

Identification of Risk Factors for Type 2 Diabetes Mellitus: A Machine Learning Approach

 Serkan Budak¹,  Yasemin Karacan²,  İsmail Bacak¹,  Şenay Özer¹

¹Department of Health Care Services, Simav Vocational School of Health Services, Kütahya Health Sciences University, Kütahya, Türkiye

²Department of Nursing, Faculty of Health Sciences, Yalova University, Yalova, Türkiye

Abstract

Introduction: Type 2 diabetes mellitus (T2DM) is a chronic metabolic disease that causes serious health problems worldwide. Multiple risk factors contribute to the development of this disease. Recently, researchers have used artificial intelligence and machine learning (ML) methods to identify these risk factors. This study aims to evaluate the risk factors for T2DM using ML methods.

Methods: This analytical study was conducted over a 2-month period. Data were collected through face-to-face interviews using a personal information form. The obtained data were analyzed using different ML models and performance parameters such as F1 score, accuracy (ACC), and area under the curve (AUC), which represents the area under the receiver operating characteristic curve.

Results: In this study, the most important risk factors for T2DM were identified as age, gender, high blood pressure, genetic predisposition, and education status. Moreover, seven different ML models were analyzed using F1 score, ACC, and AUC parameters, and support vector machine, random forest (RF), and logistic regression (LR) models provided the highest performance.

Discussion and Conclusion: Accurate classification of T2DM risk factors is important for disease prevention and risk assessment in clinical practice. The results suggest that RF or LR models may affect populations with different sociocultural characteristics.

Keywords: Artificial intelligence; Diabetes mellitus; Machine learning; Risk factors

Type 2 diabetes mellitus (T2DM) is a chronic, progressive metabolic disease characterized by hyperglycemia resulting from impaired insulin secretion or insulin action. ^[1] It causes serious health problems worldwide and reduces quality of life through complications such as cardiovascular disease, nephropathy, neuropathy, and retinopathy.

^[2] According to the World Health Organization and the International Diabetes Federation, the global prevalence of diabetes is expected to reach approximately 1.3 billion adults by 2050.^[3,4] In Türkiye, recent epidemiological data show that 13% of adults have T2DM, particularly in urban areas, indicating the need for urgent preventive strategies.^[5]

Cite this article as: Budak S, Karacan Y, Bacak İ, Özer Ş. Identification of Risk Factors for Type 2 Diabetes Mellitus: A Machine Learning Approach. Lokman Hekim Health Sci 2026;6(2):00–00.

Correspondence: Serkan Budak, M.D. Kütahya Sağlık Bilimleri Üniversitesi, Simav Sağlık Hizmetleri Meslek Yüksekokulu, Sağlık Hizmetleri Bölümü, Kütahya, Türkiye

E-mail: serkan.budak@ksbu.edu.tr **Submitted:** 15.09.2025 **Revised:** 28.01.2026 **Accepted:** 07.02.2026 **Available Online:** 21.05.2026



OPEN ACCESS This is an open access article under the CC BY-NC license (<http://creativecommons.org/licenses/by-nc/4.0/>).



T2DM develops through the interaction of multiple genetic, environmental, and lifestyle factors. Genetic predisposition, family history, and certain ethnic origins increase susceptibility.^[6] However, lifestyle-related factors, especially abdominal obesity, sedentary behavior, poor diet, and insufficient sleep, are increasingly contributing to insulin resistance and the rising prevalence of T2DM. Comorbid conditions such as hypertension and dyslipidemia also accelerate disease onset and worsen complications. Therefore, identifying high-risk individuals based on a combination of these variables is crucial for effective prevention and management.^[7,8]

Artificial intelligence (AI) and machine learning (ML) have recently enabled more precise disease prediction by identifying hidden patterns in large datasets. These techniques outperform traditional statistics in uncovering non-linear relationships and improving early diagnosis.^[9] In multifactorial diseases such as T2DM, ML models integrate genetic, clinical, and lifestyle data to produce accurate, individualized risk estimates. Among these, supervised learning models are particularly effective in developing clinical decision support systems and guiding early interventions.^[10]

Studies applying ML to T2DM risk prediction in Türkiye are still scarce, and most rely on limited datasets or single-algorithm analyses. Therefore, this study aimed to identify key risk factors associated with T2DM and to evaluate their predictive importance using multiple ML algorithms. By comparing seven different models, this study sought to determine the most accurate and clinically interpretable model for identifying individuals at high risk of T2DM.

Materials and Methods

Study Design

This study was conducted using analytical methods.

