

1 Diabetes Identification Using Machine 2 Learning(ML)

3
4
5
6 Vika Li^a, Yixin(Crystal) Luo^a, and Vaishnavi Singh^a
7
8
9
10
11 Diabetes | Machine Learning | Big Data Analysis |
12

13 1. ABSTRACTION 14

15 This study addresses the critical challenge of diabetes prediction by leveraging
16 machine learning models to analyze behavioral and health-related data. Our research
17 utilizes a comprehensive dataset from the 2014 Behavioral Risk Factor Surveillance
18 System (BRFSS), which includes 436,668 records post-oversampling to ensure class
19 balance. We implement and compare five machine learning algorithms—Logistic
20 Regression, Decision Tree, K-Nearest Neighbors (KNN), Gaussian Naive Bayes,
21 and Multi-layer Perceptron (MLP)—to identify the most effective model based on
22 the top five features determined through regression analysis. The MLP emerges as
23 the superior model, demonstrating the highest accuracy and balanced performance
24 across various metrics, including precision, recall, F1-score, and an ROC-AUC
25 score of 0.7286. Our findings suggest that integrated computational approaches can
26 significantly enhance the early detection and management of diabetes.
27
28
29
30

2. INTRODUCTION

31 Diabetes is one of the most common endocrine diseases, with more than 133
32 million Americans currently living with diabetes or in the pre-diabetic stage (CDC,
33 2022). While there are many genetic factors that can influence one's diagnosis of
34 diabetes, the surge in diabetes cases is inextricably linked to contemporary human
35 lifestyles (Hu, 2021).

36 Excessive calorie intake and lack of physical activity, coupled with smoking and
37 alcohol abuse, have largely contributed to the rising incidence of this metabolic
38 disease (Hu, 2021; Magkos et al., 2009). In other words, with proper health
39 management, diabetes can be controlled in an ideal state (Ganie et al., 2022).

40 In order to better understand what factors predict diabetes, machine learning
41 (ML) techniques have been harnessed for the diagnosis of diabetes, which is famous
42 for its efficient and predictive nature (Dewangan Agrawal, 2015; Ganie et al.,
43 2022). Utilizing advanced computational algorithms, these ML approaches facilitate
44 early and accurate detection, providing valuable insights that enable healthcare
45 professionals to effectively manage and treat diabetes as well as advice on how to
46 reverse prediabetes. As society grapples with the multifaceted impact of lifestyle
47 choices on health, integrated technologies such as ML present a promising avenue
48 to address and reduce the growing burden of diabetes.

49 Therefore, in this data science report, we seek to understand how to use machine
50 learning for diabetes identification. Our research question is: How do different
51 machine learning models compare in their ability to accurately predict outcomes
52 based on a subset of five key features within the diabete dataset?

53 3. Dataset Description and Workflow

54 **Data source.** The dataset for our research, retrieved from the UCI Machine Learning
55 Repository, comprises 253,680 records with 20 features, after applying oversampling
56 to balance the data, the dataset for our research now comprises 436,668 records
57 with 21 features, an increase from the original 253,680 records These features cover
58 a range of categories, including demographics (race, sex), personal information
59 (income, education), and health history (drinking, smoking, mental health, physical
60 health),
61
62

63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124

Significance Statement

In the face of the diabetes epidemic, early detection is paramount for effective management and treatment. This research applies machine learning to identify crucial predictive factors of diabetes, harnessing a dataset representative of a broad U.S. population. By comparing five different algorithms, our study not only determines the most effective model for diabetes prediction but also underscores the potential of artificial intelligence in public health. The Multi-layer Perceptron (MLP) Classifier, recognized for its accuracy and comprehensive performance, stands out, marking a significant step toward integrating advanced analytics in healthcare.

	Name	Role
0	Diabetes_binary	Target
1	HighBP	Feature
2	HighChol	Feature
3	CholCheck	Feature
4	BMI	Feature
5	Smoker	Feature
6	Stroke	Feature
7	HeartDiseaseorAttack	Feature
8	PhysActivity	Feature
9	Fruits	Feature
10	Veggies	Feature
11	HvyAlcoholConsump	Feature
12	AnyHealthcare	Feature
13	NoDocbcCost	Feature
14	GenHlth	Feature
15	MentHlth	Feature
16	PhysHlth	Feature
17	DiffWalk	Feature
18	Sex	Feature
19	Age	Feature
20	Education	Feature
21	Income	Feature

Fig. 1. Features and Target

Notably, this dataset originates from the 2014 Behavioral Risk Factor Surveillance System (BRFSS), which has undergone significant changes in 2011, particularly in the weighting methodology (raking) and the inclusion of cell phone only respondents. The aggregated BRFSS data, combining both landline and cell phone responses, encompasses data for all 50 states, the District of Columbia, Guam, and Puerto Rico, providing a comprehensive view of the subjects' socio-economic status and health-related behaviors. This diverse dataset offers a rich foundation for our data analysis.

