



IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 13 Issue: III Month of publication: March 2025 DOI: https://doi.org/10.22214/ijraset.2025.67905

www.ijraset.com

Call: 🕥 08813907089 🔰 E-mail ID: ijraset@gmail.com



Investigating Gender and Age Variability India Betes Prediction

Dr.P.EdithLinda¹, Dr.R.Srividhya², S.Vigneshwaran³ Dr.G.R. Damodaran College of Science, Coimbatore-641014

Abstract: Diabetesisawidespreadhealthconcern, affecting millionsglobally and posing significant risks of complications such as heart disease, hypertension, and kidney damage. Early detection and intervention are crucial for managing the disease and preventing severe health issues. This project aims to design a robust systemforanalyzingdiabetesrelateddataandpredictingdiabetesbasedonageandgenderclassification. The dataset used in this study includes key health features such as age, gender, BMI, blood glucose levels, and lifestyle factors like smoking history, providing a solid foundation for effective analysis. The current system leverages traditional machine learning models, including Random Forest, K-Nearest Neighbors (KNN), and Logistic Regression, to classify individuals based on the presence or absence of diabetes. These models have shown satisfactory performance, particularly for structured data, thanks to their simplicity and efficiency. By leveraging deep learning, the CNN model can identify subtle correlations between variables such as age, gender, BMI, and blood glucose levels, enhancing the overall predictive accuracy compared to traditional machinelearningtechniques. The Kaggle diabetes prediction datasets erves as the foundation for training and testing the models. Data preprocessing steps, including handling missing values, encoding categorical variables, and normalizing numeric features, are applied to ensure the dataset is prepared for analysis. The CNN model is designed with layers optimized for feature extraction and classification, using activation functions like ReLU and Softmax to ensure accurate predictions. The models are evaluated based on metrics such as confusion matrices, accuracy, precision, recall, and F1-score, providing a comprehensive view of their

performance. This project highlights the potential of combining advanced machinelearning and deeplearning

techniquesforimprovingdiabetesprediction.BycomparingtraditionalmethodswiththeproposedCNNmodel, the study aims to provide better diagnostic accuracy and contribute to the field of healthcare analytics. The resultsofthisresearch couldhelphealthcare providersmakemoreinformed, data- drivendecisions,facilitating earlydetectionandmanagementofdiabetes, andultimately supporting thedevelopmentofscalable, impactful predictive systems for healthcare.

I. INTRODUCTION

Diabetes mellitus, particularly Type 2 diabetes, has become a significant global health concern, affecting millions of individuals across various age groups and demographics. The growing prevalence of diabetes, driven by factors such as sedentary lifestyles, poor diet, and increasing obesity rates, emphasizes the need for early detection and effective prediction models to identify those at risk.

Age and gender are two demographic factors that are consistently associated with the risk of developing diabetes, making them crucial variables to consider inpredictive models. Age is a well-established risk factor, with the likelihood of developing Type 2 diabetes increasing as individuals age, especially beyond 45 years. However, diabetes does not discriminate based solely on age, as the condition can develop at any point in life, influenced by other variables such as genetics, lifestyle, and environment.

Gender, too, plays an important role in the development of diabetes. Research suggests that while men may developdiabetes atyoungerages, womenare atheightenedrisk post-menopause,potentially duetohormonal changesandalteredfatdistributionpatterns. Additionally,womenwithahistory of gestational diabetes face a greater risk of developing Type 2 diabetes later in life.

Despitetheknownassociationsbetweenage,gender,anddiabetes,thevariabilityinhowthesefactorsinfluencetheriskacrossdifferentpopulati onsre mainsanareaofactiveresearch. Understandinghowageandgendercontributetothepredictionofdiabetes iscriticalindesigning personalizedhealthinterventionsandimproving earlydetectionsystems. Machinelearningandstatisticalmethods offerpowerful toolstoanalyze largedatasets and identify patterns in how age and gender influence diabetes onset.

This study aims to explore the relationship between age, gender, and diabetes risk by utilizing a dataset containing relevant health metrics, such as blood sugar levels, BMI, and demographic data. Through the applicationofvarious predictive modeling techniques, we seek to gain insights into how these variables interact and affect diabetes prediction.



International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue III Mar 2025- Available at www.ijraset.com

By investigating these relationships, the study aims to enhance our understanding of diabetes risk and inform the development of more effective, tailored prediction models for different demographic groups.