Population and Sample of the Study

The study population consisted of participants who applied to a hospital in Kütahya, Türkiye, from May 20, 2025, to July 20, 2025. There were two groups in this study. The first group consisted of those diagnosed with T2DM, and the second group consisted of those without T2DM. All data were collected from patients who were followed in the internal medicine outpatient clinic. Both groups consisted of individuals aged 40–64 who were able to communicate verbally. Individuals with any chronic disease other than T2DM, mental disability, or terminal-stage condition were excluded from the study.

A power analysis was performed to determine the sample size for the study. As a result of this analysis, it was planned to include at least 478 patients/group in the study, with a 99% confidence level and a 5% margin of error. Accordingly, the sample consisted of 585 participants with T2DM and 553 participants without T2DM. Simple random sampling was used as the sampling method.

Research Questions

The research questions are as follows:

- What are the risk factors for T2DM?
- What is the level of impact of these risk factors on T2DM?
- Which ML models can be used for T2DM?
- Which evaluation parameters can be used for these models?
- What interventions can be recommended for the prevention of T2DM?

Data Collection

A personal information form was used to collect data for the study. This form consists of two sections and 19 questions. The first section contains five questions about patients' individual characteristics, whereas the second section contains 14 questions about their health-related characteristics.

The study's data collection phase was conducted through face-to-face interviews with the participants. Following an explanation of the study, participants were presented with an informed consent form, and their consent was obtained. Each interview and form completion process took an average of 30 min. Data were collected directly by the researchers through the forms and recorded with due consideration for participant privacy.

Data Assessment

For the analysis of the study, the Statistical Package for the Social Sciences (SPSS) 25 and Google Colaboratory were used. In the SPSS 25 program, the Kolmogorov–Smirnov and Shapiro–Wilk tests were applied to determine whether the continuous variables were normally distributed. Given that all continuous variables in the investigation showed a normal distribution, parametric tests were used.

Subsequently, regression analyses were performed using the ML approach in the Google Colaboratory program with Python software language. In this analysis, K-Nearest Neighbors (KNN), support vector machine (SVM), decision

tree, random forest (RF), artificial neural network (ANN), naive bayes, and logistic regression (LR) models were used. These models were evaluated using the area under the receiver operating characteristic curve (area under the curve [AUC]), accuracy (ACC), F1 score, precision, and recall parameters.

Before applying the ML algorithms, a correlation analysis was conducted to assess potential multicollinearity among the independent variables. No strong correlations ($r > 0.80$) were detected; therefore, all variables were retained in the analysis. In addition, the recursive feature elimination (RFE) method was applied to confirm the relevance of predictors and to enhance the robustness and interpretability of the models. This approach ensured that only independent and meaningful variables were included in the final models, improving overall performance. This methodological approach aligns with current recommendations in the literature, emphasizing proper management of multicollinearity and the use of RFE to optimize variable selection in classification models.^[11,12]

To evaluate model interpretability, feature importance values from the RF model were examined to determine the influence of each predictor variable on classification performance. For the LR model, odds ratios (OR) and their 95% confidence intervals (CI) were calculated to assess the magnitude and direction of associations between predictors and T2DM risk. This complementary approach allowed for comparison of variable importance across models and enhanced clinical interpretability.

Software Programs

SPSS 25 is currently owned by IBM Corporation, headquartered in Armonk, New York, USA. In contrast, Google Colaboratory is a cloud-based platform developed by Google LLC, with its corporate headquarters in Mountain View, California, USA. Since its introduction in 2017, Colaboratory has gained widespread adoption, particularly in AI and ML, owing to its accessibility, collaborative features, and seamless integration with widely used programming environments.

Declaration of AI-Assisted Technologies

The authors declare that no AI-assisted technologies (such as large language models, chatbots, or image generators) were used in the preparation, writing, or editing of this manuscript. All content has been developed entirely by the authors.

Table 1. Individual characteristics of participants (n=1138)

Variables	n	%
Gender		
Female	578	50.8
Male	560	49.2
Education status		
Illiterate	97	8.5
Primary education	584	51.3
Secondary education	291	25.6
Higher education	166	14.6
Marital status		
Married	983	86.4
Single	155	13.6
Economic status		
Good	289	25.4
Bad	849	74.6
DM presence		
Yes	585	51.4
No	553	48.6
Age		
Mean±SD	54.29±9.79	
Minimum-Maximum	40–64	

Categorical variables were summarized as frequency and percentage (n, %), while continuous variables were presented as mean±standard deviation (Mean±SD) and range (minimum–maximum). DM: Diabetes mellitus.