Unit of observations. The primary unit of observation for the BRFSS is individuals residing in the U.S. The data is collected through monthly telephone interviews and includes both landline and cellular telephone subscribers.

Methods. This data flow diagram outlines the process of a machine learning workflow for our diabetes dataset. Initially, the data undergoes preprocessing to clean and prepare it for analysis. Then, feature extraction is performed, identifying 21 perceptive features which are then narrowed down to the top 5 most significant features for model training.

These top features are used as inputs into several basic machine learning classifiers: Logistic Regression, Decision Tree, KNN (k-nearest neighbors), Naive Bayes, and Multiple Layer Perceptron. Each classifier will be trained to predict outcomes based on the processed diabetes data.

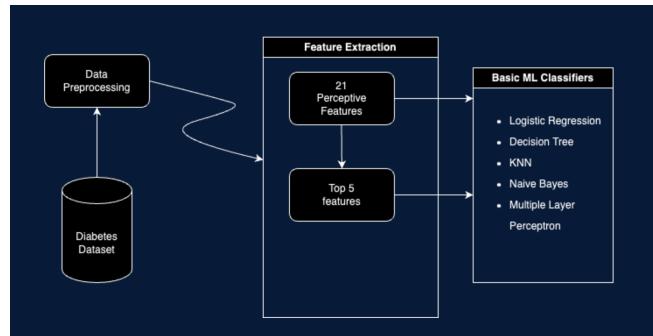


Fig. 2. Workflow Diagram

4. Feature Engineering and Exploratory Data Analysis (EDA)

Missing data. We find the missing data using following function - `print(df.isnull().sum())`. From the given results, it is clear that we do not have missing data in this data set.

HighBP	0
HighChol	0
CholCheck	0
BMI	0
Smoker	0
Stroke	0
HeartDiseaseorAttack	0
PhysActivity	0
Fruits	0
Veggies	0
HvyAlcoholConsump	0
AnyHealthcare	0
NoDocbcCost	0
GenHlth	0
MentHlth	0
PhysHlth	0
DiffWalk	0
Sex	0
Age	0
Education	0
Income	0
Diabetes_binary	0
dtype:	int64

Fig. 3. Missing Values Check

Correlation matrix of the data set. As shown in the figure below, the correlation matrix of our dataset reveals the interrelationships between the 20 feature variables. Notably, there is a strong positive correlation between "self-described general health" and "physical health", implying that those who are optimistic about their general health are also likely to report good physical health. We also observe a negative correlation between general health and income, suggesting that improvements in general health may accompany declines in reported income.

Class Balance. Given the substantial imbalance between diabetes and nondiabetic values, we've employed oversampling as a strategy to address this issue. The process starts by segregating the majority and minority classes. Subsequently, we upsample the minority class to align with the size of the majority class. The final step involves concatenating the upsampled minority class with the majority class. The result is a dataset with 218,334 instances of non-diabetic and

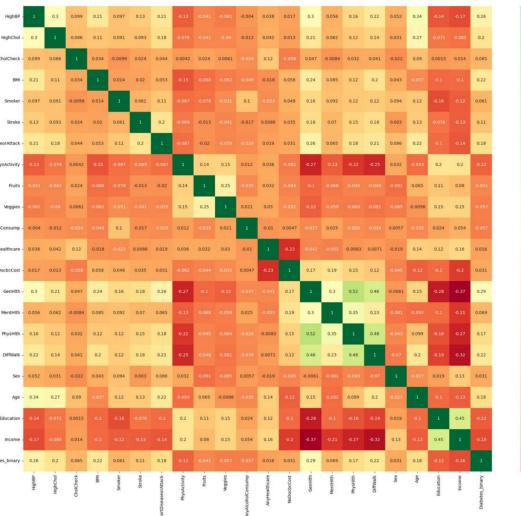


Fig. 4. Correlation Heatmap

218,334 instances of diabetic values, offering a more balanced representation of the two classes.

