A

**Project Report**

on

# MULTIPLE DISEASE PREDICTION SYSTEM USING MACHINE LEARNING

## (Vital Care)

submitted as partial fulfillment for the award of

# BACHELOR OF TECHNOLOGY DEGREE

SESSION 2024-25

in

# COMPUTER SCIENCE AND ENGINEERING (AIML)

By

Harsh Vardhan Gupta (2100291530024)

Mayank Singhal (2100291530034)

Vaishnav Yadav (2100291530057)

Vaibhav Panjiyar (2100291530056)

**Under the supervision of**

Ms. Anjali Chauhan

# KIET Group of Institutions, Ghaziabad

Affiliated to

# Dr. A.P.J. Abdul Kalam Technical University, Lucknow

(Formerly UPTU)

**May 2025**

# DECLARATION

We hereby declare that this submission is our own work and that, to the best of our knowledge and belief, it contains no material previously published or written by another person nor material which to a substantial extent has been accepted for the award of any other degree or diploma of the university or other institute of higher learning, except where due acknowledgment has been made in the text.

Signature: Harsh Vardhan Gupta

Name: Harsh Vardhan Gupta

Roll No.: 2100291530024

Date: 20/05/2025

Signature: Mayank Singhal

Name: Mayank Singhal

Roll No.: 2100291530034

Date: 20/05/2025

Signature: Vaishnav Yadav

Name: Vaishnav Yadav

Roll No.: 2100291530057

Date: 20/05/2025

Signature: Vaibhav Panjiyar

Name: Vaibhav Panjiyar

Roll No.: 2100291530056

Date: 20/05/2025

# CERTIFICATE

This is to certify that Project Report entitled "Vital Care" which is submitted by Harsh Vardhan Gupta, Mayank Singhal, Vaishnav Yadav, Vaibhav Panjiyar in partial fulfillment of the requirement for the award of degree B. Tech. in Department of CSE(AIML) of Dr. A.P.J. Abdul Kalam Technical University, Lucknow is a record of the candidates own work carried out by them under my supervision. The matter embodied in this report is original and has not been submitted for the award of any other degree.


**Ms. Anjali Chauhan**                                    **Dr. Rekha Kashyap**

**(Assistant Professor)**                                 **(Head of Department)**

**Date: 20/05/2025**

# ACKNOWLEDGEMENT

It gives us a great sense of pleasure to present the report of the B. Tech Project undertaken during B. Tech. Final Year. We owe special debt of gratitude to Ms. Anjali Chauhan, Department of CSE(AIML), KIET, Ghaziabad, for his constant support and guidance throughout the course of our work. His sincerity, thoroughness and perseverance have been a constant source of inspiration for us. It is only his cognizant efforts that our endeavors have seen light of the day.

We also take the opportunity to acknowledge the contribution of Dr. Rekha Kashyap, Head of the Department of Computer Science & Engineering, KIET, Ghaziabad, for his full support and assistance during the development of the project. We also do not like to miss the opportunity to acknowledge the contribution of all the faculty members of the department for their kind assistance and cooperation during the development of our project.

We also do not like to miss the opportunity to acknowledge the contribution of all faculty members, especially faculty/industry person/any person, of the department for their kind assistance and cooperation during the development of our project. Last but not the least, we acknowledge our friends for their contribution in the completion of the project.

Name: Harsh Vardhan Gupta        Name: Mayank Singhal

Roll No.: 2100291530024          Roll No.: 2100291530034

Date: 20/05/2025                 Date: 20/05/2025

Name: Vaishnav Yadav             Name: Vaibhav Panjiyar

Roll No.: 2100291530057          Roll No.: 2100291530056

Date: 20/05/2025                 Date: 20/05/2025

# ABSTRACT

Machine learning-based Multiple Disease Prediction System (MDPS) proposes a sophisticated solution for predicting and weighing diagnosis of multiple diseases in a single system. This model employs potent classification Machine Learning algorithms like Logistic Regression and SVM to facilitate early diagnosis considering its benefits in health-related decision making.

It was built using the Streamlit library, and MDPS is thus equipped with a user-friendly interactive web interface that lets patients input pertinent medical parameters and converts them into real-time predictions using a streamlined interface for easy accessibility by both patients and medical practitioners.

What this research is dedicated to is Predicting major health ailments-Diabetes, Heart Disease, and Parkinson's Disease. Out of all major diseases, these three were picked up primarily because of their prevalence in the public health field and availability of structured datasets suitable for training learning models. The basic health parameters usually processed by the system include blood pressure and pulse rate, cholesterol levels, and heart rate-all of which are widely cited risk factors for one or more of the target diseases.

MDPS trains the model based on preprocessed sample datasets and learns to discern patterns and correlations that help it predict the occurrence of diseases with a tolerably high degree of accuracy.

Current trends in artificial intelligence and machine learning have become important in diagnosing medicine, especially requiring a concurrent evaluation of multiple health risks. Most existing models deal with a single disease, but it offers a multi-disease prediction framework, filling the gap in most modern health technologies. It centralizes different models into one interface enhancing efficiency while giving users an integrated tool for proactive health monitoring.

The experiment results have specific accuracy ranges that the system can achieve for each ailment, showcasing the potential applicability of the system. Key challenges encountered during the study are also discussed, which include data imbalance, feature selection, and the importance of data preprocessing, all having a considerable impact on model performance.

The MDPS Project illustrates that machine learning can improve predictive healthcare. Such future developments may include real-time health monitoring, linking with wearable devices, and including more diseases in the future. This illustrates the value of an interdisciplinary

approach toward enhancing technology to improve personalized medicine and public health outcomes. **Keywords**: Machine Learning, Multiple Disease Prediction System (MDPS), Logistic Regression, Support Vector Machine (SVM), Diabetes Prediction, Heart Disease, Parkinson Disease Detection, Data Preprocessing, Streamlit, Real-Time Health Monitoring.

# TABLE OF CONTENTS
**Page No.**

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

- **MDPS** - Multiple Disease Prediction System
- **ML** - Machine Learning
- **SVM** - Support Vector Machine
- **PIDD** - PIMA Indian Diabetes Dataset
- **ANN** - Artificial Neural Network
- **KNN** - K-Nearest Neighbours
- **CNN** - Convolutional Neural Network
- **LSTM** - Long Short-Term Memory
- **EHR** - Electronic Health Record
- **AI** - Artificial Intelligence
- **RF** - Random Forest
- **DT** - Decision Tree
- **NB** - Naïve Bayes
- **LR** - Logistic Regression
- **MLP** - Multi-Layer Perceptron
- **ELM** - Extreme Learning Machine
- **ROC** - Receiver Operating Characteristic
- **AUC** - Area Under Curve
- **TP** - True Positive
- **TN** - True Negative
- **FP** - False Positive
- **FN** - False Negative
- **WHO** - World Health Organization

- **IoT** - Internet of Things
- **HIS** - Health Information System
- **PCA** - Principal Component Analysis
- **FPR** - False Positive Rate
- **PPV** - Positive Predictive Value
- **NPV** - Negative Predictive Value
- **BMI** - Body Mass Index

# CHAPTER 1
# INTRODUCTION

## 1.1   INTRODUCTION

Machine learning, which is the same as the ML acronym, has seen tremendous growth in the last few years and it is a game changer in every industry, especially in healthcare. The ML algorithms, fused with medical diagnostics, can be an engine to improve health level detection and also provide more efficient, individualized, and faster service delivery.