Ethics of Study

This study was approved by the Kütahya Health Sciences University Ethics Committee (Date: 06.05.2025, Decision no: 2025/06-40) and the participating hospital before the study. The study was conducted in accordance with the ethical principles of the Declaration of Helsinki.

Results

Individual Characteristics of the Participants

The study included a total of 1138 patients (585 with T2DM and 553 without diabetes). Table 1 presents the individual characteristics of the participants. Accordingly, 50.8% of the participants were female, 51.3% had completed primary education, 86.4% were married, 74.6% had poor economic status, 51.4% had a diagnosis of T2DM, and the average age was 54.29±9.79.

Graph of OR for Risk Factors

The ORs for these variables were calculated to assess the effect of risk factors on the presence of T2DM (Fig. 1). These factors were ranked in order of highest effect size as follows:

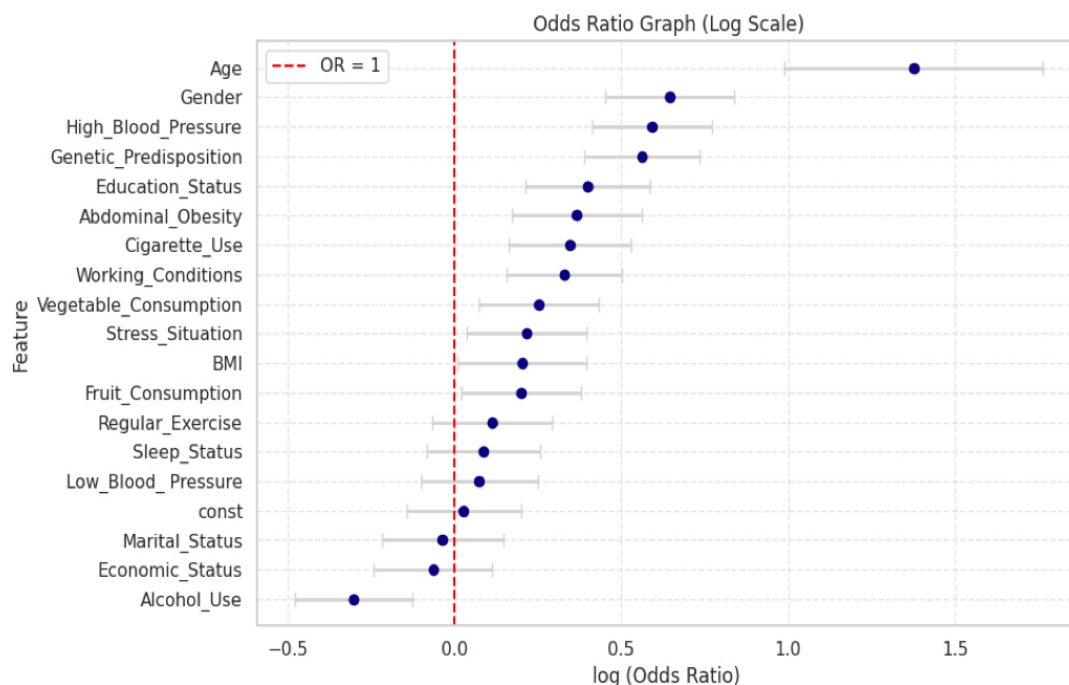


Figure 1. Graph of odds ratios for risk factors. Horizontal bars represent 95% confidence intervals derived from logistic regression analysis. The vertical dashed line (odds ratios=1) indicates the null value. Variables positioned to the right of the line indicate increased risk of Type 2 diabetes mellitus, whereas those to the left suggest a protective effect.

age, gender, high blood pressure, genetic predisposition, education status, abdominal obesity, cigarette use, working conditions, vegetable consumption, stress situation, body mass index (BMI), fruit consumption, regular exercise, sleep status, low blood pressure, marital status, economic status, and alcohol use.

Age was identified as the most influential variable (OR=3.96, 95% CI: 2.69–5.83, $p < 0.001$), followed by gender (OR=1.91, 95% CI: 1.57–2.31, $p < 0.001$), high blood pressure (OR=1.81, 95% CI: 1.51–2.16, $p < 0.001$), and genetic predisposition (OR=1.75, 95% CI: 1.47–2.08, $p < 0.001$). Educational status (OR=1.49, 95% CI: 1.23–1.79, $p < 0.001$), abdominal obesity (OR=1.44, 95% CI: 1.19–1.75, $p < 0.001$), cigarette use (OR=1.41, 95% CI: 1.18–1.70, $p < 0.001$), and working conditions (OR=1.39, 95% CI: 1.17–1.65, $p < 0.001$) also showed significant associations.

