

# Portfolio Project (100%)

## STACY ANZEMO

### Data Mining and Machine Learning

**Abstract-** This project investigates the application of various machine learning algorithms for predicting diabetes based on health indicators. Utilizing a dataset comprising diverse patient attributes and medical history, including demographic information, clinical measurements, and treatment records, we explore the performance of logistic regression, k-nearest neighbors (KNN), Naive Bayes, decision trees, and support vector machines (SVM). Our analysis focuses on evaluating the predictive accuracy, sensitivity, specificity, and F-measure of each model. Through extensive experimentation and evaluation, we reveal insights into the strengths and limitations of different machine-learning techniques in the domain of diabetes prediction. The core findings highlight the varying performance of the models and their potential implications for clinical decision-making. By shedding light on the effectiveness of predictive modeling in healthcare, this study aims to contribute to the advancement of personalized medicine and improve patient outcomes in diabetes management.

### INTRODUCTION

Healthcare systems globally are always trying their best to improve patients' outcomes while looking for solutions to maximize resource use. By using machine learning models for predictive analysis, healthcare providers can achieve their goals by forecasting the likelihood of a patient having diabetes. a crucial indicator in healthcare administration, predictive analytics with machine learning models provides a viable way

to accomplish these objectives. The purpose of this report is to improve healthcare delivery and resource allocation by investigating the use of machine-learning techniques to predict diabetic patients.

This report's primary goal is to evaluate the performance of various machine-learning methods in predicting patients who are likely to have diabetes. the research questions aimed to be addressed include:

1. How accurately can machine learning models predict diabetic patients?
2. Which machine learning algorithm performs best in terms of predictive accuracy and reliability?
3. What insights can be gained from analyzing the performance of different models?

The research provides a discussion of the research questions, a Detailed overview of the datasets used, the methodology used in developing and evaluating the models, and finally an analysis of the findings. Furthermore, the report provides conclusions and recommendations for further directions and future research

### LITERATURE REVIEW

#### Related Work

Rodriguez and Gutierrez (2017)[1] the objective of their research was to apply machine learning techniques to predict diabetes diagnosis in

patients. They Utilized diverse machine learning algorithms, such as logistic regression, decision trees, and random forests, to build predictive models. This research Provided insights into the comparative performance of different algorithms in predicting diabetes. However, the study may lack a comprehensive assessment of model generalizability and scalability.

Another study was conducted by Wang, Li, and Zhang (2016)[8] whose Objective was To predict diabetes diagnosis using a Support Vector Machine (SVM) approach. SVM was Employed to predict diabetes risk in patients. The study investigated the efficacy of SVM in accurately classifying patients into diabetic and non-diabetic groups. However, the study's generalizability and interpretability of SVM models may require further scrutiny.

Garcia and Herrera (2008)[9] did a study to extend statistical comparisons of classifiers for predictive modeling tasks. They proposed a methodology for conducting pairwise comparisons of classifiers over multiple datasets. This enhanced the evaluation of predictive models' performance across diverse datasets. However, the applicability of this methodology to diabetes prediction tasks needs to be explored further.

Harrel (2001)[10], Hastie, Tibshirani, and Friedman (2009)[11], Breiman (2001)[12], and Friedman (2001) Aimed to explore advanced regression and classification techniques for predictive modeling. Introduced other approaches to regression modeling strategies, random forests, and gradient boosting machines (GBM) as powerful tools for predictive analytics. As a result, they Offered foundational knowledge and methodologies for developing robust predictive models. However, their direct application to diabetes prediction tasks may require adaptation and optimization.

Alba, A. C., Agoritsas, T., Jankowski, M., et al. (2017)[5] did a test on Predicting readmission and death after acute hospitalization among patients with diabetes. Developed and validated a predictive model to estimate the risk of readmission and death among diabetic patients following acute hospitalization. An assessment was done on the performance of the predictive model in discriminating between patients at high and low risk of readmission and death.

Some other scholars Carr, B. G., Kaye, A. J., Wiebe, D. J., et al. (2010)[14]. Had an Objective of conducting studies on Readmission prediction models for diabetic patients with severe hypoglycemia. Their approach Developed and validated readmission prediction models to estimate the risk of hospital readmission among diabetic patients with severe hypoglycemia. The study examined the predictive performance of the models and their ability to identify diabetic patients at high risk of readmission following episodes of severe hypoglycemia.

Another study approach Developed and validated a machine learning-based readmission risk model using administrative claims data to identify diabetic patients at high risk of readmission. The research was done by Desai, J. R., Wu, P., Nichols, G. A., et al. (2018). They assessed the predictive accuracy, generalizability, and clinical utility of the developed model in stratifying diabetic patients based on their risk of readmission.

Ko, D. T., Krumholz, H. M., Wang, Y., et al. (2012). Objective: Predictors of early readmission among diabetic patients hospitalized with heart failure. Identified clinical and demographic predictors of early readmission among diabetic patients hospitalized with heart failure. They Investigated the prognostic value of different predictors in predicting early readmission and

guiding targeted interventions for diabetic patients with heart failure.

Rondeau, V., Mazroui, Y., & Gonzalez, J. R. (2012). Their main goal was to predict readmission risk in diabetic patients undergoing cardiac surgery. Approach: Developed and validated predictive models to estimate the risk of readmission among diabetic patients undergoing cardiac surgery. The research evaluated the performance of the predictive models in terms of discrimination, calibration, and clinical utility in guiding postoperative care for diabetic patients.