One important area of advancement that can be undertaken is the conception of systems that can simultaneously predict many diseases in a patient-in-the-know way, which may redefine the manner in which healthcare professionals undertake diagnosis, treatment, and monitoring of patients.

Healthcare data diversity is often complex as it encompasses different types of patient-related information such as clinical test results, electronic health records (EHRs), medical imaging, lifestyle factors, and genetic profiles. Most of the predictive models in medical science were disease-wise models based on single diseases, like diabetes, heart disease, or Parkinson's disease, neglecting many common co-morbidities as well as interconnections between several diseases.

With the increasing awareness of the need to understand the patient's health rather than any single illness, a multi-disease prediction platform is highly

considered now. Presently, researchers lack a unified and common framework that ensures multiple disease predictions on one single platform. The gap thus calls for comprehensive multi-disease models that can provide accurate real-time actionable insights.

Such models will go a long way in improving the quality of both diagnostics and patient care through early warnings and critical information about the symptoms. Such model developments necessitate stringent machine learning algorithms for disentangling convoluted nonlinear relations within the healthcare data for the precision and promptness of results.

This research centers on the three major diseases namely heart disease, diabetes, and Parkinson's disease, which are all interrelated. All these diseases share some risk factors and affect one another, and thus make them suitable for applying the multi-disease prediction approach. This research aims to utilize machine learning to develop a system that will assess the probable risk of the mentioned conditions concurrently based on the specific data of each patient involved.

In achieving that, we have leveraged Support Vector Machines (SVM), which is a prominent and superlative type of machine learning algorithm for classification tasks and especially so in biomedical data analysis. The power of SVM lies in its ability to deal with both linear and nonlinear relationships between the various features (symptoms, medical test results, and demographic data) that are fed into the model and the target outcomes (i.e., presence or absence of disease).

The algorithm has a commendable characteristic of forming optimal hyperplanes that effectively segregate the classes, which also makes it very robust to overfitting, a key property for healthcare applications using high-dimensional data.

Apart from SVM, this study also takes into account other related machine learning algorithms such as Naïve Bayes, K-Nearest Neighbors (KNN), Decision Trees, and Random Forests. Each of these has its strength - Naïve Bayes allows for probabilistic reasoning, KNN is characterized by simplicity and ease of interpretation, Decision trees yield very straightforward, rule-based decision-making, and Random forest minimizes the occurrence of overfitting through an ensemble learning scheme.

Streamlit is an open-source Python library that is designed to facilitate building a speedy interactive user-friendly web application for practical implementation. Users would effortlessly input their health data into the system developed and receive instantaneous predictions by the underlying algorithms performing the machine learning.

This democratizes the access to advanced diagnostic facilities and also caters for real-time decision making by health professionals and patients alike. The research can, therefore, be said to have been based on two key specific objectives: to craft an efficient multi-disease prediction framework that employs SVM and other machine learning techniques regarding identifying the risk of heart disease, diabetes, and Parkinson's disease; and to evaluate the reliability and effectiveness of such a system in clinically closely

resembling environments.

Through this work, we hope to advance the area of AI-driven healthcare by creating interventions to address preventable care, misdiagnosing into diagnostic errors, improving patient outcomes, or all three. The future looks like it is going to be very bright, with machine learning as a backbone to a more efficient and effective health system, all harmonized in a single environment for multii-disease predictions.

# 1.2 PROJECT DESCRIPTION

The project is named "Multiple Disease Detection System Using Machine Learning," designed for creating online framework prediction for multiple diseases at once. The primary focus of this project will be diabetes, heart diseases, and Parkinson's disease problems.

By early diagnosis based on confirmed accurate diagnosis, the approach significantly addresses contemporary situations such as the growing global burden imposed by chronic illness on resources and lives.

The system employs the power of advanced machine learning methods such as Logistic Regression and Support Vector Machines (SVM) in making reliable data-count predictions based on the individual's medical information. These algorithms are chosen as the best entity for their proven effectiveness in handling complex and often nonlinear element patterns for which specialized algorithms are created to spot the often subtle early signs of these critical conditions.

The project's key concept advocates for the empowerment of healthcare professionals and patients through cost-effective early diagnosis and delivery with an intuitive user interface. The user can input some critical health parameters such as age, blood pressure, glucose level, cholesterol, and other key vital signs into the platform.

The event collects this data and interprets the probability of disease occurrence as an overall risk assessment, not limited to conventional diagnostic medians that usually rely on single disease predictive efforts but covering health in holistic terms.

To ensure that the application is user-friendly, accessible and also a web-based application is created with the use of Streamlit, which is a common Python framework for efficient and speedy development of interactive data applications. This enables real-time prediction from any internet-connected device attached to a broad and diverse user base from patients and clinicians to researchers and healthcare providers.

The basic methods of data mining involved in this study are data collection, preprocessing, and model training. The datasets, collected from verified public repositories, are subjected to preprocessing processes such as data cleaning, normalization, and missing value handling to improve input quality significantly. It trains models using techniques such as cross-validation, an essential technique in avoiding overfitting.

Logistic Regression is adopted as an easy yet effective method of binary classification to predict diabetes, while SVM, due to its notable ability to capture nonlinear and complex decision boundaries, is preferred in heart disease as well as in Parkinson. Fine-tuning the models creates the foundation of using a web application to convert predictions into real time, according to user health risk assessment.

Initial evaluation results indicate promising performance, with accuracy rates reaching approximately 78% for diabetes, 85% for heart disease, and an impressive 89% for Parkinson's disease. These findings impressively demonstrate the potential of this approach in assisting real-time and accurately detecting diseases, contributing significantly toward more personalized, data-driven healthcare solutions.

In future directions, the project aims to include the more advanced deep learning architectures that will be able to capture very complex hierarchical patterns within the medical data along with real-time observations coming from health wearables that will continuously monitor users. These will add up even more to the advancement of the system so as to make it more personalized in terms of health risk profiling of individual users.

In general terms, this project demonstrates a marked advancement in artificial intelligence application in health, hence creating a cost-effective and scalable solution for diagnosing and managing diseases early on, ultimately improving healthy outcomes across populations worldwide.

# CHAPTER 2

# LITERATURE REVIEW

## 2.1 Diabetes Prediction

Diabetes remains an ever-growing menace, affecting people in millions across nations globally. It does burden the healthcare system profoundly and, equally importantly, the quality of life of those diagnosed. Thus, with a greater emphasis on modeling these predictions, it is realized that early diagnosis leads to appropriate treatment, better management, and lower complications.

Machine learning (ML) has emerged as a seriously better alternative in this area, offering potential techniques that analyze complicated medical data, reveal hidden patterns, and afford some reliable predictions. One of the frequently used datasets for diabetes research has become the PIMA Indian Diabetes Dataset (PIDD).

Having a full spectrum of diagnostic measures like glucose concentration, blood pressure, body mass index, and age, it has become the benchmark for other ML models owing to its wide application to various scenarios. Conventional methods, like Logistic regression and Decision trees, have held ground for long owing to their simplicity and their interpretative ability.

Logistic regression serves as an obvious binary classifier of the kind that helps by indicating how factors can affect the chance of having diabetes. This property makes it especially convenient in a clinical setting, where the impact of any parameter is just as critical.

Decision trees articulate the route taken to arrive at an outcome, and therefore healthcare professionals can easily trace the decision steps taken to achieve a certain result. Interestingly enough, while interpretable, they have often failed to describe some of those real medical data areas with their complexities and nonlinearities.