Lower vegetable consumption (OR=1.29, 95% CI: 1.07–1.54, $p = 0.006$), higher stress levels (OR=1.24, 95% CI: 1.04–1.48, $p = 0.018$), and higher BMI (OR=1.22, 95% CI: 1.01–1.48, $p = 0.041$) were additional contributors. Fruit consumption (OR=1.22, 95% CI: 1.02–1.46, $p = 0.031$) had a smaller but statistically significant effect.

On the other hand, regular exercise (OR=1.12, 95% CI: 0.94–1.34, $p = 0.215$), sleep status (OR=1.09, 95% CI: 0.92–1.29, $p = 0.324$), low blood pressure (OR=1.08, 95% CI: 0.90–1.28, $p = 0.409$), marital status (OR=0.96, 95% CI: 0.80–1.16, $p = 0.696$),

and economic status (OR=0.94, 95% CI: 0.79–1.12, $p = 0.476$) were not statistically significant predictors. Interestingly, alcohol use demonstrated a protective association with T2DM (OR=0.74, 95% CI: 0.62–0.88, $p < 0.001$).

Comparison of F1 Score, ACC, Precision, and Recall Across the Models

Figure 2 displays the F1 score, ACC, precision, and recall values for each model. Taking all parameters into account, the SVM, RF, and LR models yield the best results.

Comparison of AUC Across Models

Figure 3 displays the AUC values for each model. Considering all parameters, the SVM, RF, LR, and ANN models yielded the best results, respectively.

Discussion

This study comprehensively analyzed the risk factors for T2DM using an ML approach, and the results revealed that the ML models and evaluation parameters used demonstrated high performance and reliability. These findings confirm that ML-based risk prediction models can serve as valuable tools for early identification of individuals at risk for T2DM and can support data-driven clinical decision-making in preventive healthcare. The integration of such models into routine nursing and clinical practice

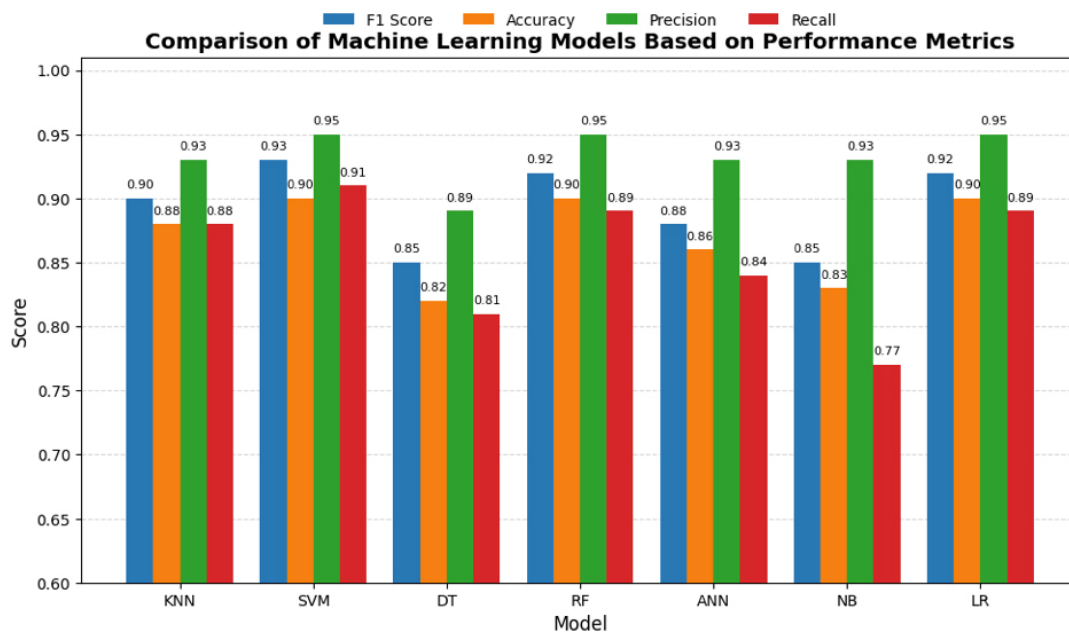


Figure 2. Comparison of machine learning models based on F1 score, accuracy, precision, and recall. Bar chart illustrating the performance comparison of seven machine learning models (K-Nearest Neighbors, support vector machine, decision tree, random forest, artificial neural network, naive bayes, and logistic regression) used to identify risk factors for Type 2 diabetes mellitus. The chart displays four key evaluation metrics: F1 score, accuracy, precision, and recall. Higher values indicate better model performance.

could help optimize screening strategies, particularly in populations with limited access to healthcare services.

