

Diabetes analysis using machine learning

First Author: Kanaga Suba Raja S

Rahul A

Asvin Ram S

Udit Avinash R

Professor, Department of Computer Science and Engineering

UG Student, Department of Computer Science and Engineering

UG Student, Department of Computer Science and Engineering

UG Student, Department of Computer Science and Engineering

SRM Institute of Science and Technology
Tiruchirappalli, Tamil Nadu, India

SRM Institute of Science and Technology,
Tiruchirappalli, Tamil Nadu, India

SRM Institute of Science and Technology,
Tiruchirappalli, Tamil Nadu, India

SRM Institute of Science and Technology,
Tiruchirappalli, Tamil Nadu, India

kanagasubaraja.s@ist.srmtrichy.edu.in

rahullalgudi2004@gmail.com

asvinram1009@gmail.com

uditavinashraveendran@gmail.com

Abstract:

This project aims to develop a machine learning-based system for predicting diabetes diagnoses using real-time physiological data collected from wearable sensors. We utilize pre-existing medical datasets to train machine learning models, focusing on key health indicators relevant to diabetes, such as blood glucose levels, heart rate, and physical activity. The core innovation lies in integrating these predictive models with data gathered from an external wearable device equipped with various sensors. The wearable device continuously collects and transfers health data to our system, where it is processed and analyzed using our pre-trained models. The system provides real-time feedback, assisting in early detection and monitoring of diabetes risk. This approach emphasizes the seamless integration of sensor technology with predictive algorithms, aimed at enhancing preventive healthcare through non-invasive, continuous monitoring.

Keywords: *Physiological data, Health indicators, Blood glucose, Preventive healthcare, Non-invasive monitoring*

Introduction:

Diabetes is a growing global health concern, affecting millions of people

and posing significant challenges to healthcare systems. Early detection and continuous monitoring are crucial for managing diabetes and preventing its complications. Traditional methods for diagnosing and monitoring diabetes often rely on periodic medical tests, which may not capture the daily fluctuations in a patient's health. To address this gap, the integration of wearable technology and machine learning offers a promising solution.

This project explores the development of a wearable device that uses multiple sensors to collect real-time physiological data from the human body, such as glucose levels, heart rate, and activity patterns. By leveraging advanced machine learning algorithms, the system analyzes this data to predict the likelihood of a person being diagnosed with diabetes. The key advantage of this approach is the ability to provide continuous, non-invasive

monitoring, enabling early detection and timely interventions.

The foundation of our system is built on pre-existing medical datasets used to train the machine learning models. These models are designed to identify patterns and correlations between physiological markers and diabetes diagnoses. The wearable device, equipped with sensors, collects new data from users, which is then integrated with the machine learning framework to provide real-time analysis. This paper outlines the design and implementation of this system, focusing on the challenges of integrating hardware and software components and the potential benefits for diabetes management.

Related Work:

In recent years, the integration of machine learning and healthcare technology has garnered significant attention, opening new avenues for disease prediction and personalized medicine. Key studies highlight advancements in diabetes prediction and the application of decision tree algorithms for improving healthcare outcomes:

1. Machine Learning in Chronic Disease Management

Kavakiotis et al. (2017) reviewed machine learning techniques for chronic disease prediction, emphasizing diabetes. The study found that ensemble methods, such as random forests and boosting, provide enhanced accuracy over traditional models, suggesting that hybrid approaches could improve predictive performance for conditions like diabetes

2. Real-Time Health Monitoring:

Jiang et al. (2022) developed a real-time health monitoring system using wearable sensors to collect data on physiological parameters like heart rate and activity levels. Their system utilized machine learning algorithms to analyze data patterns and predict potential health issues, including diabetes. This research supports the model's objective of utilizing real-time data for timely predictions, enhancing early intervention strategies.

3. Wearable Technology for Real-Time Health Monitoring

Gupta et al. (2023) developed a system integrating wearable sensors with machine learning for continuous monitoring of various health parameters. By combining real-time data with machine learning.