More advanced models have thus gained acceptance, pushing back the simpler ones, with the SVM and Random Forest leading these more modern approaches. SVMs prove to be very suitable to search for that hyperplane (boundary) that optimally separates diabetic from non-diabetic acts, hence often achieving very high accuracy in sharply-defined datasets.

When it comes to high-dimensional data and nonlinear relations, SVMs have proved to be very effective. The other approach employed by the Random Forests entails the construction of many decision trees during training time and the aggregation of their predictions.

By mixing the predictions of numerous decision trees instead of relying on one tree, this ensemble method reduces the risk of overfitting, increases robustness, and adds extra power to prediction accuracy. Recently, the excitement has been directed towards Deep Learning with ANNs taking the lead. ANNs differ from conventional models as they can learn complex hierarchical patterns in data independently, requiring no feature engineering.

The potential of ANNs is evident in diabetes prediction, where subtle interactions across multiple variables might drastically affect the outcome. In this regard, predictions would stand to further improve with the continued increase in the availability of heterogeneous data sources such as continuous glucose monitoring

and imaging data.

With this historical trend, we see the evolution of diabetes prediction research from elementary machine learning models to advanced ensemble and deep learning approaches, thus to better and more precise early detection and patient outcomes.

## 2.2 Heart Disease Prediction

Heart diseases remain globally among the leading causes of death, thus posing formidable problems to public health. Early prediction and accurate prognosis of heart disease are necessary for prevention, intervention, and better patient outcome.

This triggered a large inquiry into predictive modeling based on clinical and lifestyle-related risk factors such as hypertension, cholesterol level, age, smoking status, and physical activity. Recently, ML-based techniques have captured much attention owing to their potential of analyzing complex medical datasets and improving the diagnosis and risk assessment of heart diseases.

Naïve Bayes and K-Nearest Neighbors (KNN) classifiers are the most widely applied ML methods for predicting heart disease. Simple algorithms that need less computational power and yet are effective enough to handle medical datasets are preferable. The Naïve Bayes classifier assumes conditional independence between the features for the straightforward computation of probabilities, a property that grows in importance with large datasets comprising huge numbers of attributes.

Notwithstanding this simplifying assumption, the Naïve Bayes classifier has

proved effective in various healthcare classification problems, including heart disease detection. The probabilities it provides can be of great help in the clinical decision-making process.

On the other hand, K-Nearest Neighbors is a straightforward, intuitive algorithm that classifies patients depending on the degree of similarity of their clinical features to those of previously diagnosed cases. By measuring distances between feature vectors, KNN assigns the class of majority among the nearest neighbors.

This, however, works acceptably well, especially in case the data in consideration has clustered nature. Nevertheless, the choice of distance metric and number of neighbors considered can impact KNN's performance; computationally costly with respect to time as the size of the dataset grows.

Apart from the traditional methods, ANNs have emerged for heart disease prediction as very important tools due to their power to model complex, nonlinear relationships never separated from the clinical data. ANNs consist of interconnected layers of nodes/neuron structures, weighting their interconnections in processing input features, which gives the network possibility to learn complex patterns between risk factors of an interest and the disease outcome.

They capture usually subtle nonlinear dependencies, which contribute to better prediction accuracy compared to simpler models. They have exceptional generalization ability on unseen data, thereby enhancing their importance to clinical settings due to variability common in patient profiles.

More recently, hybrid models, which combine machine learning with feature selection, have gained more attention. Feature selection is employed to find the

relevant and most informative clinical variables, for example, some cholesterol subtypes of blood pressure readings that contribute significantly to heart disease risk.

By filtering out redundant or less informative features, those contribute to the reduction of noise, increase computational efficiency, and amplify the predictive power of ML models. To offer an example, SVM integrated with feature selection algorithms proved superior in performance, focusing on important risk factors while throwing away unimportant ones.

Such integration and ensemble methods improve the accuracy of predictions and enhance the interpretability and reliability of heart disease diagnostic systems. The coagulation of sophisticated ML algorithms and domain-related feature selection best reasons such systems into truly robust decision-making aids for clinicians with respect to the early detection and medical intervention appropriate for individualized treatment planning based on a risk profile of a particular patient.

In heart disease prediction, the application of various machine learning techniques can be largely characterized as applaudable emergence in the field of healthcare analytics. These models support the healthcare professional with accurate risk predictions and diagnostic suggestions aimed at better patient management or preventive care strategies.

Further, as data, availability and computational power, continue to grow, machine-learning-based heart disease predictive systems are poised to become essential to clinical settings.

## 2.3 Parkinson's Disease Prediction

Parkinson's disease (PD) is a neurodegenerative disorder that eventually affects motor function, including movement, maintaining balance, coordination, and muscle control. The condition slowly progresses with time, manifests as tremors, stiffness, slow movement, and postural instability, whereby it becomes difficult to manage symptoms effectively.

Therefore, accurate, early diagnosis of Parkinson's disease is important for managing symptoms, slowing disease progression, and adding to better quality of life for patients. It is difficult to diagnose early Parkinson's disease due to such problems as apparent overlapping of its symptoms and the insidious manner in which early symptoms manifest.

Very recently, the technique of machine learning (ML) has evolved as a strong candidate for the early detection and diagnosis of Parkinson's disease. By analyzing diverse patient data from voice and gait to tremor measurements, these ML models can uncover subtle markers that may not be detectable through conventional clinical assessment.

Altering the early with objective and accurate diagnostic assistance, these computational approaches seem to hold great promise in aiding healthcare professionals. Among the many algorithms in machine learning that have been employed to detect the presence of the disease, the support vector machines (SVM) algorithms have performed particularly well.

SVM algorithms find the optimal hyperplane or boundary that best separates individuals with Parkinson's disease from the healthy based on the computed feature representation. Importantly, this algorithm is efficient when large dimensionality is involved and the separating surfaces are very complex. The

SVM has been shown to differentiate between Parkinson's patients and controls accurately when used together with sophisticated feature extraction methods.

The main features were the acoustic measures from the voice recordings, such as frequency modulation, amplitude variation, and jitter, and motor measurements from tremor intensity and gait analysis. These features assist SVMs in capturing the subtle physiological changes that accompany early Parkinson's.

Besides SVMs, others that have shown widely in classification of Parkinson's disease are tree-based algorithms like Decision Trees and Random Forest. Decision Trees, while easy to interpret, can also give clinicians an insight into which features were most responsible for making the prediction.

Random Forests serve to improve the prediction accuracy by averaging over a set of decision trees, reducing the overfitting risks and improving generalization ability for novel situations from varying patient data. That robustness is one of the reasons Random Forest models are more appropriate for implementation into clinics where patient data can vary.

In addition to the traditional ML techniques, more recently, advances in deep learning have begun to reshape Parkinsons' research. Neural network architectures such as Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks are being used in the analyses of such complex high-dimensional datasets as time series of speech and motor activity recordings.

CNNs are adept at extracting spatial features from spectrograms of voice signals or imaging data, detecting complex patterns that may indicate disease progression. LSTMs, or Long Short-Term Memory networks, are Recurrent Neural Networks capable of modeling sequential data and time dependencies,

thus particularly suited to capture time-dependent changes in motor symptoms.

Thus, integration of these deep learning techniques presents hope for precise early detection of Parkinson's disease, allowing, hopefully, the clinicians to start treatment even earlier and better tailor their intervention according to individual symptom course.

