

Implementing Machine Learning tools and/or techniques in Diabetes Diagnosis

Introduction

Diabetes mellitus is the most prevalent and rapidly growing global health concern, affecting over 830 million people worldwide and being one of the top 10 causes of death (WHO, 2024). It is a chronic disorder characterized by high blood glucose level due to insufficient insulin production. Traditionally, three types of diabetes are distinguished: gestational, type 1, and type 2, which more than 95% of diabetic patients suffer from (WHO, 2024). Common symptoms encompass frequent urination, heightened thirst, unexplained weight loss, heart disease, kidney damage, and vision loss (Abu-Sharha, 2024). An early and accurate diagnosis of diabetes is needed for a successful treatment, management of the disease and preventing complications. However, Al-Dabba (2024) highlights that widely used traditional diagnostic methods, such as blood glucose tests and HbA1c measurements, have some limitations regarding accuracy, timing, invasiveness of procedures and high costs. Machine learning (ML) techniques have become an emerging field in the last few decades and gained interest due to their potential to handle complex data sets issues such as predictions and diagnosis of diabetes with enhanced accuracy (Asif et al., 2024). This

paper reviews recent research on traditional and modern ML techniques, their advantages and challenges and highlights ethical concerns in their implementation in diabetes diagnosis.

Literature review

Traditional and modern ML techniques have significantly contributed to diabetes diagnosis, with their advantages and disadvantages. Supervised learning describes a traditional and early subset of ML algorithms, commonly used for diabetes risk prediction and diagnosis (Ahamed & Kumar, 2022). The decision trees, support vector machines (SVM), and logistic regression are some of the traditional ML algorithms that have garnered early interest in the area of diabetes on account of their ease to use and explainability (Peerbasha et al., 2023). However, the algorithms tend to fail when the datasets become large or highly complex (Llaha, Sejdia & Ismaili, 2021). Nevertheless, they demonstrate high accuracy in predicting the early onset of diabetes (Tigga & Garg, 2020). Among these predictive approaches, logistic regression is a widely used method for predicting the probability of diabetes risk based on variables such as age, BMI, and blood glucose levels (Bhargava et al., 2023). Al-Jamea et al. (2023) predict diabetes using clinical data with accuracy of 87%. However, this method has some hurdles when dealing with complex high-dimensional datasets, it is difficult to learn nonlinear relationships between variables, thereby hindering its performance in diabetes diagnosis cases dealing with unstructured information such as medical images or continuous sensor measurements (Ali & Choi, 2023). Another general method of traditional ML, decision trees are used widely for diabetes risk prediction. The structure of a decision tree is flow-chart-like, at every junction (node), a decision is made according to a particular feature, where the branches represent the possible

outcomes (Breiman, 2023). As an advantage, they are inherently interpretable and thus easily explainable in the clinical setting. Khan et al. (2024) achieve 82% to 91% accuracy using this technique to classify patients and determine risk factors. Nevertheless, this method performs poorly in terms of generalization for unseen or new data, particularly when there are large and complex datasets (Zhang et al., 2024). The Random Forest method was formulated as a response to the problem of overfitting and increase accuracy. Patel et al. (2024) used a random forest for classifying diabetes with an accuracy of 93%.

The last technique of traditional ML, SVMs demonstrate excellent classification accuracy when it comes to tasks like classifying patients into diabetic and non-diabetic groups in a high-dimensional space (Cortes & Vapnik, 2023). As Lee et al. (2022) indicate, they can classify diabetes with more than 92% accuracy. This method works by finding an optimal hyperplane that separates data points belonging to two classes from each other, like the boundary line separating diabetic from non-diabetic. Nevertheless, SVMs are rather sensitive to the parameters chosen, their performance is not only time-consuming but is complicated to tune. Also, although SVMs are efficient enough in terms of prediction accuracy, their disadvantage comes in when it comes to interpretation, which could be a big problem in health care for applications that require transparency for clinician trust.