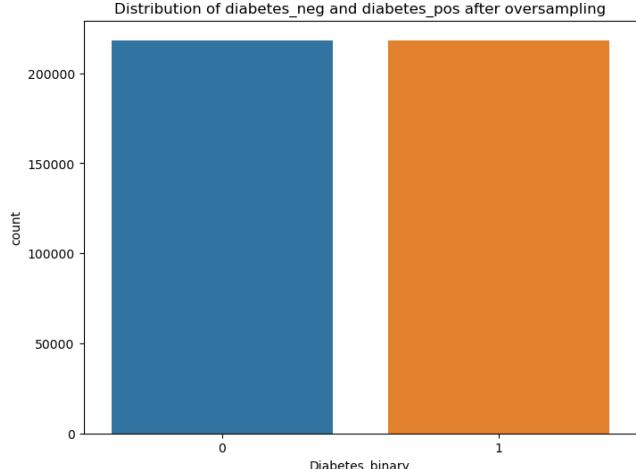


Fig. 5. Instances of Class attribute

5. Results

A. Regression Analysis. The dummy feature “Diabetes.binary” will be used as the target feature for the regression analysis. We performed the regression using the OLS model followed by t-test analysis, F-test analysis, and confidence interval analysis, and based on the t-test score and the confidence intervals values, we picked our top 5 features (BMI, Education, Age, Income, GenHlth).

B. Classification Analysis.

B.1. Logistic Regression. A Logistic Regression model was trained to predict the diabetes binary outcome based on a subset of the top 5 features from the dataset. The model’s maximum iteration parameter was set to 1000 to ensure convergence.

OLS Regression Results						
Dep. Variable:	Diabetes_binary	R-squared:	0.257	311		
Model:	OLS	Adj. R-squared:	0.257	312		
Method:	Least Squares	F-statistic:	2.717e+04	313		
Date:	Tue, 12 Dec 2023	Prob (F-statistic):	0.00	314		
Time:	18:39:30	Log-Likelihood:	-2.2698e+05	315		
No. Observations:	393001	AIC:	4.538e+05	316		
Df Residuals:	392995	BIC:	4.539e+05	317		
Df Model:	5	Covariance Type:	nonrobust	318		
	coef	std err	t	P> t	[.025	.975]
const	-0.5518	0.006	-95.324	0.000	-0.563	-0.540
BMI	0.0149	0.000	148.628	0.000	0.015	0.015
Education	-0.0096	0.001	-12.612	0.000	-0.011	-0.008
Age	0.0402	0.000	163.651	0.000	0.040	0.041
Income	-0.0108	0.000	-28.936	0.000	-0.012	-0.010
GenHlth	0.1307	0.001	186.554	0.000	0.129	0.132
Omnibus:	121897.496	Durbin-Watson:	2.006	322		
Prob(Omnibus):	0.000	Jarque-Bera (JB):	19109.722	323		
Skew:	-0.123	Prob(JB):	0.00	324		
Kurtosis:	1.948	Cond. No.	277.	325		

Fig. 6. OLS report for the top 5 features

The accuracy of the model, which represents the proportion of true results (both true positives and true negatives) among the total number of cases examined, was found to be approximately 72.49%. This suggests that the model is reasonably effective at classifying the instances correctly, although there is room for improvement.

Model Evaluation The Mean Squared Error (MSE), which provides a measure of the quality of an estimator—it is always non-negative, and values closer to zero are better—was calculated to be 0.2561. The MSE value here indicates that on average, the square difference between the predicted values and the actual values is about 0.2561, suggesting that the model predictions are moderately close to the actual values.

Confusion Matrix The confusion matrix for the model presented the following results. See the figure for Logistic Regression.

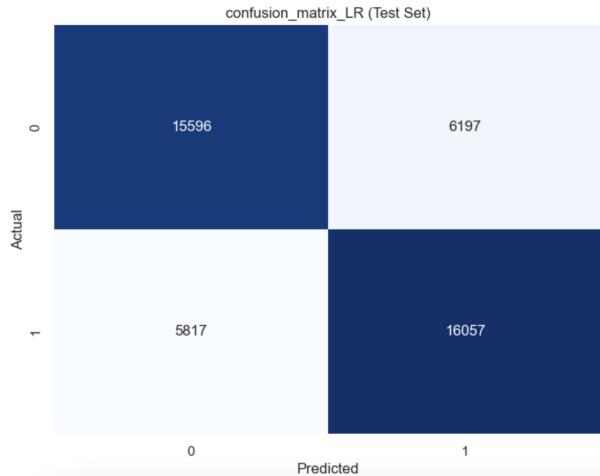


Fig. 7. Confusion Matrix for the Logistic Regression Model

From these values, we can infer that the model is more effective at identifying true positives than true negatives. However, there is a significant number of both false positives and false negatives, indicating potential areas of model improvement.