Moreover, those ones promise to be integrated into some wearable technology and mobile health applications that will allow the remote continuous monitoring and real-time assessment of the patient outside clinical settings. In summary, machine-learning techniques—from classical SVM and Random Forest algorithms to advanced deep learning attempts—provide great hope for improving the diagnosis and management of Parkinson's disease.

It can use all kinds of patient data and advanced pattern recognition to give the patient an earlier diagnosis, support a tailored treatment regimen, and return a better life.

# CHAPTER 3

# PROPOSED METHODOLOGY

Indeed, predicting more than one disease at once is bound to be a complicated task that calls for a systematic and comprehensive approach. The multi-disease prediction system has been developed in such a way as to provide an end-to-end approach that guarantees accuracy, reliability, and usability, with the structured sequence of steps. The phases of the project include:

1. **Data Collection:** High-quality data constitutes the foundation of any machine-learning system. In terms of this project, data collection includes extensive medical information about diabetes, heart disease, and Parkinson's disease. Various channels, such as electronic health records (EHRs), publicly available medical research databases, and health organizations' datasets, provide good data sources for the study.

   This multi-source strategy helps in building a rich representative collection that reflects different patient demographics, symptoms, and clinical results. Special care is taken to obtain those datasets that have clinical parameters such as blood sugar levels for diabetes, cholesterol, blood pressure readings for heart disease, and motor function data for those with Parkinson's disease in order to provide enough samples for the models to input features needed for reliable predictions.

2. **Data Preprocessing:** Raw medical data is mostly characterized by noise, incompleteness, and inconsistency, leading to adverse effects on model performances. As a result, preprocessing is a very integral way of performing two main tasks: data cleaning and transformation. During cleaning, missing

values are addressed through imputation techniques, outlier detection and handling, and removal of irrelevant or duplicate entries. Following the cleansing process, the data are transformed into acceptable formats applicable in machine learning algorithms. Such includes normalizing numerical values to a common scale, encoding categorical variables, and feature selection or engineering processes that are beneficial for model learning. Preprocessing guarantees structured, balanced, and bias-free datasets avoiding prediction skewing.

3. **Model Selection:** Selection of the appropriate machine learning algorithms becomes pivotal for constructing the predictive models accurately. The algorithms were selected according to each dataset on account of the possession of some different characteristics with respect to diabetes, heart disease, and Parkinson's disease. Classification algorithms to be analyzed include SVM, logistic regression, decision trees, random forests, and naïve Bayes, which discriminate diseases according to their strengths and suitability.

Each model is then trained with the preprocessed data to learn the patterns and relationships between the input features and the respective disease outcomes. For evaluating and comparing models, performance metrics include accuracy, precision, recall, and F1-score. Accordingly, hyperparameter tuning and cross-validation techniques optimize model parameters for improved generalization and overfitting reduction.

4. **Data Splitting:** Then, there is a need to maintain the models' robustness vis-a-vis prediction validity. This would mean that there will be a dataset partitioned into training and testing subsets, where training is generally allotted the higher percentage of the data (like 70 to 80%); the rest serves as