Theexisting system utilizeswell-established machinelearning algorithmssuch asRandom Forest,KNN, and LogisticRegression to classify individualsbased on theirdiabetesstatus.Thesemodels serveasbenchmarks, offeringacomprehensiveevaluation of traditional techniques for predicting the presence or absence of diabetes as indicated by the dataset's "classfactor" column. While effective to an extent, these models have certain limitations in handling nonlinear relationships and complex patterns in the data.

To address these limitations, the proposed system introduces a Convolutional Neural Network (CNN) SequentialModel.Thedeeplearningapproachisdesignedtouncoverintricatepatternsandrelationshipswithin

the dataset that traditional methods might over look. The CNN's ability to model complex dependencies allows for improved prediction accuracy and robustness, making it avaluable addition to the existing methodologies.

Thisprojectnotonlyhighlightsthestrengthsandweaknessesoftraditionalanddeeplearningmethodsbutalso underscores the importance of innovative approaches in healthcare analytics. By offering a comparative analysisofmultiplepredictivemodels, theprojectprovidesacomprehensiveframeworkfordiabetesprediction and risk assessment, contributing to advancements in medical research and public health strategies. Another keyobjective is evaluate the performance of traditional machinelearning models such as Random Forest, KNN, and Logistic Regression for diabetes prediction. These models are benchmarked for their accuracy, efficiency, and ability to classify individuals based on the dataset's "classfactor" column, which indicates the presence or absence of diabetes. By assessing the strengths and limitations of these techniques, this project aims to highlight the effectiveness of established methods while identifying areas where improvements are needed.

The project also seeks to implement and demonstrate the advantages of a deep learning approach through a CNNSequentialModel.TheCNNisdesignedtoovercomethelimitationsoftraditionalmethodsbycapturing

complex,nonlinearpatternsinthedataset,thusenhancingthepredictiveaccuracyandreliabilityofthesystem. By comparing theresultsof theproposed deeplearning modelwith traditionalmodels,thisobjectiveseeksto establish a more robust and scalable solution for diabetes prediction, contributing to advancements in healthcare analytics and preventive medicine.

II. METHODOLOGY

The methodology for this diabetes prediction system is structured around the design, development, and evaluation of two types of models: traditional machine learning models and a deep learning-based ConvolutionalNeuralNetwork(CNN)model. Themethodologyforthisstudyfollowsasystematicapproach to data collection, preprocessing, model selection, training, and evaluation to ensure robust diabetes prediction based on age and gender classification.

The key stages include dataset acquisition preprocessing , machine learning model implementation, deep learningmodelandevaluation. This approachaims to compare the performance of the semodels in predicting diabetes, using a comprehensive dataset from Kaggle that includes key health features such as age, gender, BMI, blood glucose levels, and lifestyle factors. The methodology follows these primary steps:

- DataPreprocessingModule
- MachineLearningModelsModule
- DeepLearning(CNN)ModelModule
- ModelEvaluationandVisualizationModule

1) Data Preprocessing Module

Thedatapreprocessing module is responsible for preparing the dataset for analysis. This module involves several key steps such as handling missing values, encoding categorical variables, and normalizing or scaling numeric features.

In the case of the diabetes dataset, this module will process features like age, gender, BMI, blood glucose levels, and lifestyle factors.

Additionally, outliers and anomalies will be detected and handled appropriately to ensure the data quality and prevent any biases in the model training process.

2) Machine Learning Models Module

This module involves the implementation and training of traditional machine learning models like Logistic Regression, K-Nearest Neighbors (KNN), and Random Forest Classifier. These models will be trained on the preprocessed dataset to predict the presence or absence of diabetes.



Themodulealsoincludesperformanceevaluation of these models based on accuracy, precision, recall, F1-score, and confusion matrices, which will help in assessing their effectiveness in the classification task.

3) Deep Learning(CNN)Mode lModule

ThedeeplearningmodulefocusesonimplementingtheConvolutionalNeuralNetwork(CNN)usingaSequential Model.The CNN is designed to learn complex patterns in the data by utilizing layers like convolutional layers, activation functions (ReLU), and fully connected layers.

This module will also involve tuning, training the model, and evaluating its performance against traditional machinelearningmodels. The goal of this module is to provide a more accurate and robust approach to diabetes prediction based on the given dataset.

4) Model Evaluation and Visualization Module

Once the models are trained, this module handles the evaluation of model performance through various metrics suchasaccuracy, precision, recall, F1-score, and AUC. This module will also visualize the results using confusion

matrices,ROCcurves,andotherperformancemetricstcomparethetraditionalmachinelearningmodelswiththe CNN model.

