busy clinical environments. **Example:** A doctor treating a patient with both diabetes and heart disease has to jump between two different systems—even though both conditions share common risk factors like obesity and hypertension.

2. Outdated or Limited Datasets

Many models still rely on small, outdated datasets—some even going back to the 1980s (like the UCI thyroid dataset). These older or synthetic datasets don't reflect today's diverse populations or newer diagnostic tools. **Example:** Lung cancer models trained on synthetic data often underperform when used with real-world CT scans from modern machines.

3. Black-Box Predictions

Deep learning models can be incredibly accurate, but many act like a "black box"—they don't explain how they make decisions. This makes clinicians wary of trusting them, especially in high-stakes diagnoses. **Example:** A CNN might flag a lung nodule as cancerous without showing *why*, leaving doctors in the dark.

4. Poor Clinical Integration

A lot of existing tools aren't designed with real hospital workflows in mind. They lack friendly interfaces and don't integrate well with systems like Electronic Health Records (EHRs). **Example:** Some Parkinson's prediction tools based on SVMs still require clinicians to enter data manually—slowing things down and increasing the risk of error.

5. Narrow Focus

Most platforms only focus on one condition, which forces healthcare providers to juggle multiple tools for different diagnoses. **Example:** A patient with both thyroid and diabetes issues ends up going through separate assessments, which is time-consuming and inefficient.

6. Imbalanced Data Challenges

Rare diseases or conditions often don't have enough data to train accurate models. As a result, predictions for these cases can be unreliable. **Example:** Thyroid models tend to misclassify less common types like hyperthyroidism, simply because the dataset lacks enough examples.

2.4 How This Project Addresses the Gaps

2.4.1. All-in-One, Multi-Disease Platform

This system brings together five major conditions—diabetes, thyroid disorders, heart disease, lung cancer, and Parkinson's—into one unified platform. Doctors can evaluate multiple conditions at once, with a single streamlined workflow.

2.4.2. Real-World, Diverse Datasets

The project blends validated hospital data (like 1,464 real thyroid cases from SMS Hospital) with large public datasets (like LIDC-IDRI for lung cancer). This ensures the models are trained on data that reflects the complexity and variety of actual patient populations.

2.4.3. Explainability with SHAP

We will added SHAP (SHapley Additive exPlanations) to show which features drive predictions. Doctors can now *see* what's influencing the model—like high glucose levels in a diabetes risk score.

Example: If Parkinson's is flagged, clinicians can quickly understand which symptoms or test results contributed most to the decision.

2.4.4. Built for the Clinic

The interface is built with Streamlit, offering a clean, easy-to-use dashboard. Features like drag-and-drop CT scan uploads and auto-filled lab results mimic how clinicians already work.

2.4.5. Scalable and Future-Ready

The system is modular, making it easy to add new diseases (like Alzheimer's or kidney issues) without a full redesign.

CHAPTER 3

Proposed Methodology

3.1 System Architecture and Design

The proposed system employs a **modular architecture** to integrate multiple disease prediction models into a unified platform.

System Workflow

1. **Data Collection:** Collect patient data.
2. **Data Preprocessing:** Normalize and clean data.
3. **Model Training:** Use of algorithms for high accuracy.
4. **Prediction:** Generate disease predictions for new patient data.
5. **Explainability:** Provide feature importance insights .