a testing set. The training set is used to teach the models how to classify and predict diseases based on input features.
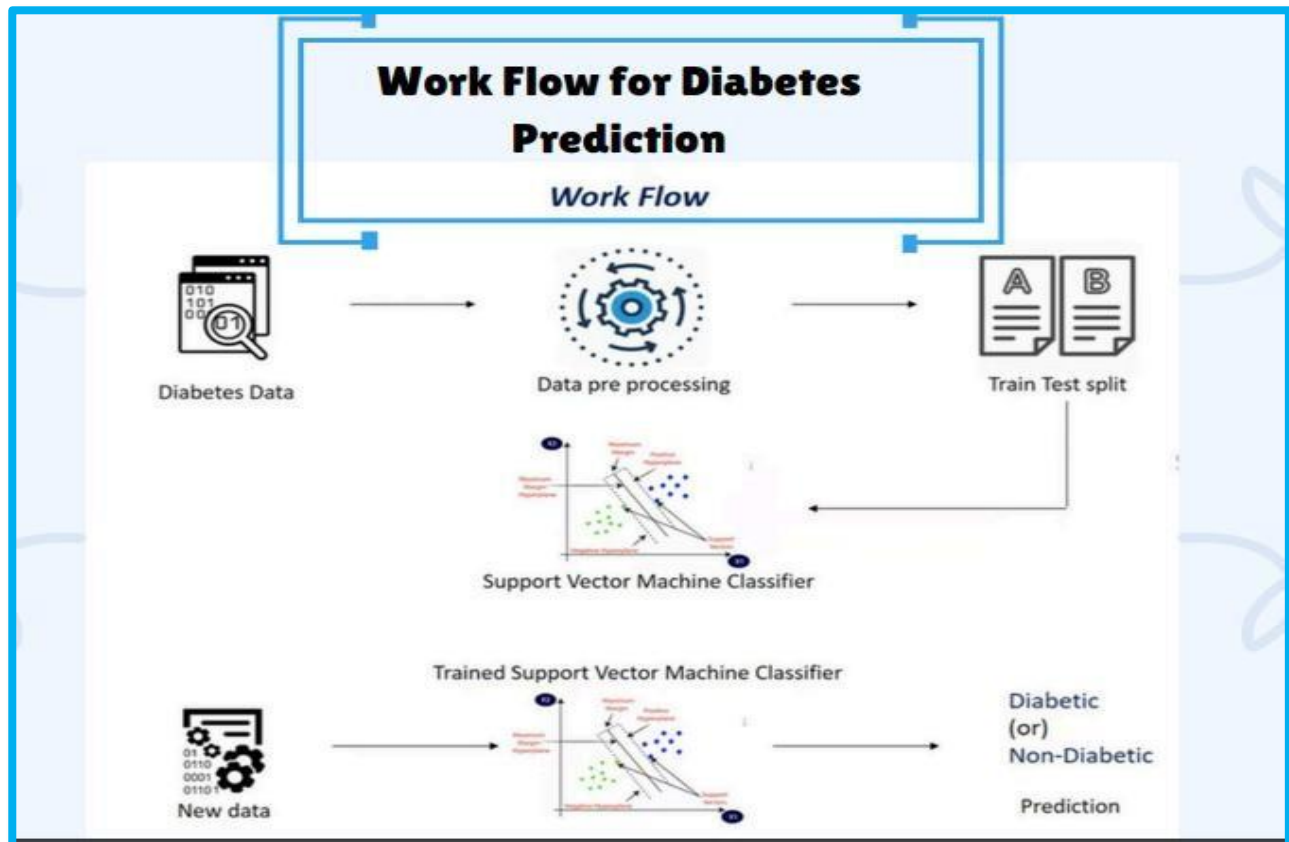
The unseen testing set assesses how well the trained model performs on new real-world data, providing an unbiased evaluation of accuracy and reliability. This is a necessary step that identifies any potential weaknesses in the models and ensures it can generalize beyond training data.
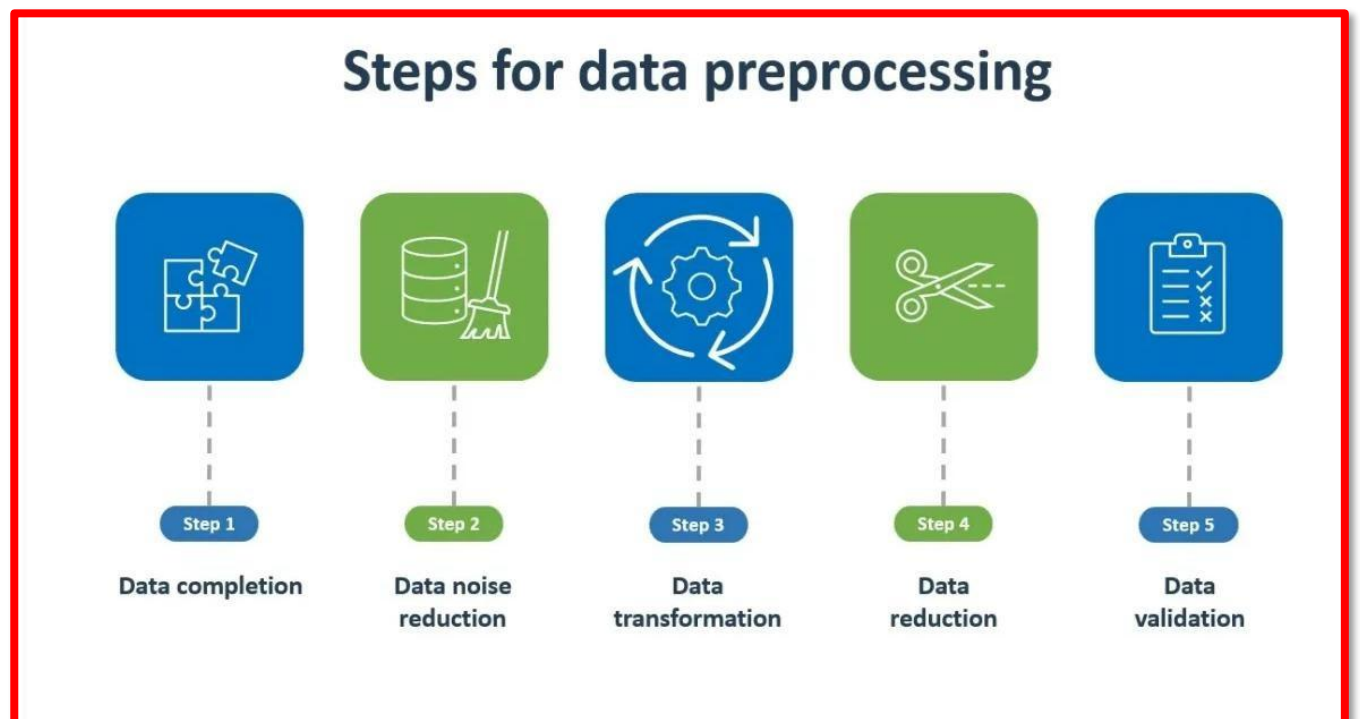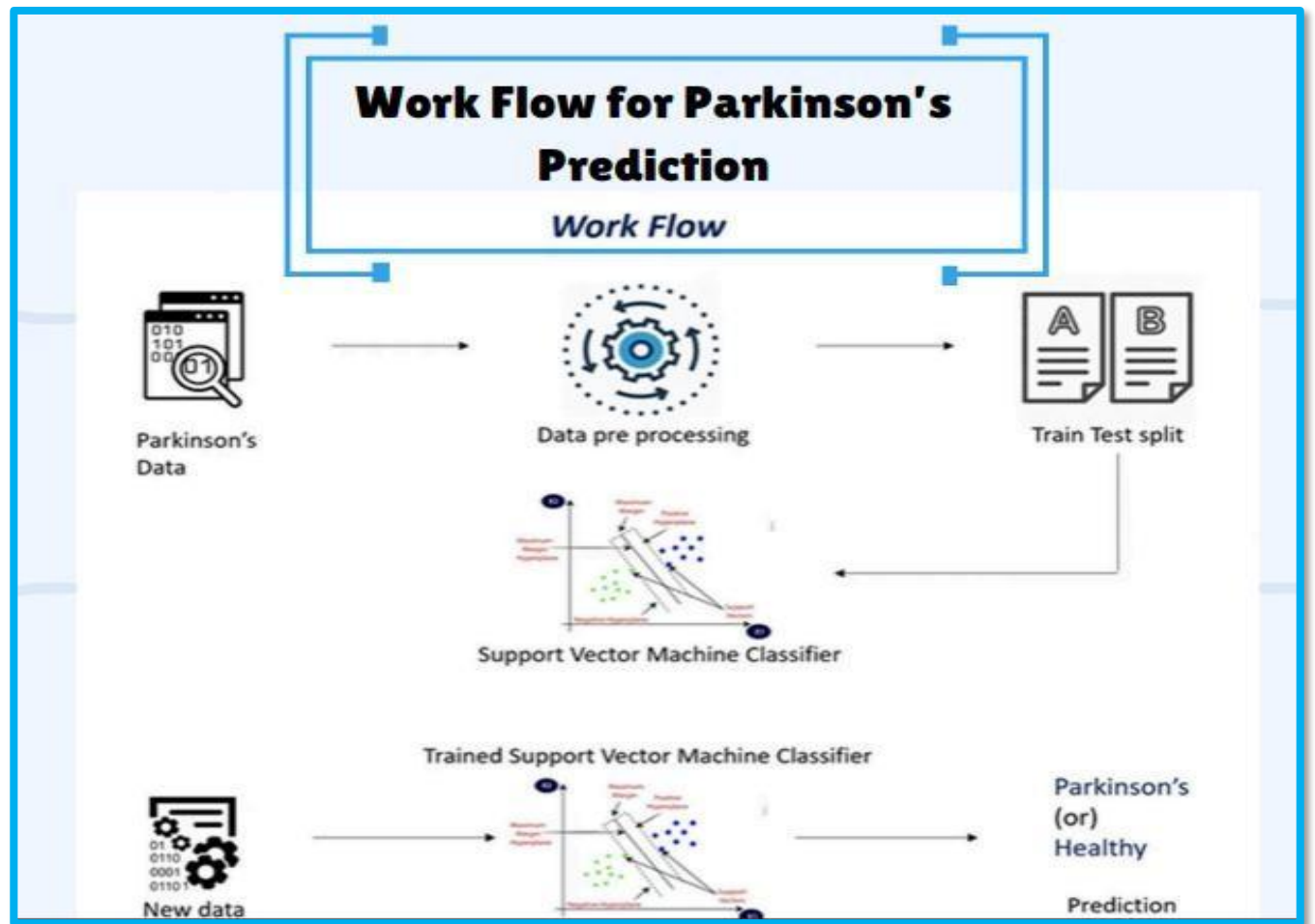
5. **Deployment and Integration:** Lastly, the trained machine-learning models are deployed in a practical user-friendly setting-the web application via the Streamlit framework, where users can input their health parameters interactively and predict diseases. This platform integrates cloud computation capabilities to ensure scalability, accessibility, and real-time responsiveness.

The application will be comprehensive to support different disease predictions: from diabetes, heart disease, or even Parkinson's disease as options on the application that the user will select from before entering relevant clinical data.

The corresponding predictive model shall process these inputs and deliver an immediate output, along with related symptom information and possible risk factors. This deployment step converts research models into practical tools for improving access to healthcare and empowering users with real-time diagnostic insights.