Risk Factors

Risk factors for T2DM can be classified into two groups: modifiable and non-modifiable. Analysis of our study results demonstrates the influence of risk factors in both groups. Accordingly, the non-modifiable risk factors with the highest level of influence are age, gender, and genetic predisposition, respectively. Among modifiable risk factors, the factors with the highest to lowest impact levels are as follows: High blood pressure, educational status, abdominal obesity, cigarette use, working conditions, vegetable consumption, stress situation, BMI, fruit consumption, regular exercise, sleep status, and low blood pressure.

Of all these risk factors, age has the highest impact level. In this regard, Fazeli et al.^[13] emphasized that the risk of developing T2DM increases after the age of 40. The gender factor comes second. Kautzky-Wilker et al.^[14] reported that the prevalence of T2DM is 17.7 million higher in men than in women worldwide. High blood pressure is the third factor. Hezam et al.^[15] stated that individuals with high blood pressure have an increased risk of T2DM. Genetic predisposition is the fourth factor. Bonnefond et al.^[16] found that genetic predisposition plays an important role in the emergence of T2DM. Education status is the fifth factor. Yan et al.^[17] indicated that individuals with a higher

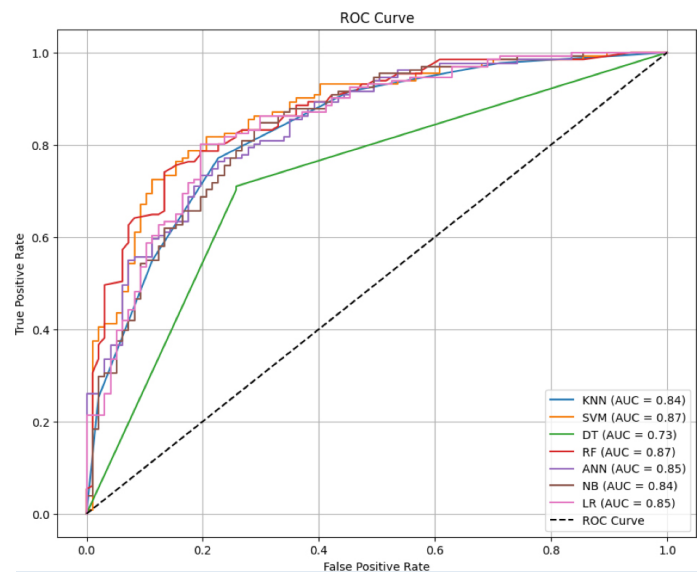


Figure 3. Comparison of receiver operating characteristic (ROC) curves and area under the curve (AUC) values for machine learning models. ROC curves illustrate the classification performance of seven machine learning models (K-Nearest Neighbors, support vector machine, decision tree, random forest, artificial neural network, naive bayes, and logistic regression) for predicting Type 2 diabetes mellitus. The diagonal dashed line represents the reference line (AUC=0.5), indicating no discrimination. Curves farther from the diagonal and closer to the upper-left corner demonstrate better model performance.

level of education have greater awareness of T2DM. The findings of the present study are consistent with previous research on this topic.

In addition to confirming known risk factors, this study contributes to the literature by quantifying their relative importance in a Turkish population, providing evidence that sociocultural and lifestyle factors may modify the strength of these associations. From a nursing and public health perspective, modifiable factors such as obesity, smoking, and sedentary behavior should be targeted through structured education, counseling, and preventive intervention programs.

Precision, Recall and F1 Score

Precision is a measure that indicates how many of the examples a model predicts as positive are actually positive. In a model, precision is critical when false positives can lead to serious consequences. Recall indicates how many of the models correctly predict as positive. Recall is more important when false negatives are likely to cause problems. The F1 score is a performance metric obtained by taking the harmonic mean of precision and recall values. It provides a balanced evaluation in situations where both false positives and false negatives are important. In imbalanced datasets, ACC can be misleading, so the F1 score provides a more realistic measure of success.^[11,12]

Since it is important to detect both false positives and false negatives in T2DM, the F1 score is particularly relevant. In our study, the models with the highest F1 scores were SVM, RF, and LR, with F1 scores of 0.93, 0.92, and 0.92, respectively. In studies conducted on this subject, the model with the highest F1 score was the RF model.^[18,19] The results obtained in our study are consistent with the literature. The F1 score is considered an important metric for evaluating the balance between the model's sensitivity and specificity, especially in medical datasets with imbalanced class distributions.^[20]

This result demonstrates that the performance metrics obtained are clinically meaningful rather than merely statistical. In medical contexts, a balanced F1 score indicates that the model can minimize both false alarms (unnecessary testing) and missed diagnoses, thus optimizing patient safety and resource utilization.