ROC Curve The Receiver Operating Characteristic (ROC) curve and Area Under the Curve (AUC) were used

to evaluate the model's discrimination capacity. The ROC curve (attached as a separate figure) demonstrates the trade-off between the true positive rate and the false positive rate. The AUC for our model was calculated to be 0.80. An AUC of 0.80 indicates a good discriminative ability and suggests that the model has a high chance of distinguishing between the positive class and the negative class.

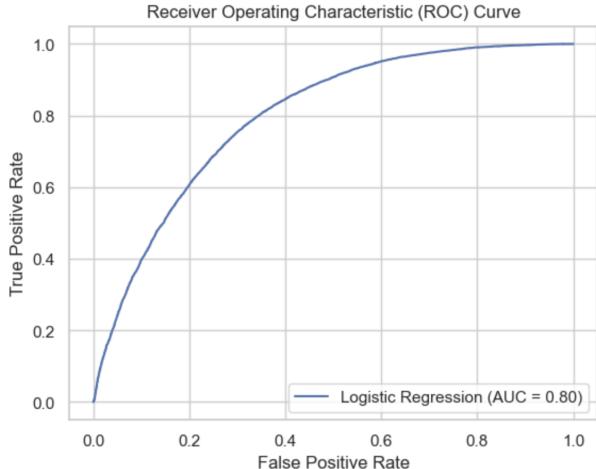


Fig. 8. ROC Curve for the Logistic Regression Model

In conclusion, while the Logistic Regression model demonstrates good predictive ability as indicated by the AUC of 0.80, there is room for improvement in terms of accuracy and reducing the number of false predictions. Further refinement of the model could involve feature engineering, hyperparameter tuning, or trying different algorithms to improve the performance.

B.2. Decision Tree. We employed a comprehensive approach that involved pre-pruning and post-pruning techniques to enhance the model's performance and generalization ability.

Pre-Pruning Pre-pruning was performed using GridSearchCV to find the optimal combination of hyperparameters. Our parameter grid included various criteria, depths, minimum samples for leaf nodes and splits, and feature counts. The process aimed to prevent the decision tree from overfitting by limiting its growth during training. The optimal parameters obtained through pre-pruning were as follows: Criterion: entropy Max Depth: 6 Max Features: None Min Samples Leaf: 3 Min Samples Split: 2 With these parameters, the pre-pruned model achieved an accuracy score of 0.7285. The Receiver Operating Characteristic (ROC) curve yielded an Area Under the Curve (AUC) of 0.80, indicating a good predictive performance.

Post Pruning Post-pruning was conducted after the decision tree was fully grown. It involved reducing the complexity of the tree by adjusting the regularization parameter, which simplifies the model and helps in avoiding overfitting. The goal was to find a value that retains model accuracy while simplifying its structure. The post-pruning process identified the best value as 0.0, suggesting that no additional regularization was necessary. The post-pruned model showcased an improved accuracy of 0.7798 on the test set. The corresponding ROC curve demonstrated an AUC of

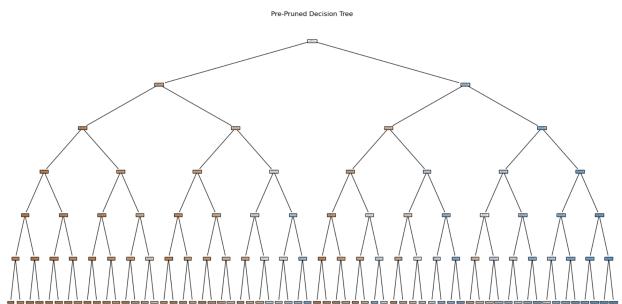


Fig. 9. Prepruned Decision Tree

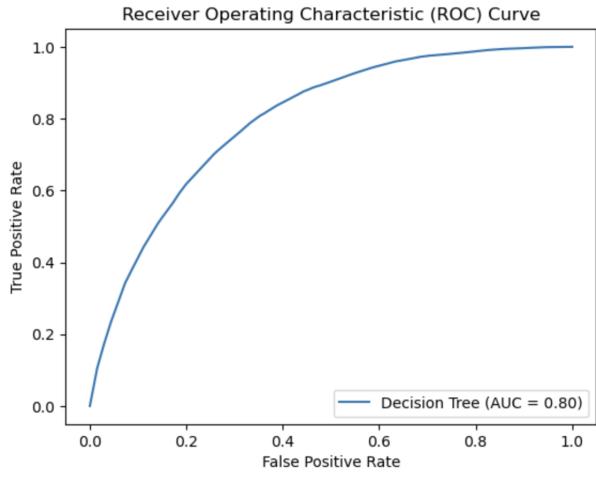


Fig. 10. ROC Curve for the Pre Decision Tree Model

0.86, indicating a superior ability to distinguish between the classes compared to the pre-pruned model.