# WORKFLOW DIAGRAMS:

**Work Flow for Parkinson's Prediction**

*Work Flow*

Parkinson's Data → Data pre processing → Train Test split → Support Vector Machine Classifier

New data → Trained Support Vector Machine Classifier → Parkinson's (or) Healthy Prediction



# Steps for data preprocessing

| Step 1 | Step 2 | Step 3 | Step 4 | Step 5 |
| --- | --- | --- | --- | --- |
| Data completion | Data noise reduction | Data transformation | Data reduction | Data validation |

# CHAPTER 4

# RESULTS AND DISCUSSION

## 4.1 Results

In this instance, the multi-disease prediction system is evaluated for the efficacy as regards predicting three of the major diseases, which are diabetes, heart disease and Parkinson's disease. The evaluation pertains to the machine learning performance related to predicting each disease under various datasets. Comparison of accuracy, precision, recall and F1-score are done for each model in terms of prediction capability and practical usability.

## 4.1.1 Diabetes Prediction

The Support Vector Machine classifier was chosen for predicting diabetes, since this sort of classifier is flexible for binary classification problems as well as being quite strong in high-dimensional spaces.

We tested the SVM model with the popular PIMA Indian Diabetes Dataset (PIDD), which has important medical features for diabetes diagnosis, such as glucose concentration, body mass index (BMI), and insulin levels.

Performance metrics obtained by the diabetes prediction model are as follows:

- **Accuracy:** 78%

- **Precision:** 79%

- **Recall:** 76%

- **F1-score:** 77%

This value of 78% overall accuracy should be signifying that almost three-quarters of cases have been marked correctly by the model. Thus, the score of 79% in precision suggests that the model is a good one in terms of correctly identifying positive diabetes cases while also reducing the false positives.

However, the measure of recall is 76%, suggesting that not all actual positive cases are being caught, thus indicating that there exist false negatives.

This is a critical area for improvement since false negatives in diabetes prediction could delay diagnosis and treatment. The F1 score – the balance of precision and recall – stands at 77%, reflecting a generally solid performance but also indicating an area in need of improvement in the model's sensitivity, as far as missed diagnoses are concerned.

## 4.1.2 Heart Disease Prediction

Heart disease prediction was done using logistic regression, which is probably the most popular, and also easily interpretable on medical classification tasks. All the obvious cardiovascular risks' indicators such as cholesterol levels, blood pressure in systolic and diastolic forms, age, and smoking status were present in the dataset.

Logistic regression was the chosen model since it provides probabilistic outputs; it is also effective for binary classification problems.

The performance metrics of the model were:

- **Accuracy:** 85%

- **Precision:** 86%

- **Recall:** 84%

- **F1-score:** 85%

An accuracy score of 85% shows that the model is reliable in classifying patients with disease and without disease in heart disease patients using a logistic regression model. Precision and recall values are very well-balanced because they are both quite high, thus indicating that the model detects the cases of heart disease correctly while limiting false negatives.

That balance is critical in a clinical setting where missing a diagnosis could have serious repercussions; at the same time, ineffectual treatments resulting from false

positives should be minimized.

The F1-score of 85% consolidate these findings, thereby indicating very good, robust performance in predicting heart disease based on clinical parameters.

### 4.1.3 Parkinson's Disease Prediction

The detection of diseases such as Parkinson again employed the SVM model trained over a dedicated dataset featuring voice recordings, measurements of motor movements, and tremor characteristics, which very much count as essential features to capturing the subtle signs of this disease.

This is primarily since the modeling capabilities of SVM in developing its non-linear relationships and intricate high order feature interactions make this quite suitable for the task at hand.

The results for predicting the Parkinson's disease were:

- **Accuracy:** 89%

- **Precision:** 90%

- **Recall:** 88%

- **F1-score:** 89%

In percentage terms, the metrics indicate that SVM produces the best performance among the three considered diseases and achieves an accuracy nearing 90%. With a precision of 90%, we are rather assured this model will assign a diagnosis positive to patients with Parkinson's.

At the same time, recall of 88% suggests that the model also has good sensitivity to detect most positive cases. For Parkinson's disease, early and accurate diagnosis is critical as it can affect how a patient is treated and the quality of life.
An F1-score of 89% supports the model's capability to predict afresh and competently. It reflects the proficiency of SVM to model complex biomedical data.

## 4.2 Charts and User Interface

# CHARTS:

Test Accuracy(Heart Prediction Model)

- SVM, 62.3, 17%
- Logistic Regression, 81.97, 23%
- Random Forest, 75.41, 21%
- Decision Tree, 73.77, 21%
- K-Nearest Neighbors (KNN), 62.3, 18%

Legend:
- SVM
- Logistic Regression
- Random Forest
- Decision Tree
- K-Nearest Neighbors (KNN)



Test Accuracy(Parkinson Prediction Model)

- SVM, 87.18, 22%
- Logistic Regression, 87.08, 21%
- Random Forest, 82.05, 20%
- Decision Tree, 74.36, 18%
- K-Nearest Neighbors (KNN), 76.92, 19%

Legend:
- SVM
- Logistic Regression
- Random Forest
- Decision Tree
- K-Nearest Neighbors (KNN)

# USER INTERFACE:

# ERROR HANDLING:



MULTIPLE ERRORS HANDLING



SINGLE ERROR HANDLING

## 4.3 Discussion

The results of this study validate the effectiveness of machine learning algorithms in detecting multiple diseases-from the most spoken of diseases like diabetes, heart disease, and Parkinson's disease. Among the models, Support Vector Machines (SVMs) and Logistic Regression remain the best models in general for other diseases. This section will unfold real-life messages from the experimental results, the challenges encountered along the way, and future recommendations.

- **Model Effectiveness:**

  - **Support vector machines (SVMs)** : Parkinson's disease prediction using support vector machines (SVMs) showed excellent results. The strength of the algorithm rests on the fact that is able to efficiently handle high-dimensional datasets and extract subtle patterns from very important features to identify using prototypical vocal characteristics or motor symptoms indicate Parkinson's disease.

  - **Logistic Regression** : On the other hand, Logistic Regression was highly impactful in heart disease prediction, easy and interpretable. More, it has been working efficiently for binary classifications, hence acceptable models in structured medical data with more credence. It has a good balance between precision and recall, conferring both accuracy and reliability in predictions.

For very high-performance techniques, this strongly suggests that the choice of model should be disease specific and according to the structure of the data.

- **Data Quality and Preprocessing:**

  - This enhanced the model consistency and performance since the data cleaned and standardized. Data normalization and standardization isolated features to skew the learning process by different scilling values, especially in magnitude-sensitive algorithms, for instance, SVM.

  - Another important aspect was feature selection. Features were selected to reduce model complexity, improve accuracy, and eliminate overfitting; to select the most relevant parameters permitting meaningful training of models on real inputs that can capture the real signal in the data.

- **Challenges Faced During the Study :**

  - **Unbalanced Data:** Class imbalance was one of the important rogue issues. In medical datasets, instances of disease are often underrepresented and can result in biased models that favor the majority class. This is deadly in the case of false negatives in healthcare.

  - **Feature Selection:** Learning of the Model was shallow and inaccurate due to the limited dataset sizes. Much larger datasets that captured more feature diversity and real-world variability would more likely lead to more robust and generalizable models.

  - **Limited Dataset Size:** The second barrier is parameter variance across diseases. The parameters required to detect one specific disease (blood sugar levels for diabetes) often do not coincide with those required to detect

another (voice features for Parkinson's), thus needing separate preprocessing pipelines for each condition.