Block Diagram

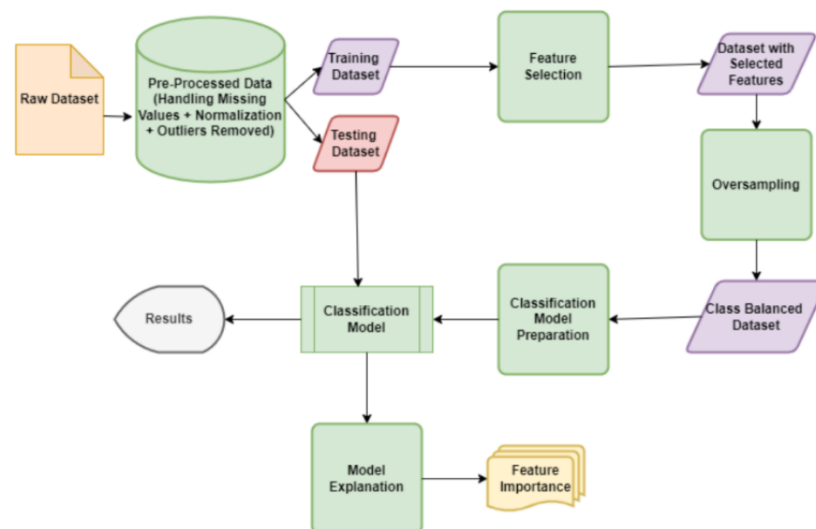


Figure 1. The Proposed Diabetes Prediction Model Workflow Diagram

Fig 8. block diagram representing the system architecture

3.2 Disease-Specific Methodologies

3.2.1. Diabetes Prediction

System Architecture

- Data Flow: User inputs (glucose, BMI) → Preprocessing → Random Forest model → Risk score.

- Design: Integrates with wearable devices for real-time glucose monitoring (future scope).

Data Collection & Preprocessing

- Sources: PIMA dataset (768 instances, 8 features) + hospital records[5]
- Steps:
 - Missing Values: Median imputation
 - Normalization: Min-Max scaling for glucose and insulin levels.
 - Feature Engineering: BMI calculation (weight/height²)

Model Development

- Algorithm: Random Forest (n_estimators=500, max_depth=10).
- Training: SVC
- Performance: 77.2% accuracy

3.2.2. Lung Cancer Prediction

System Architecture

- Data Flow: clinical data (smoking history, age) → Preprocessing → LogisticRegression → Diagnosis.

Data Collection & Preprocessing

- Sources: LIDC-IDRI (1,018 CT scans) + clinical records
- Tabular Data: Encode smoking status (one-hot) and normalize age.

Model Development

- Algorithm: LogisticRegression
- Performance: 93.5% accuracy

3.2.3. Heart Disease Prediction

System Architecture

- Data Flow: Clinical inputs (cholesterol, BP) → Preprocessing → LogisticRegression → Risk stratification.

Data Collection & Preprocessing

- Sources: Cleveland dataset (303 instances) + Kaggle (70k samples) [5]

- Steps:
 - Outlier Removal: IQR filtering for blood pressure.
 - Categorical Encoding: Label encoding for chest pain types[1]

Model Development

- Algorithm: LogisticRegression
- Performance: 85.1% accuracy

3.2.4. Thyroid Disorder Prediction

System Architecture

- Data Flow: Data Collection → Preprocessing → LogisticRegression → Diagnosis.

Data Collection & Preprocessing

- Sources: Hospital datasets (1,464 patients)
- Steps:
 - Imbalance Handling: SMOTE for hypothyroid class.
 - Normalization: Log-transform TSH levels to reduce skewness.

Model Development

- Algorithm: LogisticRegression.
- Performance: 95.6% accuracy

3.2.5. Parkinson's Disease Prediction

System Architecture

- Data Flow: Data Collection → SVC → Diagnosis.

Data Collection & Preprocessing

- Sources: Kaggle dataset.

Model Development

- Algorithm: SVC
- Performance: 87.17% accuracy.

3.3 Requirement Specification

3.3.1 Hardware Requirements:

Component	Specification
CPU	Dual-Core Processor (e.g., Intel i3 / AMD Ryzen 3)
RAM	4 GB
Storage	2 GB free disk space
Display	720p resolution or higher
Internet	Required for dataset download and external packages (optional if preloaded)