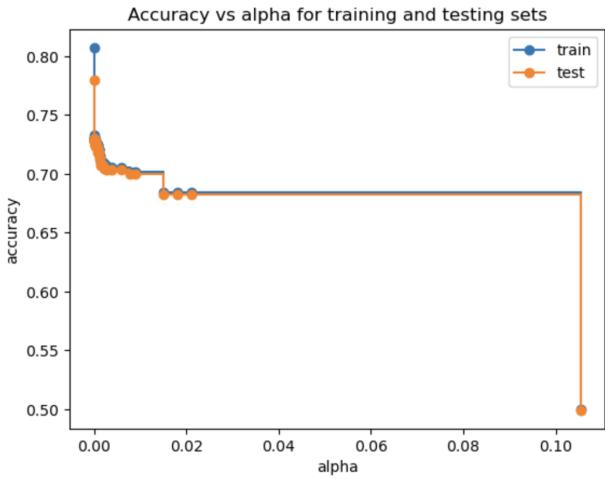
The increase in AUC from 0.80 to 0.86 through post-pruning indicates that the model's discriminatory capacity was enhanced while maintaining a low false-positive rate. The post-pruned model not only outperformed the pre-pruned model in terms of accuracy but also showed a better balance between sensitivity and specificity.

B.3. K-Nearest Neighbors (KNN). The KNN model was implemented with `n_neighbors` set to 5, which defines the number of nearest neighbors to consider for making predictions. The model was trained on a subset of features deemed most significant, referred to as `top_5_features`.

The KNN model achieved an accuracy of 0.7442 on the test dataset, indicating a strong predictive capability. The MSE for the test set was calculated to be 0.2557, which is relatively low. See figure 12 for the matrix.

The ROC curve for the KNN model has an AUC of 0.82. This value is a measure of the model's ability to discriminate between the positive and negative classes. An AUC of 0.82 indicates a good classification performance.

B.4. Gaussian Naive Bayes. The Gaussian NB model achieved a test accuracy of 0.7111, indicating a moderate level of predictive performance. The mean squared error (MSE) for the test set was calculated as 0.2888, which is higher compared



Best alpha for post-pruning: 0.0
Accuracy of the best post-pruned tree: 0.7798108411386173

Fig. 11. Accuracy vs Alpha Decision Tree

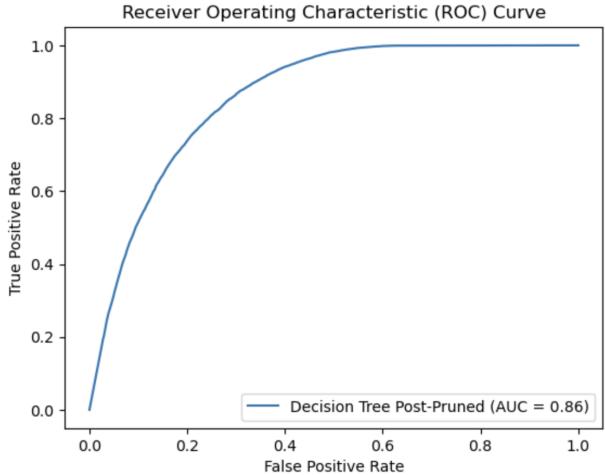


Fig. 12. ROC Curve for the Pre Decision Tree Model

to more complex models, signifying potential misclassification errors.

The ROC curve for the Gaussian NB model exhibited an AUC of 0.78. Although this indicates an acceptable ability to discriminate between the positive and negative classes, it is less than the AUCs achieved by the other models. And here is the confusion matrix:

B.5. Neural Network: Multi-layer Perceptron Classifier(MLPClassifier). MLPs consist of multiple layers of neurons, allowing them to learn and represent complex, non-linear relationships in data. The training process includes feedforward architecture, and backpropagation (implicit in scikit-learn) is used during the training process to update the weights of the neural network based on the error between predicted and actual values.