ACC

ACC is a basic performance metric that expresses the ratio of correctly predicted examples to the total number of examples. In clinical applications, ACC is a critical indicator of the model's overall classification success. However, it can be misleading, especially in health data with imbalanced class distributions.^[11,12]

This study found that the best ACC values among the ML models used to identify risk factors for T2DM were

obtained with the SVM, RF, and LR models, each achieving an ACC value of 90%. This result indicates that these models perform well in terms of overall ACC.^[21] In previous studies, Bhat et al.^[22] and Laila et al.^[23] used RF models, Talukder and Hossain^[24] employed LR, and Negi and Jaiswal^[25] used SVM in their analyses. Pradhan et al.^[26] utilized ANN, Alpan and İlgi^[27] applied KNN, and Islam et al.^[12] employed gradient boosting machine models, all reporting high ACC performance. According to our findings, while some results are consistent with previous studies, others differ, which may be attributed to the sociocultural characteristics of the populations studied. The 90% ACC rate obtained in this research demonstrates that ML methods are effective and reliable tools for predicting T2DM risk. These findings may contribute significantly to the development of clinical decision support systems and early intervention strategies.

The consistency of these results across multiple ML algorithms reinforces the reliability of the predictive relationships identified in this study. In addition, using face-to-face collected clinical data rather than public datasets strengthens the validity of the ACC estimates and supports the model's applicability in real-world healthcare environments.

AUC

The AUC is a graphical method that evaluates the balance between sensitivity and specificity under different threshold values for a classification model. The area under this curve is a powerful performance metric that summarizes the model's ability to distinguish between positive and negative classes. As the AUC value approaches 0.5, the model's discriminative power decreases, while as it approaches 1.0, it indicates a strong discriminative model. In clinical applications, the AUC plays a critical role in model selection, especially when the costs of false positives and false negatives differ. In T2DM diagnosis, a false negative result can lead to the disease being overlooked, whereas a false positive result can cause unnecessary treatment and increased patient anxiety.^[11,12]

This study found an AUC value of 0.87 for the SVM and RF models and 0.85 for the LR and ANN models. These results indicate that the models generally have high discriminatory power. In the literature, Kaur and Kumari^[28] reported an AUC of 0.90 for the SVM model, Islam et al.^[29] found an AUC of 0.60 for the Bagged Chart model, and Kopitar et al.^[30] reported an AUC of 0.85 for the GLMNET model. Our findings are consistent with those of Kaur and Kumari^[28] but differ from other studies, which may be attributed to the use of preprocessed or ready-made datasets in those analyses.

The AUC value provides a more balanced and meaningful assessment than one-dimensional criteria such as ACC, especially in medical datasets where class distribution may be imbalanced. Therefore, AUC is considered a reliable and effective tool for the development of clinical decision support systems.

The high AUC values obtained in this study confirm that the models possess strong discriminatory ability even with clinically collected data, supporting their potential for integration into hospital information systems and risk-based screening programs. Moreover, combining AUC with other evaluation metrics such as F1 and ACC enables a more balanced and multidimensional understanding of model performance.

Overall, the combination of high F1, ACC, and AUC values demonstrates that ML algorithms, particularly SVM, RF, and LR, can be effectively applied for accurate T2DM risk prediction and may serve as robust components of clinical decision support systems.

Strengths and Limitations

The study's use of the ML approach for analysis is one of its major strengths. The use of AI and ML rather than traditional statistical methods in the study provides strong evidence for the clinical validity of the research. Furthermore, while data in the majority of similar studies in the literature are obtained from ready-made data sets, in this study, data were collected by researchers through face-to-face interviews. This method significantly enhances the validity and reliability of the data set. The study evaluated 18 risk factors for T2DM. No other study in the current literature examines this number of risk factors together. The research was conducted at a single center. Therefore, the sample was limited to patients in a specific region and time period. This restriction limits the generalizability of the results to different sociocultural groups.