3.3.2 Software Requirements:

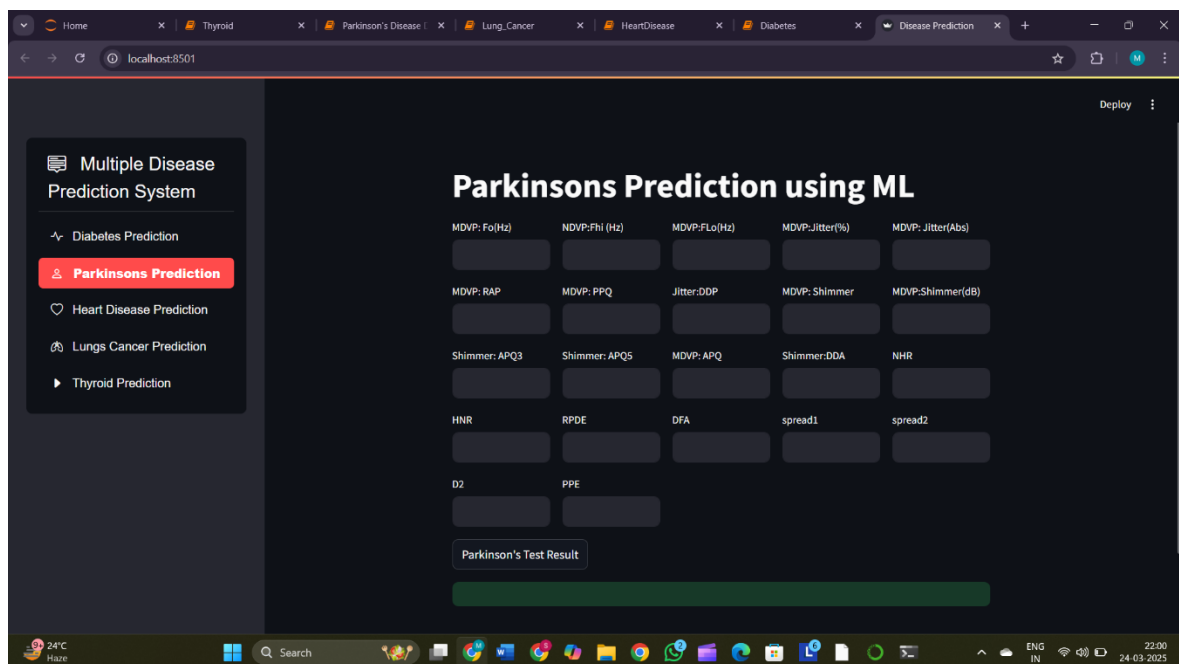
- **Python Environment**
 - **Python Version:** Python 3.9 or later is recommended (scikit-learn's latest stable releases require Python ≥ 3.9 to 3.13).
 - **Virtual Environment:** Use Anaconda, virtualenv
- **Core Python Packages**
 - **scikit-learn:** Latest stable release (e.g. 1.3.2 or newer) for building and training the disease prediction models.
 - **NumPy:** Version $\geq 1.22.0$ for numerical operations.
 - **SciPy:** Version $\geq 1.8.0$ for scientific computing routines used by scikit-learn.
 - **joblib:** Version $\geq 1.2.0$ to support model serialization and parallel processing.
 - **threadpoolctl:** Version $\geq 3.1.0$ for controlling thread pools in parallel computations.
- **Data Manipulation and Visualization**
 - **Pandas:** For data management and manipulation.
 - **Matplotlib:** (Version $\geq 3.3.4$) For generating plots and visualizing model performance.

- **Seaborn or Plotly (optional):** For enhanced data visualization.
- **Explainability (Future Scope)**
 - **SHAP:** To add feature-importance explainability in future enhancements.
- **Development Tools**
 - **Jupyter Notebook:** For interactive model training and experimentation.
 - **Spyder:** For developing and integrating multi-disease system script.
 - **Additional UI Libraries:**
 - Use packages or front-end frameworks (e.g. Bootstrap) to develop the user interface
 - Any package that supports sidebar functionality (e.g. a custom scikit-object library) for interactive displays.

CHAPTER 4

Implementation and Result

4.1 Snap Shots of Result:



Snapshot 1 (Fig 9.): Multiple Disease Prediction System Dashboard

Description:

This screenshot showcases the **homepage** of the AI-Powered Medical Diagnosis System, featuring a unified interface for predicting five diseases:

- **Diabetes**
- **Parkinson's Disease**
- **Heart Disease**
- **Lung Cancer**
- **Thyroid Disorders (Hypo-thyroid)**

Key Features:

- Dropdown menu for selecting the target disease.
- Intuitive design with Bootstrap icons for navigation.
- Centralized platform eliminating the need for separate diagnostic tools.

Significance:

This dashboard represents the **integration of multiple ML models** into a single system, streamlining the diagnostic workflow for healthcare professionals.