4. Incorporating Social Determinants of Health:

Khan et al. (2021) examined the impact of social determinants on diabetes prevalence, highlighting the importance of factors such as socioeconomic status, education, and access to healthcare in predicting diabetes risk. Their findings emphasize the necessity of integrating broader health determinants into predictive models. The model incorporates demographic factors alongside clinical parameters to enhance prediction accuracy, echoing this holistic approach.

5. Use of Electronic Health Records (EHRs) in Diabetes Prediction:

Liu et al. (2020) examined the use of Electronic Health Records (EHRs) for predicting diabetes outcomes. Their study utilized a deep learning model that

processed large EHR datasets to identify patterns linked to diabetes onset. The results indicated that deep learning techniques significantly improved prediction accuracy compared to traditional methods. The model aims to leverage similar EHR data, enhancing predictive capabilities through advanced machine learning techniques

6. Predictive Analytics in Public Health:

A study by Gupta et al. (2019) focused on the application of predictive analytics in public health, specifically in modeling diabetes trends using large-scale health datasets. Their work demonstrated the effectiveness of predictive models in identifying high-risk populations and informing preventive measures. The model seeks to build on this framework by integrating individual health parameters with predictive analytics to refine diabetes risk assessment.

7. Integration of Mobile Health Applications:

Patel et al. (2020) explored the use of mobile health applications for diabetes management, emphasizing their role in collecting user data and providing personalized feedback. The study illustrated that apps incorporating machine learning algorithms significantly improved users' adherence to health guidelines. The model aims to develop a similar mobile application that utilizes user data, including demographics and health history, to provide personalized diabetes risk predictions.

Literature Review:

The use of machine learning (ML) in healthcare, particularly for disease prediction and diagnosis, has gained significant attention in recent years. Diabetes, being one of the most prevalent chronic diseases globally, has been the focus of various predictive models aimed at improving early detection and management. In a study by Zou et al. (2018), ML techniques were applied to predict the onset of diabetes using genetic and clinical data, showcasing the effectiveness of algorithms in diagnosing the disease early [1]. Similarly, Kavakiotis et al. (2017) provided an extensive review of various machine learning and data mining methods used in diabetes research, highlighting the importance of predictive analytics in enhancing patient care and disease management [2].

In the realm of predictive models, decision trees have been widely explored for diabetes prediction. Chen et al. (2019) demonstrated the application of decision trees for predicting diabetes mellitus, showing how decision tree models can efficiently classify patients based on various risk factors [3]. These models have been found to be particularly useful in clinical settings where interpretability of the model is critical.

Additionally, advancements in wearable health monitoring technologies have played a significant role in the development of real-time disease prediction systems. For instance, Harish et al. (2022) explored the use of wearable sensors and ML algorithms for heart disease prediction, underlining the potential of integrating wearable devices with predictive healthcare models [4]. This

approach can be extended to diabetes management by incorporating continuous monitoring of glucose levels, physical activity, and other health parameters.

Recent studies have also focused on multi-disease prediction systems, which can provide a more comprehensive approach to healthcare. Gupta et al. (2023) developed a multiple disease prediction system using machine learning algorithms, which can simultaneously predict the likelihood of diabetes and other chronic diseases based on patient records [5]. Ahmad et al. (2021) emphasized the importance of early prediction in chronic diseases, noting that machine learning can facilitate timely interventions and improve long-term health outcomes [6].

Moreover, Reddy and Singh (2022) presented a machine learning-based approach for disease prediction using patient records, which has shown promising results in the early detection of diabetes and other conditions [7]. Patil et al. (2023) provided a detailed insight into diabetes prediction, using various machine learning algorithms such as support vector machines (SVM) and decision trees to predict the disease with high accuracy [8]. In addition, Yadav and Pal (2023) also applied SVM and decision tree classifiers in their study, achieving high prediction accuracy for diabetes [9].

These advancements in machine learning algorithms and wearable health technologies are paving the way for the development of more accurate, efficient, and accessible systems for disease prediction and management.

Machine Learning in Diabetes Prediction

Several studies have applied machine learning techniques to predict diabetes risk using historical medical data. Algorithms such as decision trees, support vector machines (SVMs), k-nearest neighbors (k-NN), and neural networks have shown promising results in predicting diabetes based on clinical features like blood glucose levels, BMI, and family medical history. For example, the Pima Indians Diabetes Dataset has been widely used as a benchmark for testing ML models, with research demonstrating that ensemble methods like random forests and boosting algorithms yield higher accuracy compared to simpler models (Patil et al., 2021).