Conclusion

This study identified age, gender, high blood pressure, genetic predisposition, and education level as the most influential risk factors for T2DM. Among the seven ML models applied, SVM, RF, and LR achieved the highest performance. SVM demonstrated strong generalization capability for complex and high-dimensional data, whereas RF combined high ACC with interpretability by providing insights into variable importance. LR remains a clinically preferred method due to its ease of interpretation and practical applicability. Integrating ML-based prediction

models into clinical workflows may help nurses and healthcare professionals detect at-risk individuals earlier, prioritize preventive counseling, and allocate resources more efficiently.

Model selection in clinical research should consider not only predictive performance but also interpretability and feasibility. Although SVM may be complex for clinical decision-making, RF offers advantages in identifying and understanding risk factors, whereas LR is more easily interpretable and applicable for healthcare professionals. Therefore, for accurate classification of T2DM risk factors and integration into clinical decision support systems, the use of RF or LR models is recommended. Future studies should focus on validating these models across different regions and populations and exploring hybrid or ensemble approaches that combine ML with traditional epidemiological methods to improve precision and clinical utility.

Ethics Committee Approval: This study was approved by the Kütahya Health Sciences University Ethics Committee (Date: 06.05.2025, Decision no: 2025/06-40).

Informed Consent: Written informed consent was obtained.

Conflict of Interest: None declared.

Financial Disclosure: The authors declared that this study has received no financial support.

Use of AI for Writing Assistance: None declared.

Authorship Contributions: Concept: SB, YK, İB, ŞÖ; Design: SB, YK, İB, ŞÖ; Supervision: SB, YK, İB; Resource: İB; Materials: ŞÖ; Data collection and/or processing: ŞÖ; Analysis and/or interpretation: SB; Literature review: SB; Writing: SB, YK, İB, ŞÖ; Critical review: SB, YK.

Peer-review: Double blind peer-reviewed.

References

1. Ardahanlı İ, Aslan R, Çelik M, Akgün O, Akyüz O. Effects of empagliflozin on carotid intima-media thickness and epicardial fat tissue volume in patients with type-2 diabetes mellitus. *Lokman Hekim Health Sci* 2021;1(3):74-80. [\[CrossRef\]](#)
2. Bilgehan T, Inkaya BV, Şendur EG. Türkiye's First Diabetes Nurse Coaches' Opinions on Diabetes Nurse Coaching: A Qualitative Study. *Lokman Hekim Health Sci* 2025;5(2):181-9. [\[CrossRef\]](#)
3. Shi F, Zhao Q, Yang Y, Liu L, Zhang X, Kim HJ, et al. Global burden of diabetes in women from 1990 to 2021, with projections to 2050: population-based study. *BMC Med* 2025;23(1):538. [\[CrossRef\]](#)
4. Liu C, Li Y, Wang N, Wu Y, Liu J, Ding M, et al. Trends and comparisons of diabetes burden in China and the world from 1990 to 2021, with forecasts to 2050: a systematic analysis of the Global Burden of Disease Study 2021. *Diabetol Metab Syndr* 2025;17(1):309. [\[CrossRef\]](#)