Additionally, it will generate visualizations for model predictions and real-world accuracy, providing an overview of how well the system can predict diabetes based on age, gender, and other factors.



Fig1:ProposedSystemWork

III. RESULTS AND EVALULATION

The evaluation of the diabetes prediction models was carried out using several metrics, including accuracy, precision, recall, F1score, and confusion matrix, to assess the effectiveness of both traditional machine learning models and the proposed Convolutional Neural Network (CNN) model.

A. Traditional Machine Learning Models

The traditional machine learning models used for comparison included Random Forest, K-Nearest Neighbors (KNN), and Logistic Regression. These models showed satisfactory performance, with Random Forestachieving the highest accuracy of approximately 80%, followed by KNN at 75%, and Logistic Regression at 73%. The models performed well in handling structured data, providing solid baselines for comparison. However, they struggled with capturing complex, non-linear relationships between features, particularly in the presence of high-dimensional data and subtle interactions between variables such as age, BMI, and blood glucose levels.





Fig2:Gradient Boost

B. Convolutional Neural Network(CNN)

The proposed CNN model was designed to overcome the limitations of traditional machine learning models by leveraging deep learning techniques to automatically extract high-level features from the data. This model outperformed the traditional models, achieving an accuracy of approximately 85%. The CNN's ability to detect complex patterns in the data, such as non-linear relationships between the features, resulted in a notable improvement in predictive performance. The CNN also exhibited improved performance on the precision and recall metrics, which are crucial for healthcare applications where false positives and false negatives can have significant consequences.

C. Evaluation Metrics

Accuracy: TheCNNmodelachieved thehighestaccuracy, outperforming the traditional models by asignificant margin.

Precision:Precisionmeasurestheproportionoftruepositivepredictionsamongallpositivepredictions.TheCNN model showed better precision, indicating fewer false positives compared to traditional models.

Recall: Recall, or sensitivity, refers to the proportion of actual positive cases correctly identified by the model. The CNN model also demonstrated higher recall, ensuring that more diabetic cases were correctly predicted, which is critical for early intervention.

F1-Score:TheF1-score,whichbalancesprecisionandrecall,washighestfortheCNNmodel,furtherconfirming its superior performance in predicting diabetes.

=92.1%





D. Confusion Matrix

The Confusion matrix for the CNN model displayed a higher number of correct predictions (both true positives and true negatives), suggesting that the model had better overall classification performance. The traditional models, while effective, had a higher number of false negatives and false positives, indicating that their predictions were not as reliable in detecting diabetes.