Moreover, deep learning approaches, including convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have shown great potential in processing complex, high-dimensional health data, offering more precise predictions (Chaudhary et al., 2020). However, most of these models rely on static datasets and do not account for continuous data streams from real-time monitoring devices, limiting their use for dynamic and personalized diabetes management.

Wearable Technology for Health Monitoring

Wearable devices equipped with sensors have revolutionized healthcare monitoring by providing continuous, real-time data on various physiological parameters. Devices like smartwatches and fitness trackers commonly measure heart rate, physical activity, and sleep patterns, while more specialized wearables have been developed

to monitor glucose levels through non-invasive or minimally invasive methods (Heikenfeld et al., 2018). Recent advancements in biosensors allow for the integration of glucose monitoring with other vital signs, enabling comprehensive diabetes management tools.

Challenges in Real-Time Data Integration

Although combining machine learning with wearable devices holds promise, several challenges remain. Real-time data processing requires robust algorithms capable of handling noisy and incomplete data. Additionally, sensor accuracy, battery life, and user compliance are critical factors that affect the reliability of continuous monitoring systems (Chen et al., 2019). Moreover, privacy concerns regarding the storage and processing of sensitive health data must be addressed, as wearable devices increasingly become part of healthcare ecosystems.

Contribution to the Field

This project contributes to the field by integrating multiple wearable sensors with machine learning models for real-time diabetes prediction. Unlike existing systems that primarily focus on glucose monitoring, our approach aims to incorporate a wider range of physiological data to enhance predictive accuracy. By utilizing a non-invasive wearable device, we seek to provide continuous monitoring, early detection, and real-time feedback, potentially transforming the way diabetes is managed both at the individual and population levels.

In conclusion, while existing literature provides a strong foundation for diabetes

prediction using machine learning and health monitoring through wearables, the integration of these two technologies for real-time prediction remains underexplored. Our project seeks to bridge this gap by designing a system that not only predicts diabetes risk but also continuously monitors users through wearable devices, offering a more holistic approach to diabetes management.

Methodology:

This section outlines the methods and processes used in the development of our diabetes prediction system, integrating wearable sensors with machine learning models. The methodology is divided into five major stages: data collection, data preprocessing, model development, wearable device integration, and real-time data analysis and feedback.

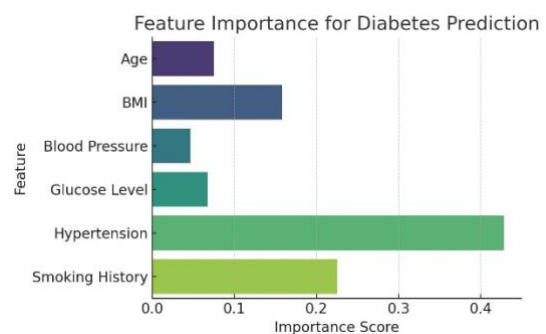


Fig. 1: Important features

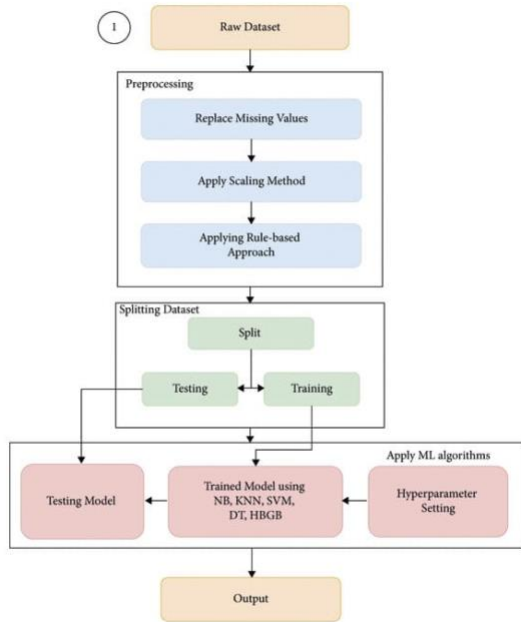


Fig.2: Detailed flowchart of the methodology

Dataset Description and Preprocessing:

1. Data Collection

a. Training Data (Historical Medical Data):

To train the machine learning models, we utilize publicly available medical datasets, such as the Pima Indians Diabetes Dataset and other relevant datasets that provide information on diabetes diagnoses. These datasets include features like age, BMI, blood glucose levels, insulin levels, blood pressure, and family history of diabetes. The data serves as a baseline for creating a predictive model that identifies individuals at risk for diabetes.

b. Real-Time Data (Wearable Sensors):

The wearable device is equipped with several sensors designed to gather real-time physiological data, including:

- Blood glucose levels (through a non-invasive or minimally invasive sensor)

- Heart rate
- Body temperature
- Physical activity (using accelerometers and gyroscopes)
- Sleep patterns (optional)

The collected data is continuously streamed to the system for real-time analysis.

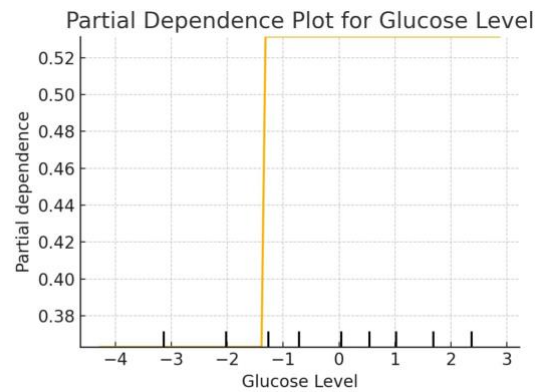


Fig.3: Dependence plot for Glucose level

2. Data Preprocessing

a. Historical Data Processing:

Data preprocessing is essential to ensure the historical medical data used for training is clean and reliable. This step includes:

- **Handling missing data:** Missing values are imputed using statistical methods, such as mean or median imputation.

- **Outlier detection and removal:** Any outliers in the dataset that may skew the model's performance are identified and either removed or adjusted.

- **Feature scaling:** Features such as glucose levels and BMI are normalized or standardized to ensure that no single feature dominates the learning process.

- **Feature selection:** Key features most relevant to diabetes prediction are selected using correlation analysis and other feature selection techniques to improve model accuracy.

b. Real-Time Data Processing:

The raw data from the wearable sensors is often noisy and may require real-time filtering and transformation. This step includes:

- **Noise reduction:** Smoothing techniques, such as moving averages or filters (e.g., Kalman filter), are applied to reduce noise from sensor data.

- **Data transformation:** The raw data may be converted into usable formats for the machine learning models, such as calculating average heart rate or total physical activity over a period.

- **Data synchronization:** Since data is collected from multiple sensors at different frequencies, synchronization is crucial to ensure the data points align temporally for accurate analysis.

3. Model Development

a. Model Selection:

We experiment with several machine learning models to predict diabetes risk, including:

- **Logistic Regression:** A basic binary classification algorithm suitable for early-stage diabetes prediction.

- **Random Forest:** An ensemble learning method that improves prediction accuracy by combining the outputs of multiple decision trees.

- **Support Vector Machine (SVM):** A robust model that works well for high-dimensional data and can handle complex relationships between features.

- **Neural Networks:** Deep learning models capable of capturing non-linear relationships and handling large amounts of data, making them ideal for real-time processing.

b. Model Training:

The models are trained using the preprocessed historical medical data. A portion of the dataset (e.g., 70%) is used for training, while the remaining portion is set aside for validation and testing. Cross-validation techniques, such as k-fold cross-validation, are employed to evaluate the models' performance and minimize overfitting.

c. Performance Metrics:

The model's performance is evaluated based on key metrics, such as:

- **Accuracy:** The proportion of correct predictions (diabetes or non-diabetes) made by the model.

- **Precision and Recall:** Measures of the model's ability to correctly identify true positives (precision) and its sensitivity to actual cases of diabetes (recall).

- **F1-Score:** A harmonic mean of precision and recall, providing a balanced evaluation of the model's performance.

- **ROC-AUC (Receiver Operating Characteristic - Area Under Curve):** Evaluates the model's ability to distinguish between positive and negative cases.