5. Turkish Statistical Institute. Türkiye Sağlık Araştırması 2022. Available at: [\[CrossRef\]](#)
6. Sami A, Javed A, Ozsahin DU, Ozsahin I, Muhammad K, Waheed Y. Genetics of diabetes and its complications: a comprehensive review. *Diabetol Metab Syndr* 2025;17(1):185. [\[CrossRef\]](#)
7. Zhou R, Li F, Chen G, Fu Q, Gu S, Wu X. Associations between general and abdominal obesity and incident diabetic neuropathy in participants with type 2 diabetes mellitus. *J Diabetes* 2021;13(1):33-42. [\[CrossRef\]](#)
8. Nouripour F, Mazloom Z, Fararouei M, Zamani A. Effect of protein and carbohydrate distribution among meals on quality of life, sleep quality, inflammation, and oxidative stress in patients with Type 2 diabetes: A single-blinded randomized controlled trial. *Food Sci Nutr* 2021;9(11):6176-85. [\[CrossRef\]](#)
9. Esteva A, Robicquet A, Ramsundar B, Kuleshov V, DePristo Mark, Chou K. A guide to deep learning in healthcare. *Nat Med* 2019;25(1):24-9. [\[CrossRef\]](#)
10. Kumar Y, Mahajan M. Recent advancement of machine learning and deep learning in the field of healthcare system. In: Srivastava R, Mallick PK, Rautaray SS, Pandey M, editors. *Computational Intelligence for Machine Learning and Healthcare Informatics*. Berlin, Germany: De Gruyter; 2020. p.77-98. [\[CrossRef\]](#)
11. Miao J, Zhu W. Precision-recall curve (PRC) classification trees. *Evol Intell* 2022;15(3):1545-69. [\[CrossRef\]](#)
12. Islam MR, Banik S, Rahman KN, Rahman MM. A comparative approach to alleviating the prevalence of diabetes mellitus using machine learning. *Comput Methods Programs Biomed Update* 2023;4(21):100113. [\[CrossRef\]](#)
13. Fazeli PK, Lee H, Steinhauer ML. Aging Is a powerful risk factor for type 2 diabetes mellitus independent of body mass index. *Gerontology* 2020;66(2):209-10. [\[CrossRef\]](#)
14. Kautzky-Willer A, Leutner M, Harreiter J. Sex differences in type 2 diabetes. *Diabetologia* 2023;66(6):986-1002. [\[CrossRef\]](#)
15. Hezam AAM, Shaghdar HBM, Chen L. The connection between hypertension and diabetes and their role in heart and kidney disease development. *J Res Med Sci* 2024;29(1):22. [\[CrossRef\]](#)
16. Bonnefond A, Florez JC, Loos RJ, Froguel P. Dissection of type 2 diabetes: a genetic perspective. *Lancet Diabetes Endocrinol* 2025;13(2):149-64. [\[CrossRef\]](#)
17. Yan Y, Wu T, Zhang M, Li C, Liu Q, Li F. Prevalence, awareness and control of type 2 diabetes mellitus and risk factors in Chinese elderly population. *BMC Public Health* 2022;22(1):1382. [\[CrossRef\]](#)
18. Phongying M, Hiriote S. Diabetes Classification Using Machine Learning Techniques. *Computation* 2023;11(5):96. [\[CrossRef\]](#)
19. Atif M, Anwer F, Talib F, Alam R, Masood F. Analysis of machine learning classifiers for predicting diabetes mellitus in the preliminary stage. *Int J Artif Intell* 2023;12(3):1302-11. [\[CrossRef\]](#)
20. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLoS One* 2015;10(3):e0118432. [\[CrossRef\]](#)
21. Chicco D, Jurman G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics* 2020;21(1):6. [\[CrossRef\]](#)
22. Bhat BS, Selvam V, Ansari GA, Ansari MA. Analysis of diabetes mellitus using machine learning techniques. In: 2022 5th International Conference on Multimedia, Signal Processing and Communication Technologies (IMPACT); 2022. p. 1-5. [\[CrossRef\]](#)
23. Laila UE, Mahboob K, Khan AW, Khan F, Taekeun W. An Ensemble Approach to Predict Early-Stage Diabetes Risk Using Machine Learning: An Empirical Study. *Sensors (Basel)* 2022;22(14):5247. [\[CrossRef\]](#)
24. Talukder A, Hossain MZ. Prevalence of diabetes mellitus and its associated factors in Bangladesh: Application of two-level logistic regression model. *Sci Rep* 2020;10(1):10237. [\[CrossRef\]](#)
25. Negi A, Jaiswal V. A first attempt to develop a diabetes prediction method based on different global datasets. In: 2016 4th International Conference on Parallel, Distributed and Grid Computing (PDGC). Piscataway (NJ): IEEE; 2016. p. 237-41. [\[CrossRef\]](#)
26. Pradhan N, Rani G, Dhaka VS, Poonia RC. Diabetes prediction using artificial neural network. In: Agarwal B, Balas VE, Jain LC, Poonia RC, Manisha, editors. *Deep Learning Techniques for Biomedical and Health Informatics*. Cambridge (MA): Academic Press; 2020. p. 327-39. [\[CrossRef\]](#)
27. Alpan A, Ilgi GS. Classification of diabetes dataset with data mining techniques by using WEKA approach. In: 2020 4th International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT); 2020. p. 1-7. [\[CrossRef\]](#)
28. Kaur H, Kumari V. Predictive modelling and analytics for diabetes using a machine learning approach. *Appl Comput Inform* 2020;18(1-2):90-100. [\[CrossRef\]](#)
29. Islam MM, Rahman MJ, Chandra Roy D, Maniruzzaman M. Automated detection and classification of diabetes disease based on Bangladesh demographic and health survey data, 2011 using machine learning approach. *Diabetes Metab Syndr* 2020;14(3):217-9. [\[CrossRef\]](#)
30. Kopitar L, Kocbek P, Cilar L, Sheikh A, Stiglic G. Early detection of type 2 diabetes mellitus using machine learning-based prediction models. *Sci Rep* 2020;10(1):11981. [\[CrossRef\]](#)