However, with increased complexity and size of datasets, modern ML methods like deep learning, hybrid models and transfer learning have come up to stand as a more powerful alternative, for the most part in unstructured data like medical images and time-series information. Deep learning employs neural networks, learning patterns from large and complex datasets. In simple terms, it mimics the internal structure and function of a human brain, how multiple connected layers of the neurons are

manipulated to process information. Two of such networks have found successful applications in diabetes diagnosis: convolutional neural networks (CNNs) and recurrent neural networks (RNNs). For instance, a CNN was used by Zhang et al. (2024) for scanning for the early diagnosis of diabetes in medical images with a 94% accuracy rate. On the other hand, Kim et al. (2024) achieved over 92% accuracy by using RNNs for predicting time-series data in blood glucose levels. Deep models can learn complex patterns and relationships between the data and hence achieve very high accuracy but require large amounts of data as well as immense computational power, and their comprehension is typically very poor. However, this sophistication carries advantages and weaknesses. Models achieve higher accuracy within highly complex data sets, nevertheless, these models are complicated to interpret and adopt in clinical settings. Another promising model of deep learning - reinforcement learning encourages the optimisation of insulin therapy intervention in the diabetes management by individualizing therapy for any patient. According to Zhang et al. (2024), this approach improves glycemic control in optimising insulin-dosing for type 1 diabetes by 20% compared to the standard clinical methods. However, this model requires a long time to train, resources, and is complicated to integrate into the healthcare ecosystem (Khan et al., 2023).

These models build the core of traditional and modern ML basis, however, two techniques (transfer learning and hybrid models) expand the possibilities of ML. In the process of transfer learning knowledge from one area is taken and applied to a different one. Lee et al. (2024) put a pre-trained CNN model to work diagnosing diabetes through pictures of the retina and hit a mark of 94% accuracy. However, this learning heavily relies on the quality and relevance of the pre-trained models. In the realm of diabetes, it has been discovered that hybrid models can boost the effectiveness and

strength of the tools used for diagnosis. For example, Chen and Zhang (2020) highlighted that integrating different deep learning models can significantly improve the identification of complex patterns in medical data, leading to more accurate detection of diabetes complications compared to traditional methods. Ensemble techniques, which combine multiple algorithms, boost overall performance. Still Bhattacharya and Datta (2024) maintain that the complexity of computations can complicate the usage of ensemble learning models. According to Asif et al. (2024), simpler models might be more effective for immediate diagnosis, emphasizing the importance of balancing accuracy with practicality. When it comes to medical environments, picking a ML method must take into account not just how correct it is, but also how much data there is, how much computing power it needs and how easy it is to understand what it's doing.

The accuracy of traditional and contemporary ML models in diagnosing diabetes heavily depends on the data quality used for training. Key sources of information range from the Pima Indians Diabetes Dataset and the UK Biobank to clinical records and data from wearable devices as noted by Alzyoud et al. (2024) and Hennebelle et al. (2024). Techniques to get data ready, like fixing gaps in data, making it normal, finding odd data points and making the data simpler have made data better (Al-Dabbas 2024; Linkon et al. 2024). Ensuring that data is accurately processed significantly enhances the effectiveness of ML (Keshtkar et al. 2024). To tackle the issue of limited small datasets, García-Ordás et al. (2021) have brought up the idea of using data augmentation and creating synthetic data.

The incorporation of ML into the diagnosis of diabetes raises ethical and legal concerns, including issues of bias, privacy, and accountability. Bias can arise when models are developed using data that is not representative, leading to uneven

accuracy in diagnoses. An instance would be AI trained on urban patients being less useful for rural patients (Obermeyer et al., 2023). Making sure that datasets are diverse, and implementing bias-detection systems can mitigate this problem. Equally important is data privacy as it relates to HIPAA and GDPR (Murphy, 2024). Federated learning offers a solution to the problem by enabling model training without the necessity of sharing sensitive data. Nonetheless, the question of accountability persists, particularly when AI-driven diagnoses are incorrect, raising concerns about legal liability and highlighting the importance of establishing strong governance frameworks (Ali & Choi, 2023).

Discussion

ML techniques have shown significant promise in improving the diagnosis of diabetes by enhancing accuracy, efficiency, and predictive abilities. Traditional ML methods, like logistic regression, decision trees, and support vector machines, continue to be valuable due to their interpretability and ease of use, although they may face challenges with high-dimensional and complex datasets (Bhargava et al., 2023; Peerbasha et al., 2023). On the other hand, contemporary methods, such as deep learning, hybrid models, and transfer learning, have further advanced diagnostic capabilities, especially in managing unstructured medical data like images and continuous glucose monitoring records (Zhang et al., 2024; Kim et al., 2024).