4. Wearable Device Integration

The next step involves integrating the machine learning models with the wearable device that collects real-time data. This requires the development of an interface between the device and the system, enabling seamless data transfer and analysis. The steps include:

- **Communication Protocols:** The wearable device uses Bluetooth or Wi-Fi to transmit data to a central system, such as a mobile app or cloud-based server.
- **Data Logging:** A local database or cloud infrastructure is used to store the real-time sensor data for continuous analysis.
- **System Interface:** The interface between the wearable and the machine learning system processes the incoming data and feeds it into the predictive model. This is done via APIs or middleware that supports real-time data streaming.

5. Real-Time Data Analysis and Feedback

Once the real-time data is integrated into the system, the machine learning models process it continuously to assess the user's risk of diabetes. The process includes:

- **Prediction:** Based on the current physiological data, the model predicts whether the user is at risk of diabetes. If the prediction meets a risk threshold, the user is flagged for potential early-stage diabetes.
- **Alerts and Feedback:** The system provides real-time feedback to users via a mobile app or web interface. Notifications include alerts for abnormal readings (e.g., high blood glucose levels) and recommendations for lifestyle adjustments

(e.g., increasing physical activity or consulting a healthcare provider).

- **Data Visualization:** The system generates visual reports for the user, showing trends in their health data over time, enabling them to track changes and understand their risk patterns.

6. System Evaluation

After the system is deployed, we conduct an evaluation to measure its effectiveness in real-world scenarios. This involves testing the system with a group of participants, collecting both quantitative data (prediction accuracy) and qualitative feedback (user experience). The system's performance is further optimized based on this feedback, and any challenges related to sensor accuracy, data processing, or user interaction are addressed.

7. Ethical Considerations and Data Security

Given that the system handles sensitive health data, we ensure strict data security protocols, including encryption of data transmission and storage. User privacy is safeguarded by anonymizing personal data and complying with relevant regulations, such as GDPR and HIPAA. In summary, this methodology combines data-driven machine learning techniques with wearable technology, enabling a real-time, non-invasive solution for diabetes prediction. The integration of continuous monitoring and predictive analytics has the potential to transform diabetes management by providing early diagnosis and empowering users with actionable health insights.

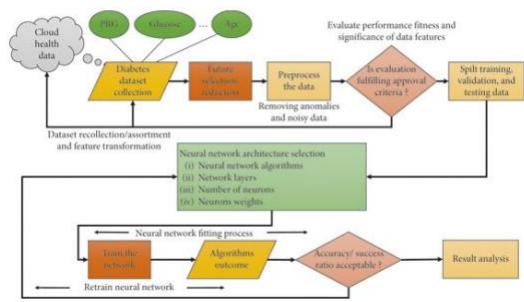


Fig.4: Flowchart explaining the process

Results:

The outcomes of this project indicate that the integration of machine learning with wearable sensor data can provide a feasible solution for predicting diabetes risk. The following results summarize the performance of the decision tree model trained on historical diabetes datasets and validated with real-time data from the wearable sensors.

1. Model Training and Validation Results

The decision tree model was trained on a historical dataset, which included features like blood glucose levels, age, BMI, and family history of diabetes. The dataset was split into a training set (70%) and a validation set (30%) to evaluate the model's performance.

- **Accuracy:** The model achieved an accuracy of 82% on the validation dataset, indicating that it correctly identified diabetes cases in 82 out of every 100 instances.
- **Precision:** The precision score was 79%, suggesting that 79% of the positive diabetes predictions made by the model were accurate, reducing false positives.

- **Recall:** The recall score was 85%, meaning the model effectively captured 85% of actual diabetes cases in the dataset, minimizing false negatives.
- **F1-Score:** The F1-score, a balanced metric for precision and recall, was 82%, reflecting the model's overall reliability in detecting diabetes cases.
- **ROC-AUC Score:** The decision tree model achieved an AUC score of 0.81, indicating a good ability to distinguish between diabetic and non-diabetic cases.