Fig4:ConfusionMatrix-KNN

Classification Report: proclaim modil (2-nour support) 8 8.8 1.98 8.98 1.074 1 8.97 8.98 8.77 1.094 Rocza eq 8.97 8.89 8.78 2.994 Rocza eq 8.97 8.89 8.98 2.994 Rolation Roundy: 19.05 Consent 8.09 Roclaim Report: proclaim modil. Consen support 10 Joint 8.97 8.99 8.98 2783 Debrin 8.98 8.45 8.7 2.997 Roczary 8.98 8.45 8.7 2.997	E conteners	energyeda T							
precision seculi 67-score support 8 8.8 1.08 8.80 12754 1 8.97 8.88 8.70 1284 scores op 8.87 8.88 1289 score op 8.87 8.88 8.89 1289 skipteni op 8.88 8.88 8.89 1289 skipteni op 8.88 8.89 1289 skipteni fogeresien fotoloo: soarsep 8.89 score op 8.70 Societien food precision 1.80 score op 8.70 Societien food Societien 8.87 8.99 8.88 27483 Dideten 8.87 8.99 8.88 27483 Dideten 8.88 8.42 8.70 2547 score op 8.58 1.88 8.80 3889	Care (Filestia	e Nepert:							
8 8.86 1.80 8.20 1 8.77 8.80 1.70 mccrs eg 8.78 1.200 mccrs eg 8.77 8.80 1.200 mccrs eg 8.77 8.80 1.200 mccrs eg 8.77 8.80 1.200 McCatation 6.85 1.2000 1.00 McCatation 6.85 1.2000 1.00 McCatation 6.85 1.2000 1.00 McCatation 6.85 1.80 1.00 Application McCatation McCatation 1.00 Application McCatation McCatation McCatation Mocaline Mc		precision	PROFE	Grown	support)				
1 8.57 8.99 8.74 209 score og 8.57 8.58 8.58 2099 skiptel og 8.58 8.56 8.58 2099 Nisteter kourup: 8.40 //hereildetidetidgenderhelitierspilter prjil.gr aptick legrender Motion toarsp: 6.90 kodi: 8.63 2.5melfautter Repet: pretider moli. f2-som segert jentider moli. f2-som segert pretider 8.51 8.99 8.90 2763 Didetes 8.57 8.99 8.90 2763 Didetes 8.57 8.99 8.90 2763		0.96	1.00	1.00	1754				
ACCENCY B. B. 300 MCT AN ALL ST B. B. B. 2000 MCT AN ALL ST B. B. B. 2000 MCDATE Reports (S. AN Mean Observation Matrice covery B. AN Mean Matrice Reports persisten Provid Process separt In Outpets B. B. B. B. B. 2000 Debates B. B. B. B. B. 2000 Debates B. B. B. B. B. 2000 ACCENCY B. B. B. B. B. 2000 ACCENCY B. B. B. B. 2000 ACCENCY B. B. B. B. 2000 ACCENCY B. B. B. B. B. 2000 ACCENCY B. B. B. 2000 ACCENCY B. B. B. 2000 ACCENCY B. B. B. 2000 ACCENCY B. B. B. 2000 ACCENCY B. B. 2000 ACCENCY B. B. B		0.07	1.50	1.74	1284				
macro seg 0.27 0.08 0.08 0.000 kijded seg 0.38 0.09 0.00 Vekeelidededapleeleelivelietsergeber pejil gy ajats hyveelee Mation coarsy 0.00 beiden	00027801			0.96	23990				
edgited og 8.56 8.56 8.56 2000 Militetier kourupy: 96.63 Agistic Agenesiae fotolos: kourup 2.600 Mediae 8.60 Mediae 8.60 percisien 1.60 mediae 8.60 percisien media. Foreare aggert percisien media. Foreare aggert percisien media. Foreare aggert southout 8.60 8.61 8.00 8.70 Didette 8.61 8.00 8.00 2700 Didette 8.63 8.64 8.72 2507 accuroy 8.58 8.65 8.50 3000	NOTE INC.	0.07	10.50	0.00	23800				
Nüstrine Kosmoy: M.KM //form/Bishetishekeelevelistisenyethen pepill.gy agistis Negrossian Matrico: Nasaline KM2 Mealine KM2 Mealine KM2 perclalar proli.fprozer seggent perclalar proli.fprozer seggent Sasalinetten K.M1 K.M9 K.M 27453 Diabetes K.M1 K.M9 K.M 27453 Diabetes K.M1 K.M 8.K5 2000	elijittel eg	0.96	10.005	0.96	1000				
//konvillubrichydianiadfweliafianygthe pepill gy agistis fagensiae fetrios: soarog 8 (18) median 1.80 median 1.80 Soare 8.70 Soarfiotiae peoil. Firsten segent peololen 8.87 8.99 8.98 2783 Dahrtes 8.87 8.99 8.98 2783 Dahrtes 8.87 8.99 8.98 2783 Dahrtes 8.88 8.45 8.72 2547 socarog 8.88 8.45 8.72 2547	kilaintan k	oang: 16.6	3						
1 Sover 8.758 Sensification Report: particles recall firenex suggert ReSolveton 8.87 8.99 8.98 77403 Solveton 8.88 8.62 8.70 2947 accessoy 8.98 8.65 30600 macro og 8.51 8.89 8.85 30600 macro og 8.51 8.89 8.85 30600	opietic App course; E.S becision: E. becision: E.	ressian Mrasi 69 80 1	n y Blan Lan Se	n del fronte de	der H				
Canadriantian Aparti precision recall Grocen separt In Didetes 8.97 8.99 8.90 2743 Didetes 8.38 8.61 8.72 2947 accuracy 8.58 3.89 micro eq 8.55 8.58 3.899 micro eq 8.55 8.58 3.899	Ci Sonne II.1	20							
precision modil 1/jointe suggert In Diadetes 8.47 8.58 8.08 27463 Diadetes 8.38 8.62 8.70 2847 eccessoy 8.58 3.889 micro eq 8.51 8.58 3.889 micro eq 8.51 8.58 3.889	Genericetta	n linguart							
In Dialwhon 8.57 8.59 8.58 27453 Dialwhon 8.58 8.45 8.77 2547 accuracy 8.58 3880 micro ang 8.55 8.58 3880 micro ang 8.55 8.58 3880		bacroom.	PROFES.	11-10-10-1	President.				
Dialwin B.B B.C B.T 257 scorecy B.B 3000 score of B.S. B.B 8.D 3000 score of B.S B.D 8.D 3000	In Distance	0.07	1.00	0.00	17483				
essney 8.8 389 Norman 8.5 8.9 8.5 389 Kiptel ang 8.5 8.5 8.5 3894	Diabettes	0.36	1.61	0.72	2547				
Marrieg I.S. I.H I.S. 2009 staffed eg I.S. I.N I.S. 2009	1002780			1.8	30000				
elpholog 8.96 8.96 8.96 3000	Micro and	0.03	18.39	11.15	10000				
	ecipities) and	0.96	11/26	0.56	3000				