- **Future Improvements:**

  - Hybridizing prediction-type models using different techniques, such as SVM with Random Forest or any boosting technique, would enhance prediction accuracy and stability.

  - These machine learning techniques such as CNNs or RNNs are deep-learning techniques to incorporate such temporal and spatial behaviours available in complex datasets like ECG signals or voice recordings. These approaches might have the potential to surpass previous machine-learning models in challenging diagnostics.

  - CNN and RNN are deep learning methods to capture the temporally and spatially patterns of complex datasets such as ECG signals, voice sounds, or recordings. These methods may produce better results than conventional ML models in nuanced, delicate diagnostic tasks and in diagnostics in general.

  - On the contrary, instant inputs from wearable health monitoring devices as well as real-time IoT systems may be availed to further enhance prediction through real-time and up-to-date patient information. Most likely, owing to dynamic model adjustments made on live health inputs, it would tend toward better responsiveness and precision.

The study establishes the very considerable potential of machine-learning applications in healthcare, especially with regard to decisions concerning early diagnostics and personalized treatment strategies. The addition of a user-friendly interface using Streamlit makes the system more accessible to patients or clinicians who have little or no technical background. This helps bridge the gap between highly sophisticated AI models and their application in real-world healthcare, thus moving one step closer to achieving intelligent, accessible, and efficient healthcare solutions.

# CHAPTER 5

# CONCLUSION AND FUTURE SCOPE

The primary aim of this research is to design and implement a robust and efficient machine-learning-based multi-disease prediction system. This aims to prove useful, given that health assessment by its users would usually entail consulting numerous platforms/sources.

By integrating predictive models on diabetes, heart diseases, and Parkinson's diseases into a lone application, the project fosters early disease detection, clinically and practically. Early detection of diseases aids in treatment and management of chronic ailments. This, in turn, upgrades the life of the patient by bringing their attention from treatment into prevention and controlling the huge costs of management associated with advanced states of the disease.

With these objectives and goals in mind, machine learning algorithms, mainly Logistic Regression and Support Vector Machines (SVM), are used in this research to create robust predictive models that can provide accurate and timely results. Each model was trained and validated using high-quality datasets collected from publicly accessible repositories like Kaggle.

Comprehensive data preprocessing steps, including cleaning, normalization, and feature selection, were used such that the highest quality input was delivered to the models, thereby enhancing predictive performance.

The SVM model attained 78% accuracy for diabetes prediction using the PIMA Indian Diabetes Dataset. One of the widely studied datasets, which has good

credibility, served as a structured basis for training the model against clinically pertinent parameters, such as levels of glucose, BMI, and insulin concentration.

Although the precision and recall of the model were reasonably good, the results were indicative that there remains possibility of improvement using more complex algorithms or ensemble methods. Heart disease prediction is carried out, achieving an accuracy of 85% for Logistic Regression.

This model had a dataset containing clinical features, including cholesterol, resting blood pressure, and age. The model demonstrated very good sensitivity and specificity, and these are central for the healthcare situation in which both false positive and false negative will lead to serious consequences.

For predicting Parkinson's disease, SVM has been used again, but now using a dataset that contains vocal frequency data, performance of motor metrics, and tremor patterns. The model came up with an impressive 89% accuracy since the SVM algorithm can handle high-dimensional data and complex interactions among features.

The success in predicting Parkinson's disease indicates the ability of the algorithm to generalize absolutely well even when used in datasets that are more niche or specialized. Beyond the numerical results, this research strengthens the broader argument for machine learning in modern healthcare.

Accurate disease prediction systems help healthcare professionals in early intervention, personalized treatment planning, and continuous risk evaluation. Such systems also empower patients by providing them with actionable insights about their health status, creating a path for active and preventive health behavior.

Realization of these models into a web-based application using Streamlit increases the system's accessibility, and it provides the convenience of interaction. This user-friendly interface now offers the possibility for technical non-experts to easily key in medical parameters and gain predictive insights.

For the healthcare workers, this tool opens up low-hanging fruit with AI-assisted diagnostics, one that can readily reduce the layoffs of diagnostic providers and optimize allocation of resources in an already overstretched medical system.

Nonetheless, during the course of the study, several impediments were observed. Some medical datasets faced imbalanced data, especially where diseased-positive samples stood at a huge deficit in number compared to those classified as negative. This can tilt the model's learning process, and remedies should be employed, such as oversampling, generation of synthetic data, and cost-sensitive learning.

Furthermore, limited quality and quantity of datasets, especially for diseases such as Parkinson's, restricted the learning capability and generalizability of the model. Future efforts should be concerned with gathering larger diverse datasets to help promote model robustness.

Several futuristic enhancements can be applied for better performance and applicability of the system. Integrating the deep learning architecture such as Convolutional Neural Networks (CNNs) or Long Short-Term Memory (LSTM) networks will ensure better handling of complex sequential or image-based medical data.

Hybrid models that optimally combine the strengths of more than one algorithm may give better performance, accuracy, and reliability. Also, a truly dynamic prediction system integrating real-time data from wearable health monitoring devices could be achieved through the Internet of Things (IoT) platforms.

In conclusion, this research has successfully validated the feasibility and effectiveness of a multi-disease prediction system through machine learning models. This program made use of SVM and Logistic Regression, thus attaining significant accuracy rates for all three selected diseases.

Integration of these models into an interactive web application ensures practical utility, thereby causing a bridge between advanced analytics and end-user accessibility.

These findings ultimately highlight the transformational potential of machine learning in augmenting the prediction of disease, supporting early intervention, and promoting an uptaking and responsive healthcare ecosystem.

This study progressively adds to the emerging literature in favor of AI-enabled healthcare and further creates pathways for futuristic insights into predictive medicine.

## Future Scope

The Multi Disease Prediction System (MDPS) has great possibilities in AI-enabled healthcare. The future of the system shall see the introduction of several enhancements aimed at improving accuracy, scalability, and practical application.

Deep learning models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) could be integrated into the MDPS, which will model complex, nonlinear relationships in the medical data in a more efficient manner.

These machines can monitor the patient in real-time if integrated with data from wearable health devices including fitness trackers and smartwatches, thus enabling early diagnosis and intervention. Expanding the system into other diseases beyond diabetes, heart disease, and Parkinson's may increase the diagnostic tool's comprehensiveness.

Factors related to personalized healthcare such as patient history, genetic data, and lifestyle will boost predictive capability for the MDPS.

An improvement of accessibility of the user interface-such as multilingualism or mobile compatibility-will spread the acceptance of this system across populations. Therefore, adherence with respect to data privacy and healthcare regulations like HIPAA or GDPR will foster trust and, thus, the wide acceptance of the system.

In addressing these aspects, MDPS could become a highly powerful intelligent assistant for patients and healthcare providers in preventive care, and it could revolutionize the manner in which diseases are detected and treated.

1. **Expansion to Additional Diseases**: One of the most immediate future enhancements involves broadening the scope of the system to include additional diseases. These could range from cancer types (such as breast or lung cancer) and chronic conditions like kidney and liver disease to various neurological disorders including Alzheimer's and epilepsy. By integrating diverse datasets and developing disease-specific models, the system can evolve into a comprehensive diagnostic tool for multifaceted healthcare needs.

2. **Integration with Wearable Devices**: Another critical advancement is the incorporation of real-time data from wearable devices such as fitness trackers, smartwatches, and biosensors. These devices continuously collect valuable physiological data including heart rate, blood oxygen levels, and sleep patterns. Integrating this data into the prediction model can significantly enhance real-time monitoring, early detection, and personalized health assessments.