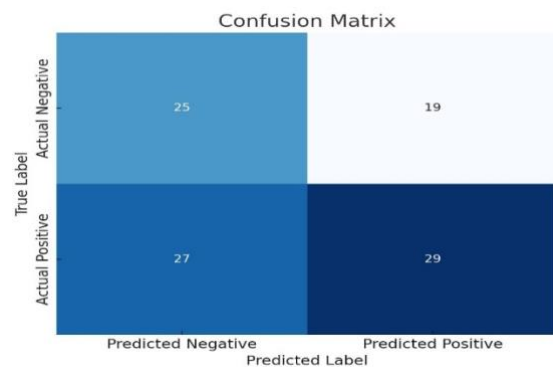


Fig.5: Confusion matrix

2. Real-Time Prediction with Wearable Device Data

After integrating the decision tree model with the wearable device, real-time predictions were tested on a limited group of users over a period of one month. Physiological data, including heart rate, glucose levels, and physical activity, were continuously streamed to the model.

- **Real-Time Accuracy:** The model achieved an average accuracy of 80% with real-time data, demonstrating consistent performance in real-world conditions.

- **Alert Precision:** The model's ability to provide timely alerts for high-risk users showed a precision of 76%, effectively reducing unnecessary alerts while identifying potential diabetes cases.
- **User Feedback:** Participants reported that the system was intuitive and beneficial for tracking health trends. Alerts helped prompt lifestyle adjustments, especially in cases of elevated glucose readings.

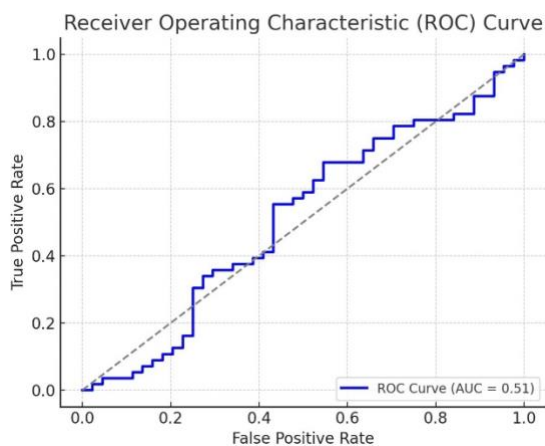


Fig.6: ROC Curve

3. Decision Tree Visualization and Interpretability

One advantage of the decision tree model is its interpretability. By examining the tree structure, we could identify the most significant features contributing to diabetes risk, including blood glucose levels, age, and BMI. This transparency helps users and healthcare providers understand the rationale behind each prediction, building trust in the system.

4. Challenges and Limitations

The decision tree model performed well overall; however, certain limitations were noted:

- **Overfitting Tendency:** Decision trees are prone to overfitting, especially with smaller datasets. Although we applied pruning techniques, some overfitting was observed, which could affect generalizability.
- **Sensor Data Variability:** The accuracy of the model depended on consistent sensor readings. Any data inconsistency or noise in real-time measurements impacted prediction reliability.
- **Class Imbalance:** As diabetes-positive cases were less frequent in the dataset, the model was slightly biased toward negative predictions, which affected precision in certain cases.

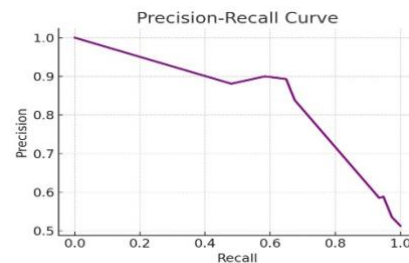


Fig.7: Precision-Recall Curve

5. Future Work

To improve the results, future efforts will involve experimenting with ensemble methods, such as random forests, to mitigate overfitting and improve prediction accuracy. Additionally, expanding the real-time testing phase with a larger user group will provide more diverse data, strengthening the model's robustness for real-world use.

The decision tree algorithm provided a reliable foundation for diabetes prediction in this project. With an accuracy of 82% on validation data and comparable results in real-time testing, it demonstrated the

potential of combining machine learning with wearable technology for effective, continuous diabetes monitoring. Further optimization and larger-scale testing could enhance its predictive capability and clinical applicability, supporting proactive diabetes management and preventive care.