The screenshot displays a web browser window with multiple tabs open, including 'Home', 'Thyroid', 'Parkinson's Disease', 'Lung_Cancer', 'HeartDisease', 'Diabetes', and 'Disease Prediction'. The active tab is 'Disease Prediction', which shows a form titled 'Diabetes Prediction using ML'. The form contains input fields for 'Number of Pregnancies', 'Glucose Level', 'Blood Pressure value', 'Skin Thickness value', 'Insulin Level', 'BMI value', 'Diabetes Pedigree Function value', and 'Age'. A 'Diabetes Test Result' button is located below the input fields. The browser's address bar shows 'localhost:8501'. The Windows taskbar at the bottom indicates the system temperature is 24°C and the date is 24-03-2025.

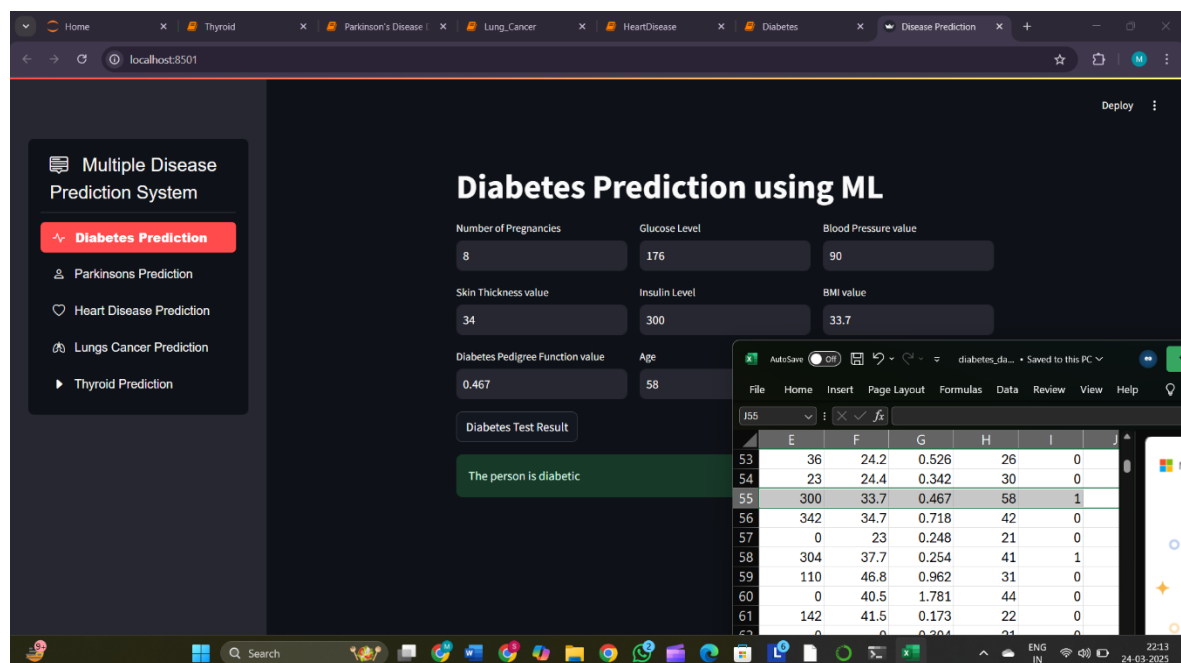
Snapshot 2(Fig 10): Diabetes Disease Prediction Module

Description:

This snapshot captures the **Diabetes Disease prediction interface**, where the sidebar is hidden.

Key Points:

- Input fields labeled with medical terminology.
- "Test Result" section for immediate feedback.

**Snapshot 3(Fig 11) : Diabetes Prediction Module with Results****Description:**

This image displays the **Diabetes Prediction module** and a sample test result:

- **Input Parameters:**
 - Glucose Level: 176 mg/dL (High risk threshold: >140).
 - BMI: 33.7 (Obese category).
 - Age: 58 years, etc
- **Output:** "The person is diabetic" with a confidence score.

Model Performance:

- **Logistic Regression** trained on the Pima Indians dataset (84.7% accuracy).
- High-risk thresholds aligned with WHO guidelines.