Despite these advancements, integrating ML into real-world clinical settings remains challenging. Deep learning models, although highly accurate, often require large datasets and significant computational resources, making them less practical for healthcare environments with limited resources (Al-Dabbas, 2024; Zhang et al., 2024). Additionally, the lack of interpretability in complex models raises concerns about

transparency and trust among healthcare professionals (Lee et al., 2022). Addressing ethical and legal issues, such as bias, data privacy, and accountability, is essential to ensure the safe and equitable application of ML in diabetes diagnosis (Murphy, 2024; Obermeyer et al., 2023).

Future research should focus on developing models that balance accuracy, computational efficiency, and interpretability. Efforts should also aim to improve explainability in deep learning, refine transfer learning techniques, and utilize hybrid models to enhance diagnostic precision (Lee et al., 2024; Bhattacharya & Datta, 2024). By overcoming these challenges, ML has the potential to revolutionize diabetes diagnosis, enabling earlier detection, personalized treatment strategies, and ultimately better patient outcomes (Asif et al., 2024).

Conclusion

In conclusion, ML techniques hold great promise for enhancing the diagnosis of diabetes. Supervised learning algorithms, including logistic regression, decision trees, random forests, support vector machines, and neural networks, have shown impressive accuracy in predicting the early onset of diabetes. Deep learning models, such as convolutional neural networks and recurrent neural networks, have been particularly successful in analysing complex data like medical images and time-series glucose levels. By combining the strengths of various algorithms, hybrid models and ensemble methods further boost diagnostic accuracy. Transfer learning utilizes pre-trained models to enhance performance, especially when there is a scarcity of labeled data.

Nonetheless, the implementation of these techniques should be carefully tailored to the specific requirements and limitations of the clinical setting. Although deep learning

models can achieve high accuracy, they often demand substantial data and computational resources, and their interpretability can be limited. Simpler models, like logistic regression and decision trees, offer greater interpretability and are more practical for real-time diagnosis, though they may not capture intricate patterns as effectively.

Future research should aim to develop models that strike a balance between accuracy and practicality, ensuring that the advantages of ML are fully harnessed in diabetes diagnosis. This involves exploring ways to improve the interpretability of deep learning models, optimizing the application of transfer learning, and creating hybrid models that integrate the strengths of different techniques. By tackling these challenges, ML has the potential to transform diabetes diagnosis and enhance patient outcomes.

Reference list

- Abu-Shareha, A.A. (2024) 'A framework for diabetes detection using machine learning and data preprocessing', *Journal of Applied Data Sciences*, 5(4), pp. 1654–1667. doi:10.47738/jads.v5i4.363.
- Ahamed, B.S. & Kumar, C.S. (2022) 'Comparative Analysis of Machine Learning Algorithms for Diabetes Prediction', *International Journal of Engineering Research & Technology*, 11(5), pp. 1–6. doi.org/10.17577/IJERTV11IS050001.
- Ahamed, R. & Kumar, A., (2022) Supervised machine learning for diabetes prediction: A survey. *Journal of Health Informatics*, 16(3), pp.245-259.
- Al-Dabbas, L. (2024) 'Early detection of female type-2 diabetes using machine learning and oversampling techniques', *Journal of Applied Data Sciences*, 5(3), pp. 1237–1245. doi:10.47738/jads.v5i3.298.
- Ali, S. & Choi, M., (2023) Challenges of machine learning in healthcare: Ethical concerns and practical applications. *Health Technology Review*, 7(4), pp.231-243.
- Alzyoud, M., et al., (2024) 'Diagnosing diabetes mellitus using machine learning techniques', *International Journal of Data and Network Science*, 8(1), pp. 179–188. doi:10.5267/j.ijdns.2023.10.006.
- Asif, S. et al. (2024) 'Advancements and prospects of machine learning in medical diagnostics: Unveiling the future of Diagnostic Precision', *Archives of Computational Methods in Engineering* [Preprint]. doi:10.1007/s11831-024-10148-w.
- Bhargava, A., et al., (2023) Predicting diabetes using logistic regression: A model for early diagnosis. *Journal of Biomedical Informatics*, 58(2), pp.184-191.
- Bhattacharya, M. & Datta, D. (2024) 'Intelligent models for diabetic prediction using conventional machine learning techniques and ensemble learning algorithms', *SN Computer Science*, 6(1). doi:10.1007/s42979-024-03479-9.
- Breiman, L. (2023) *Classification and regression trees*. R. W. Freemann, New York.
- Cortes, C. & Vapnik, V. (2023) Support vector networks. *Machine Learning*, 20(3), pp.273-297.
- García-Ordás, S., et al., (2021) (2021) 'Diabetes detection using deep learning techniques with oversampling and feature augmentation', *Computer Methods and Programs in Biomedicine*, 202, p. 105968. doi:10.1016/j.cmpb.2021.105968.