Fig5:ClassificationMetrics

IV. CONCLUSION

The Diabetes Data Analysis and Prediction project successfully leverages both traditional machine learning algorithms and deep learning techniques to predict the likelihood of diabetes based on age, gender, and other relevant health features. Theuse of models like Logistic Regression, K-Nearest Neighbors (KNN), and Random Forest proved effective in classifying diabetes, but the proposed deep learning approach using a Convolutional Neural Network (CNN) provided a superior method for capturing complex, nonlinear patterns within the data. The comparison between these models highlighted theadvantages of deep learning in improving prediction accuracy,offeringamorerobustsolutionfortheearlydetectionofdiabetes.

Byintegratingtheseadvanced techniques, the project emphasizes the potential of machine learning and deep learning in the health caredomain, particularly for diabetes management. The results underscore the importance of data-driven decision-making in health care, enabling more accurate, scalable, and efficient diagnostic systems. This research provides valuable insights into the potential of AI in health care, paving the way for further improvements and innovations in the field.



International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue III Mar 2025- Available at www.ijraset.com

REFERENCES

- [1] Breiman, L. (2001).Random forests. Machine Learning, 45(1), 5-32. <u>https://doi.org/10.1023/A: 1010933404324</u>
- [2] Cover, T., & Hart, P. (1967). Nearestneighborpatternclassification. IEEE Transactionson Information Theory, 13(1), 21-27. <u>https://doi.org</u> /10.1109/TIT.1967.1053964
- [3] Hosmer, D.W., Lemeshow, S., & Sturdivant, R.X. (2013). Applied logistic regression (3rded.). Wiley.
- [4] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deeplearning. Nature, 521 (7553), 436-444. https://doi.org/10.1038/nature14539
- [5] Kingma, D.P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980. <u>https://arxiv.org/abs/1412.6980</u>
- [6] Kaggle.(n.d.).Diabetespredictiondataset.Retrievedfrom<u>https://www.kaggle.com/</u>
- [7] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. Proceedings of theIEEE Conference on Computer Vision and Pattern Recognition (CVPR), 770-778. https://doi.org/10.1109/CVPR.2016.90
- [8] Choi, E.,Bahadori,M.T.,Schuetz, A.,Stewart, W. F.,& Sun, J.(2016).Retain: Aninterpretable predictive model for healthcare using reverse time attention mechanism. Advances in Neural Information Processing Systems (NeurIPS), 3504-3512.
- [9] Rahman,M.M., & Muniyandi, R.C. (2020). Performance analysis of machine learning techniques indiabetes prediction. IEEE Access, 8, 115512-115524. https://doi.org/10.1109/ACCESS.2020.3003731
- [10] Abhari, S., NiakanKalhori, S. R., & Ebrahimi, M. (2019). A comparison of machine learning models for diagnosing
 diabetes.
 Health

 Informatics
 Journal,
 25(3),
 984-1000. https://doi.org/10.1177/1460458218796633
 Health
- [11] Jothi, N., Rashid, N. A., &Yunus, Y. (2015). Data mining in healthcare- A review. Procedia ComputerScience, 72, 306-313. https://doi.org/10.1016/j.procs.2015.12.142
- [12] Zhang,Z.,Li,X.,&Zhang,H.(2021).AnimprovedCNNmodelfordiabetesdiagnosis. Computationaland Mathematical Methods in Medicine, 2021, 1-12. https://doi.org/10.1155/2021/5540436
- [13] Soni,J.,Ansari,U.,Sharma,D.,&Soni,S.(2011).Predictivedataminingformedicaldiagnosis: Anoverview of heart disease prediction. International Journal of Computer Applications, 17(8), 43-48.
- [14] Farooq, M.U., & Shaukat, F. (2022). Deeplearning-based prediction of diabetes mellitus. Journal of Biomedical Informatics, 128, 103988. https://doi.org/10.1016/j.jbi.2022.103988
- [15] Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deeplearning. MITPress.











45.98



IMPACT FACTOR: 7.129







INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 🕓 (24*7 Support on Whatsapp)