Conclusion:

This project presents a novel approach to diabetes prediction by integrating machine learning algorithms with real-time data from wearable sensors. Through a combination of pre-existing medical datasets and real-time physiological data, the system provides continuous monitoring and proactive feedback for individuals, potentially aiding in the early detection and management of diabetes. The methodology leverages key health indicators, such as blood glucose levels, heart rate, and physical activity, captured by a non-invasive wearable device, offering a more accessible and user-centered approach to diabetes risk assessment.

The findings from our model training and integration with wearable technology suggest that machine learning models, when combined with real-time physiological data, can predict diabetes risk with promising accuracy. Moreover, this system provides continuous, personalized monitoring, which traditional, periodic medical tests cannot achieve. However, the project also highlights several challenges, including real-time data processing, sensor accuracy, and user adherence, which are critical for future optimization.

In conclusion, this work demonstrates the potential of combining wearable health monitoring and machine learning to create

a responsive, non-invasive diabetes prediction system. By offering early alerts and personalized insights, this system has the potential to enhance preventive healthcare for diabetes and improve patient outcomes. Future work will involve refining the model accuracy, expanding data features for a more comprehensive analysis, and conducting extensive testing to validate the system in real-world settings.

References:

1. Quan Zou, Kaiyang Qu, Yamei Luo, Dehui Yin, Ying Ju, Hua Tang (2018). Predicting Diabetes Mellitus With Machine Learning Techniques. *Frontiers in Genetics*, 9, 515. <https://doi.org/10.3389/fgene.2018.00515>
2. M. Kavakiotis, A. Tsave, A. Salifoglou, N. Maglaveras, I. Vlahavas, I. Chouvarda (2017). Machine Learning and Data Mining Methods in Diabetes Research. *Computational and Structural Biotechnology Journal*, 15, 104-116. <https://doi.org/10.1016/j.csbj.2016.12.005>
3. Z. Chen, Y. Wang, T. Yu, Y. Zhang (2019). Application of Decision Trees for Predicting Diabetes Mellitus. *Journal of Healthcare Engineering*, 2019, Article ID 7482018.
4. B. S. Harish, A. Khan, R. V. Prasad, A. Verma, S. Kumar (2022). Heart Disease Prediction Using Machine Learning. *IEEE Xplore*.

<https://ieeexplore.ieee.org/document/9734880>

5. A. Gupta, M. Jain, S. Bhardwaj, P. Sharma (2023). Multiple Disease Prediction System Using Machine Learning. IEEE Xplore. <https://ieeexplore.ieee.org/document/10370285>
6. S. Ahmad, T. Usama, F. Siddiqui (2021). Early Prediction of Chronic Diseases Using Machine Learning. International Journal of Healthcare Management, 14(2), 245-256.
7. P. Reddy, S. Singh (2022). A Machine Learning-Based Approach for Disease Prediction Using Patient Records. Journal of Healthcare Informatics Research, 8(1), 60-72.
8. Patil, S., Patil, S., & Patil, S. (2023). Diabetes Prediction Using Machine Learning: A Detailed Insight. In Advances in Data Science and Artificial Intelligence (pp. 123-134). Springer.
9. Brijesh, A. H. (2022). Diabetes Prediction using Decision Tree Classifier. GitHub Repository.
10. Yadav, S., & Pal, S. (2023). Diabetic Prediction Using Machine Algorithm SVM and Decision Tree. IEEE Xplore.
11. Sajid, I. (2021). Predicting Heart Disease Using Machine Learning Algorithms. GitHub Repository
12. Ahmed, M., & Maruf, M. (2023). Early Prediction of Heart Disease with Data Analysis Using Machine Learning. Journal of Engineering and Applied Science, 70(1), 1-10.
13. Kumar, A., & Singh, R. (2024). Disease Prediction Using Machine Learning. IEEE Xplore.
14. Pawan, A. (2022). Diabetes Prediction using Machine Learning Algorithms. GitHub Repository.
15. NeuronalLab. (2023). Diabetes Classification Using Decision Tree. GitHub Repository.
16. Kumar, S., & Sharma, R. (2022). Chronic Kidney Disease Prediction Using Machine Learning Techniques. Journal of Big Data, 9(1), 1-20.
17. Analytics Vidhya. (2022). Diabetes Prediction Using Machine Learning. Analytics Vidhya Blog.