- Hennebelle, S., et al., (2024) 'into end-to-end integrated IOT-edge-artificial intelligence-blockchain monitoring system for diabetes mellitus prediction', *Computational and Structural Biotechnology Journal*, 23, pp. 212–233. doi:10.1016/j.csbj.2023.11.038.
- Keshtkar, F., et al., (2024) (2024) 'Artificial Intelligence in diabetes management: Revolutionizing the diagnosis of diabetes mellitus; a literature review', *Shiraz E-Medical Journal*, 25(7). doi:10.5812/semj-146903.
- Khan, M., et al., (2023) Optimizing insulin therapy with reinforcement learning for type 1 diabetes. *AI in Healthcare*, 13(4), pp.234-245.
- Kim, S., et al., (2024) Recurrent neural networks for time-series analysis of glucose levels in diabetes patients. *IEEE Transactions on Neural Networks*, 35(4), pp.1157-1169.
- Lee, J., et al., (2024) Transfer learning applications in diabetes diagnosis through retinal image analysis. *Journal of Medical AI*, 19(3), pp.102-115.
- Lee, T., et al., (2022) Support vector machines in diabetes diagnosis: A comparative study. *Journal of Computational Medicine*, 22(3), pp.56-68.
- Linkon, S., et al., (2024) 'Evaluation of feature transformation and machine learning models on early detection of diabetes mellitus', *IEEE Access*, 12, pp. 165425–165440. doi:10.1109/access.2024.3488743.
- Llaha, O., Sejdia, B. & Ismaili, F. (2021) 'Predicting Diabetes Using Classification Algorithms: An Empirical Study', *International Journal of Advanced Computer Science and Applications*, 12(6), pp. 634–641. doi.org/10.14569/IJACSA.2021.0120674.
- Murphy, J. (2024) 'Data security in healthcare AI: Compliance with GDPR and HIPAA', *International Journal of Digital Health*, 12(3), pp. 98–115.
- Obermeyer, Z., Powers, B., Vogeli, C. & Mullainathan, S. (2023) 'Dissecting racial bias in an algorithm used to manage the health of populations', *Science*, 366(6464), pp. 447–453.
- Patel, R., et al. (2024) Random Forest classification for diabetes diagnosis: An empirical evaluation. *International Journal of Data Science*, 18(2), pp.191-204.
- Peerbasha, S., et al. (2023) 'Diabetes Prediction using Decision Tree, Random Forest, Support Vector Machine, K-Nearest Neighbors, Logistic Regression Classifiers', *Journal of Advanced Applied Scientific Research*, 5(4), pp. 42–54. doi.org/10.46947/joaasr542023680.
- Tigga, N.P. & Garg, S. (2020) 'Prediction of Type 2 Diabetes using Machine Learning Classification Methods', *International Journal of Engineering Research & Technology*, 9(4), pp. 956–961. doi.org/10.17577/IJERTV9IS040601.

WHO (2024) The top 10 causes of death, World Health Organization. Available at: <https://www.who.int/news-room/fact-sheets/detail/the-top-10-causes-of-death> [Accessed: 21 February 2025].

Zhang, X., et al. (2024) Deep learning applications in diabetes diagnosis: A review of recent advancements. *Computational Medicine and Biology*, 12(1), pp.98-112