- Activation Functions:

Hidden Layers: ReLU (Rectified Linear Unit), tanh (hyperbolic tangent function) Output Layer (Binary Classification): Logistic (Sigmoid)

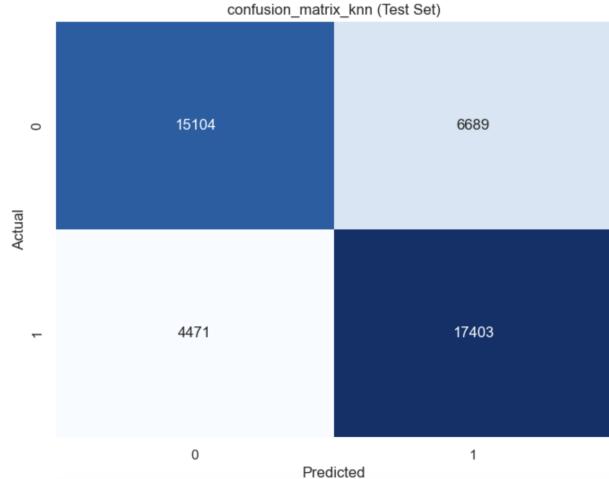


Fig. 13. Confusion Matrix for the K-Nearest Neighbors Model

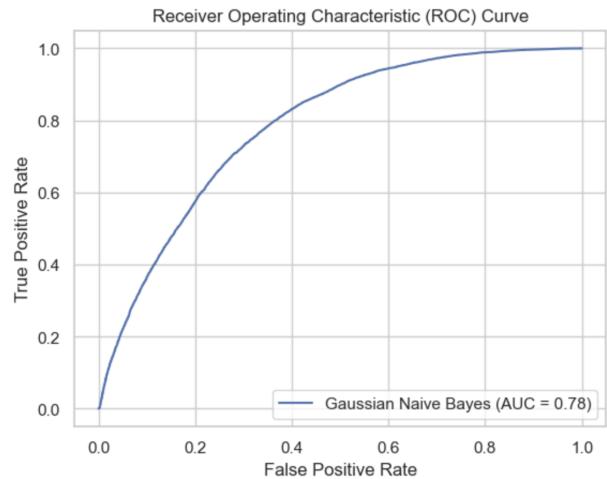


Fig. 14. ROC Curve for the K-Nearest Neighbors Model

Based on our MLP model we got the following hyperparameters as best suited for the model :

Best Hyperparameters: 'activation': 'relu', 'alpha': 0.001, 'hiddenlayer_sizes' : (60,40), 'learning_rate' : 'constant', 'solver' : 'adam'

Moreover, the test accuracy estimated by our model is : 0.7286280257402615 and the confusion matrix estimated for our model is provided in the following figure.

The confusion matrix for the MLP classifier can help in evaluating its performance by providing insights into the model's ability to correctly classify instances. It helps in assessing true positives, true negatives, false positives, and false negatives, which are essential for understanding the model's predictive accuracy and error rates.

On Analyzing the true positives, true negatives, false positives, and false negatives in the confusion matrix for the MLP classifier we got several insights: True Positives (TP) of 15,964 indicate the number of instances where the model correctly identified the positive class. True Negatives (TN) of 15,853 show the number of instances where the model correctly identified the negative class. These figures suggest that the model is comparably reliable in identifying both