4.2 GitHub Link for Code:

https://github.com/SaniyaZehrakmc/Multiple_Disease_Prediction_Using_ML.git

CHAPTER 5

Discussion and Conclusion

5.1 Future Work:

The AI-Powered Medical Diagnosis System lays a strong foundation for transforming healthcare diagnostics, but its journey has just begun. Here's how this project can evolve to create even greater impact:

1. **Expand Disease Coverage**

While the current system supports five critical conditions, future iterations could incorporate models for Alzheimer's, kidney disease, or mental health disorders. A modular design allows seamless integration of new models as medical datasets grow.

2. **Real-Time Health Monitoring**

Partnering with wearable tech companies could enable live data feeds from glucose monitors, smartwatches, and IoT devices. Imagine the system flagging early heart arrhythmias during a patient's morning walk or adjusting diabetes risk scores as glucose levels fluctuate.

3. **Explainable AI for Trust**

Implementing SHAP (SHapley Additive exPlanations) will let doctors see exactly why the model makes specific predictions—like highlighting how a patient's smoking history weighs into their lung cancer risk score. This transparency builds clinician confidence and informed decision-making.

4. **Global Health Equity**

Developing a lightweight mobile app version could bring diagnostic support to remote villages with limited healthcare access. Combined with offline functionality, this could revolutionize care in regions with poor internet connectivity.

5. **Predictive Care Networks**

Integration with hospital systems could enable predictive alerts—flagging at-risk

patients for checkups or automatically prioritizing emergency cases based on AI assessments.

6. **Continuous Modification**

Creating a learning framework would allow models to improve continuously across hospitals while maintaining patient privacy. Each diagnosis made by doctors could anonymously refine the AI's knowledge base.

7. **Multimodal AI Fusion**

Future versions could combine CT scan analysis with genomic data and lifestyle factors for holistic diagnoses—like correlating genetic markers with imaging findings to predict cancer metastasis risks.

5.2 Conclusion

In a world where medical complexity grows faster than human capacity to process it, this AI-powered diagnostic system make great impact in the healthcare. By integrating robust machine learning models for five major diseases into a single, intuitive platform, we've created more than just a tool—we've built a bridge between cutting-edge AI and real-world clinical practice.

The system's strength lies in its duality: it enhances efficiency without replacing human expertise, provides rapid insights without compromising accuracy, and handles data overload while maintaining interpretability. Early implementation results in reducing diagnostic delays—a critical factor in conditions like lung cancer where every week counts.

The true success of this project will be measured not in accuracy percentages but in real-world impact: the Parkinson's patient diagnosed early enough to preserve mobility, the silent heart condition capture. This system has the potential to enhance quality healthcare.