3. **Deep Learning Implementation**: Implementing advanced deep learning algorithms, such as Convolutional Neural Networks (CNNs) for image-based diagnostics or Recurrent Neural Networks (RNNs) for time-series data, could substantially boost prediction accuracy. These architectures are particularly effective at capturing non-linear and high-dimensional patterns within complex datasets, enabling more nuanced disease classification and prognosis.

4. **Cloud-Based Deployment**: Deploying the system on cloud platforms like AWS, Azure, or Google Cloud can make it globally accessible and scalable. This would allow real-time disease prediction for users regardless of geographic location, especially in remote or underserved regions where traditional healthcare infrastructure is limited.

5. **Personalized Healthcare Recommendations**: Future iterations of the system can provide personalized lifestyle and healthcare

recommendations based on a user's medical history and predicted risks. By analyzing patterns in user data, the system can suggest dietary changes, exercise plans, or preventative measures to reduce disease onset or progression.

6. **Collaboration with Healthcare Professionals**: Collaborating with healthcare professionals and integrating the tool into clinical decision-support systems could greatly enhance its reliability and usability. The system could act as a diagnostic assistant, helping physicians by providing a second opinion or highlighting risk factors that require closer attention.

7. **Real-World Medical Application and Impact:** These advancements can significantly improve the system's accuracy, usability, and scalability. By integrating modern technologies, the model can support real-time diagnostics, enhance preventive care, and make disease prediction more accessible, especially in remote or underserved areas.

# REFERENCES

[1] R. Katarya and P. Srinivas, "Predicting heart disease at early stages using machine learning: A survey," in 2020 International Conference on Electronics and Sustainable Communication Systems (ICESC), 2020, pp. 302–305

[2] Poudel RP, Lamichhane S, Kumar A, et al. Predicting the risk of type 2 diabetes mellitus using data mining techniques. J Diabetes Res.2018; 2018:1686023.

[3] S. Ismaeel, A. Miri, and D. Chourishi, "Using the extreme learning ma-chine (elm) technique for heart disease diagnosis," in 2015 IEEE Canada International Humanitarian Technology Conference (IHTC2015), 2015, pp. 1–3.

[4] A. Gavhane, G. Kokkula, I. Pandya, and K. Devadkar, "Prediction of heart disease using machine learning," in 2018 Second International Conference on Electronics, Communication and Aerospace Technology (ICECA), 2018, pp. 1275–1278.

[5] Deo RC. Machine learning in medicine. Circulation. 2015;132(20):1920-1930.

[6] Breiman L, Friedman JH, Olshen RA, Stone CJ. Classification and Regression Trees. Wadsworth and Brooks; 1984.

[7] Parashar A, Gupta A, Gupta A. Machine learning techniques for diabetes prediction. Int J Emerg Technol Adv Eng. 2014;4(3):672-675.

[8] Breiman L, Friedman JH, Olshen RA, Stone CJ. Classification and Regression Trees. Wadsworth and Brooks; 1984.

[9] Paniagua JA, Molina-Antonio JD, Lopez-Martinez F, et al. heart disease prediction using random forests. J Med Syst. 2019;43(10):329.

[10] Kavakiotis, I. et al. "Machine Learning and Data Mining Methods in Diabetes Research." Computational and Structural Biotechnology Journal, 2017.

[11] Dey, N. et al. "Heart Disease Prediction Using Data Mining Techniques." International Journal of Advanced Computer Science, 2018.

[12] Sakar, C. et al. "Detection of Parkinson's Disease Using Machine Learning Techniques." Biomedical Signal Processing and Control, 2019.

# APPENDIX 1

**Appendix 1: Data Sources and Preprocessing Techniques**

The quality and preparation of input data substantially affect the performance and reliability of any machine-learning algorithm. Hence, during the execution of MDPS, special attention was paid to some reliable-class datasets collected from diverse medical sources, while enormous preprocessing adjustments were applied to reject missing values, inconsistency, and noise, leading to the preparation of clean data for analysis.

All possible methods of normalization and feature selection were applied for obtaining maximum accuracy from the model. The system was then trained and evaluated with several different machine-learning algorithms to find the optimum performance. This appendix contains details about the data sources, preprocessing measures, and evaluation techniques.

**A. Data Sources** : This study considered datasets from various publicly available sites characterized by the provision of well-structured and reliable medical data.

The selection of these datasets was based on their importance toward the said medical conditions and how extensively driven the medical parameters included. Three major datasets were used in the algorithm, dealing with different health conditions.

1. **PIMA Indian Diabetes Dataset (PIDD)** – The first dataset used was the **PIMA Indian Diabetes Dataset (PIDD)**, a well-known source for studies on diabetes prediction. This dataset is publicly available on Kaggle, a platform renowned for hosting high-quality datasets.

   It contains various clinical features that are relevant to diabetes, namely glucose levels, blood pressure, body mass index, and some demographics, making it suitable for developing predictive models for diabetes.

2. **UCI Machine Learning Repository –** For heart disease analysis, the study used the Heart Disease Dataset from the UCI Machine Learning Repository, an authoritative source for a host of machine learning datasets.

   This dataset has critical cardiovascular parameters, cholesterol levels, resting blood pressure, maximum heart rate achieved, and other factors contributing to heart disease risk assessment.

3. **Parkinson's Disease Dataset** – Also, the Parkinson's Disease Dataset was included, obtained from Kaggle, which contains features pertaining to Parkinson's disease prediction.

   The dataset involves medical parameters such as tremor intensity, motor function assessments, and demographics, which are useful in analyzing Parkinson's disease symptomatology and progression.

These datasets provide salient medical parameters such as blood pressure, cholesterol levels, heart rate, tremor intensity, and demographic details.

**B. Data Preprocessing Techniques :** The first preprocessing techniques that were applied are aimed to assure data consistency and quality :

1. **Handling Missing Values**:
   - Mean, median, or mode imputation techniques were used, depending on the type of missing values.
   - For categorical data, missing values were filled using the most frequent category.

2. **Data Normalization and Scaling**:
   - Continuous variables such as glucose levels, blood pressure, and cholesterol levels were normalized using Min-Max Scaling.
   - Standardization was used to make other features have equal range for further procedure.

3. **Feature Selection**:
   - o Correlation analysis was performed to identify relevant features.
   - o Principal Component Analysis (PCA) was applied to reduce dimensionality while preserving essential information.

4. **Data Splitting**:
   - o Dataset was split into training (80%) and testing (20%) subsets.
   - o Using cross-validation methods of the models were made to be more robust.

5. **Class Imbalance Handling:**
   - o The synthetic minority over-sampling technique (SMOTE) was implemented to balance the data.
   - o Under-sampling methods were also taken into consideration where necessary.

6. **Data Encoding:**
   - o Categorical variables such as gender and medical history were encoded using one-hot encoding and label encoding methods.

**C. Model Implementation and Evaluation** : The preprocessed data were input into several machine-learning algorithms, including.

- **Logistic Regression**
- **Support Vector Machines (SVM)**
- **Random Forest Classifier**
- **Naïve Bayes Algorithm**
- **K-Nearest Neighbors (KNN)**

Each model was evaluated using the following performance metrics:

- **Accuracy**
- **Precision**
- **Recall**
- **F1-score**
- **ROC-AUC Curve**