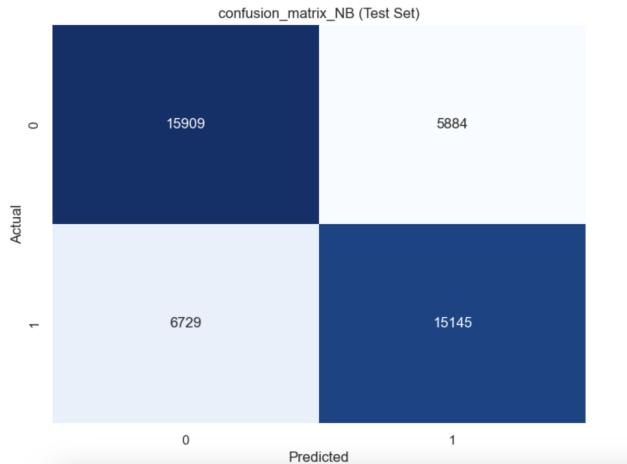


Fig. 15. Confusion Matrix for the Gaussian Naive Bayes Model

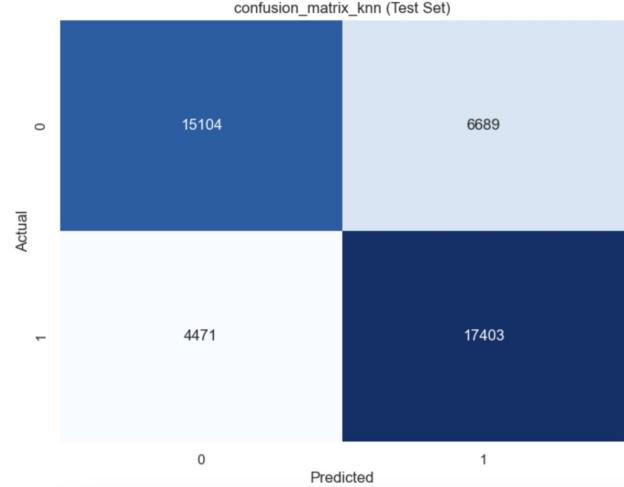


Fig. 17. Confusion Matrix for the MLPClassifier

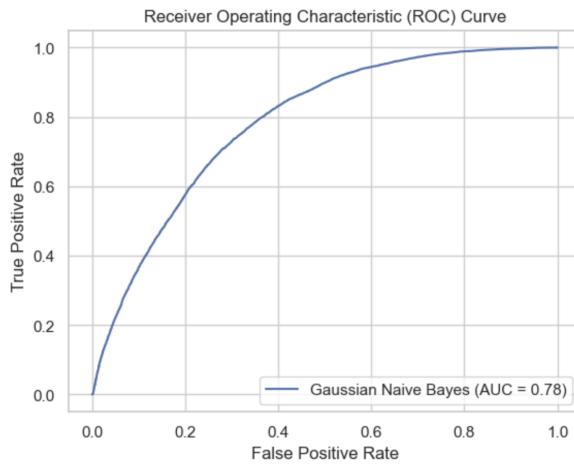


Fig. 16. ROC Curve for the Gaussian Naive Bayes Model

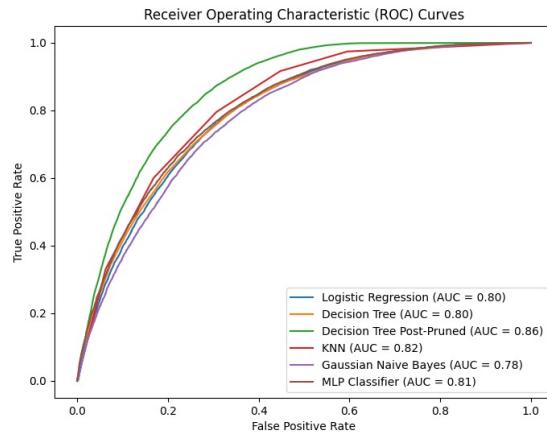


Fig. 18. ROC Curve for MLPClassifier

positive and negative classes. False Positives (FP) of 6,021 represent instances where the model incorrectly predicted the positive class. False Negatives (FN) of 5,829 are instances that were incorrectly predicted as the negative class. The quantities of false positives and false negatives suggest areas for potential model improvement, as reducing these will improve overall accuracy.

The above ROC AUC score of 0.81 for the MLP classifier indicates strong performance in discriminating between classes, which is essential for diabetes detection. It suggests that the model has a commendable ability to distinguish between positive and negative cases, making it a promising model.

C. K-fold Cross Validation of the Models: In a comparative study of five classification models for diabetes diagnosis, the models are evaluated across several metrics: accuracy, confusion matrix, precision, recall, F1-score, ROC-AUC score, and cross-validation (CV) scores using k-fold validation as presented in the table below:

6. Limitations

The study presented valuable insights but also faced certain limitations:

Unbalanced data: Despite efforts to use oversampling to balance the dataset, synthetic oversampling methods may introduce bias or overfitting. The oversampling technique does not take into account the potential complexity and underlying patterns of minorities. **Feature Selection:** The first 5 features were selected based on the provided dataset. However, other relevant features may not have been included in the study that could have improved the predictive performance of the model. **Interpretability of the model:** Although MLP classifiers show good results,