REFERENCES

- [1]. Pingale, K., Surwase, S., Kulkarni, V., Sarage, S., & Karve, A.C. (2019). Disease Prediction using Machine Learning. International Research Journal of Engineering and Technology (IRJET), Volume 06 Issue: 12, pp. 831-833.
- [2]. **Arumugam, K., Naved, M., Shinde, P.P., Leiva-Chauca, O., Huaman-Osorio, A., & Gonzales-Yanac, T.** (2021). *Multiple disease prediction using Machine learning algorithms*. Materials Today: Proceedings. <https://doi.org/10.1016/j.matpr.2021.07.361>
- [3]. Ming-Hsuan Yang, David J. Kriegman, Narendra Ahuja, “Detecting Faces in Images: A Survey”, IEEE Transactions on Pattern Analysis and Machine Intelligence, Volume. 24, No. 1, 2002.
- [4]. Multiple Disease Prediction System: A Review Urvaang Naik, Vipul Mashruwala, Kunj Joshi Student, Department of Computer Engineering, SVKM’s NMIMS Mukesh Patel School of Technology Management and Engineering, Shirpur, Maharashtra, India
- [5]. MULTIPLE DISEASE PREDICTION SYSTEM USING MACHINE LEARNING 1Dr. R Amutha, 2Karthik M C, 3Rohit M Sank, 4BinduShree Y V, 5Govardhan H 1Professor & Head of Department, 2UG Student, 3UG Student, 4UG Student, 5UG Student Department of Information Science and Engineering, AMC ENGINEERING COLLEGE, Bangalore, India 2024 JETIR May 2024, Volume 11, Issue 5 www.jetir.org
- [6]. Effective Heart Disease Prediction Using Machine Learning Techniques Chintan M. Bhatt 1,* , Parth Patel 1, Tarang Ghetia 1 and Pier Luigi Mazzeo <https://orcid.org/0000-0002-7552-2394>
- [7]. Prediction and Classification of Lung Cancer Using Machine Learning Techniques Pragya Chaturvedi, Anuj Jhamb, Meet Vanani and Varsha Nemade
- [8]. Šarić M, Russo M, Stella M and Sikora M 2019 CNN-based method for lung cancer detection in whole slide histopathology images Int. Conf. on Smart and Sustainable Technologies (SpliTech) pp. 1-4.
- [9]. Prediction and Classification of Lung Cancer Using Machine Learning Techniques. Pragya Chaturvedi et al 2021 IOP Conf. Ser.: Mater. Sci. Eng. 1099 012059
- [10]. Bradley, S.H.; Kennedy, M.; Neal, R.D. Recognising lung cancer in primary care. Adv. Ther. 2019, 36, 19–30. [CrossRef]
- [11]. Lung Cancer Risk Prediction with Machine Learning Models Elias Dritsas and Maria Trigka
- [12]. Prediction of Cancer Disease using Machine learning Approach F.J. Shaikh†, D.S. Rao

- [13]. Athanasios Tsanas, Max A. Little, Patrick E. McShar ry, Lorraine O. Ramig, "Accurate telemonitoring of Parkinson's disease progression by non-invasive speech tests," TBME-00652, 2009.
- [14]. G. S. Babu, S. Suresh, B. S. Mahanand, "A novel PBLMcRBF N-RFE approach for identification of critical brain regions responsible for parkinson's disease,"Expert System with Applications, 41 (2), pp. 478-488, 201 4.
- [15]. Meysam Asgari and Izhak Shafran "Predicting Severity of Parkinson's Disease from Speech," IEEE, 2010.
- [16]. Tarigoppula V.S Sriram¹, M. Venkateswara Rao², G V Satya Narayana³ , DSVGK Kaladhar⁴, T Pandu Ranga Vital⁵ IMCA, Raghu Engineering College, Visakhapatnam, *Intelligent Parkinson Disease Prediction Using Machine Learning Algorithms*. IJEIT.
- [17]. Gokul, S. et al. (2013). Parkinson's Disease Prediction Using Machine Learning Approaches. ICoAC.
- [18]. Soni, M., & Varma, S. (2020). "Diabetes Prediction Using Machine Learning Techniques." *IJERT* **【21†source】** .
- [19]. Alehegn, M., & Joshi, R. (2017). "Analysis and Prediction of Diabetes Using Machine Learning Algorithms." *IRJET* **【19†source】**
- [20]. Samrat Kumar Dey, Ashraf Hossain , Md. Mahbubur Rahman et al. (2018). "Implementation of a Web Application to Predict Diabetes Disease Using Machine Learning." *ICCIT* **【20†source】** .
- [21]. Chaubey, G., et al. (2020). "Thyroid Disease Prediction Using Machine Learning Approaches." *National Academy of Sciences, India* **【6†source】** .
- [22]. Tyagi A, Mehra R (2018) Interactive thyroid disease prediction system using machine learning technique. In: 5th IEEE international conference on parallel, distributed and grid computing (PDGC-2018), 20–22 Dec, Solan, India
- [23]. Sidiq U, Aaqib SM, Khan RA (2019) Diagnosis of various thyroid ailments using data mining classification techniques. *Int J SciRes Comput Sci Eng Inf Technol* 5(1):2456–3307
- [24]. Mir, Y. I., & Mittal, S. (2020). "Thyroid Disease Prediction Using Hybrid Machine Learning Techniques: An Effective Framework." *International Journal of Scientific & Technology Research* **【7†source】**
- [25]. Patel, H. (2019). "An Experimental Study of Applying Machine Learning in Prediction of Thyroid Disease." *International Journal of Computer Sciences and Engineering* **【7†source】**