Algorithm	Accuracy	Geburte Rate	Classification techniques comparison	Precision	Recall	F1-score	AUC-AUC score	CV Score
Logistic Regression	0.704972200498686	0.0008	0.704972200498686	0.704972200498686	0.704972200498686	0.704972200498686	0.704972200498686	0.704972200498686
Decision Tree	0.720080700894643	0.0002	0.720080700894643	0.720080700894643	0.720080700894643	0.720080700894643	0.720080700894643	0.720080700894643
Decision Tree Post-Pruned	0.720080700894643	0.0002	0.720080700894643	0.720080700894643	0.720080700894643	0.720080700894643	0.720080700894643	0.720080700894643
KNN	0.7440210000000001	0.0001	0.7440210000000001	0.7440210000000001	0.7440210000000001	0.7440210000000001	0.7440210000000001	0.7440210000000001
Gaussian Naive Bayes	0.7158450000000001	0.0001	0.7158450000000001	0.7158450000000001	0.7158450000000001	0.7158450000000001	0.7158450000000001	0.7158450000000001
MLP Classifier	0.7300000000000001	0.0001	0.7300000000000001	0.7300000000000001	0.7300000000000001	0.7300000000000001	0.7300000000000001	0.7300000000000001

Fig. 19. K-fold Cross Validation Table

745 neural networks are often viewed as "black boxes" with
746 limited interpretability, which can be a significant drawback
747 in clinical settings where understanding the decision-making
748 process is critical. **Generalizability:** models are trained
749 and tested on BRFSS-specific datasets. Without further
750 validation, results may not generalize to other populations
751 or datasets. **Algorithmic Complexity:** Some algorithms,
752 particularly MLP classifiers, involve complex structures that
753 may not be required for all types of prediction tasks and may
754 be computationally expensive for large-scale applications.

755 756 **7. Future Work**

757 To address these limitations and enhance the robustness of our
758 findings, we propose the following future work: Alternative
759 resampling techniques: Experiment with different resampling
760 techniques. Cross-population validation: validate models on
761 different populations and datasets to ensure generalizability
762 and applicability of models to different populations. Hybrid
763 Models: Develop hybrid models that combine the strengths of
764 various algorithms, including MLP, to improve performance
765 while balancing complexity and interpretability.

766 767 **8. Insights and Conclusion based on the cross valida-** 768 **tion of the models**

769 Upon evaluating the cross-validation scores, which are critical
770 for assessing the stability and generalization ability of the

772 773 **9. References**

- 774 Amit Kumar Dewangan, Agrawal, P. (2015). Classification
775 of Diabetes Mellitus Using Machine Learning Techniques.
776 International
777 Journal of Engineering and Applied Sciences,
778 2(5), 257905. CDC diabetes health indicators.
779 UCI Machine Learning Repository. (n.d.).
780 <https://archive.ics.uci.edu/dataset/891/cdc+diabetes+health+indicators>
781 CDC. (2019). Diabetes basics. Cen-
782 ters for Disease Control and Prevention.
783 <https://www.cdc.gov/diabetes/basics/index.html>
784 Hu, F. B. (2011). Globalization of diabetes. Diabetes
785 Care, 34(6), 1249–1257. <https://doi.org/10.2337/dc11-0442>
786 Lai, H., Huang, H., Keshavjee, K., Guergachi, A., Gao,
787 X. (2019). Predictive models for diabetes mellitus using
788 machine learning techniques. BMC Endocrine Disorders,
789 19(1). <https://doi.org/10.1186/s12902-019-0436-6>
790 Teboul, A. (2021, November 8). Di-
791 abetes health indicators dataset. Kaggle.
792 <https://www.kaggle.com/datasets/alextreboul/diabetes->
793 [health-indicators-dataset](https://www.kaggle.com/datasets/alextreboul/diabetes-health-indicators-dataset)

807 models across different subsets of the dataset, the **MLP**
808 (**Multilayer Perceptron**) **classifier** emerges as the superior
809 model for diabetes detection. The MLP has demonstrated
810 the highest mean CV score, with a range between 75.63% and
811 76.10%, with an approximate mean of 75.86%. This suggests
812 that the MLP is likely to maintain its performance across
813 various unseen data points and is less prone to overfitting
814 compared to the

815 Prioritizing CV scores over other individual performance
816 metrics is advantageous because CV scores provide insights
817 into how well a model generalizes beyond the specific data
818 it was trained on. Cross-validation involves partitioning
819 the data into subsets, training the model on some subsets,
820 and validating it on the remaining subsets. This process is
821 repeated multiple times to ensure a more reliable estimate of
822 the model's predictive performance. The high CV scores of
823 the MLP model indicate that its performance is robust across
824 different folds of the dataset, which is important for medical
825 diagnostic models that will be applied to real-world scenarios
826 where the data may vary. Such robustness is essential in the
827 medical field, as it implies that the model will consistently
828 perform well in practice, making it the preferred choice for
829 diabetes detection based on the CV score criterion.

830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868