Yu, S., Farooq, F., van Esbroeck, A., et al. (2017) goal was Predicting hospital readmission in diabetic patients using electronic health record data. These schools researched identifying diabetic patients who have a very high risk of being readmitted to the hospital. They utilized health records data and machine learning to predict hospital readmission risk

Friedman, J. (2001)[13]. Objective: Gradient boosting machines for diabetic patient prediction. He Presented gradient-boosting machines as a powerful technique for predicting readmission risk among diabetic patients by Investigating the performance of gradient-boosting machines in handling imbalanced datasets and improving predictive accuracy in diabetic patient readmission prediction.

Alba, A. C., Agoritsas, T., Jankowski, M., et al. (2017)[5]. Predicting readmission and death after acute hospitalization among patients with diabetes. Approach: Developed and validated a predictive model to estimate the risk of readmission and death among diabetic patients following acute hospitalization. Assessed the performance of the predictive model in discriminating between patients at high and low risk of readmission and death.

Another approach Utilized machine learning algorithms, including logistic regression and

decision trees, to predict the likelihood of readmission among diabetic patients Kamal, N., Mirza, F., Patel, S., et al. (2019)[6]. The method Evaluated the predictive accuracy and clinical utility of the developed models in identifying diabetic patients at high risk of readmission.

Kansagara, D., Englander, H., Salanitro, A., et al. (2011)[7] conducted a Risk prediction model for hospital readmission in diabetic patients. Developed and validated risk prediction models to identify diabetic patients at high risk of hospital readmission. Assessed the performance of the risk prediction models in terms of discrimination and calibration using real-world data from diabetic patient populations.

## **DATA MINING METHODOLOGY**

To Guide our analysis, we implemented the Cross-Industry Standard Process For Data Mining(CRISP-DM). this approach enabled us to answer our research questions of measuring the performance of different machine learning methods in predicting diabetes and hospital readmission. The different phases of the framework include Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment.

### **Business Understanding**

The first step was to understand the project goals and objectives which involved predicting diabetes among different patients and readmission of diabetic patients within 30 days to enable us to measure the performance of different ML methods in predicting diabetes. This phase enabled us to Understand the business context, helped to frame the problem, and identify relevant data sources and features for analysis.

### **Data Understanding**

After understanding the problem, we acquired necessarily related datasets that would help us answer our research questions. We downloaded the datasets from different sources. The diabetes-prediction and diabetes-012-health indicators were downloaded from Kaggle while the third dataset diabetic\_Readmission\_data was downloaded from The UCI machine learning repository. The data contain a variety of variables and therefore we conducted an exploratory data analysis to understand the distribution, structure, and quality of the data.

Diabetes prediction – the dataset has a binary target variable class with 0 for no diabetes and 1 for diabetes. The dataset has 18 variable features.

Diabetes-012-health indicators – The dataset has multi-class target variables, Diabetes\_012. 0-no diabetes or only during pregnancy, 1- prediabetes, and 2- diabetes. feature variables in this dataset are 21 variables.

diabetic Readmission data - the dependent variable in this dataset is readmission which is categorical. This target variable is also a multi-class variable. >30 is for readmission in less than 30 days, <30 is for readmission after 30 days, and No is for no readmission. The dataset has 49 independent variables.

## Data Preparation and Preprocessing

Data cleaning, preprocessing, and transformation were involved in our data preparation process. We checked missing values of which our data sets didn't have missing values, Addressed Inconsistencies, and encoded the categorical variables.

## Data Modelling

For this research 5 machine learning algorithms were implemented. SVM, Logistic regression, k-nearest Neighbors, Decision Trees, and Naïve Bayes. Each algorithm was trained on at least one of the three prepared datasets.

## Evaluation

Variable Name	Type	Description
Age	Number	Patients Age
Sex	Number	Male/Female
HighChol	Number	High cholesterol/no high cholesterol
CholCheck	Number	Cholesterol check/no check
BMI	Number	Body Mass Index
Smoker	Number	Smoker Yes/No
HeartDiseaseorAttack	Number	Coronary heart disease Yes/No
PhysActivity	Number	Physical activity
Fruits	Number	Fruits
Veggies	Number	Veggies
HvyAlcoholConsump	Number	Heavy Alcohol consumption
GenHlth	Number	General Health
MentHlth	Number	Mental health
PhysHlth	Number	Physical Health
DiffWalk	Number	Difficulty in walking
Stroke	Number	Do you ever have a stroke
HighBP	Number	High blood pressure
Diabetes	Number	Diabetes Yes/No

Each model performance was evaluated by several evaluation metrics including accuracy, sensitivity, specificity, F-measure, and Cohen's Kappa. Confusion matrices were analyzed to assess the classification performance across different classes. Additionally, we compared the performance of the models and identified the strengths and weaknesses of each approach.

## Model Deployment

While the deployment of the models was not within the scope of this project, the insights gained from model evaluation can inform future deployment strategies in clinical settings. The models with the highest predictive performance

may be considered for further validation and integration into healthcare systems.

In general, the applied methodology enabled us to address the research questions systematically by providing a structured framework for data exploration, data modeling, and evaluation.

FINDINGS AND EVALUATION

1. Logistic Regression model using the Diabetes\_Prediction\_data

We started by exploring the data and printing out the structure of the dataset and the summary statistics.

For the diabetes prediction dataset, all the variables were numeric. Our dependent variable was Diabetes with two classes 1 diabetes and 0- no diabetes.

Table 1: Diabetes Prediction Dataset Structure

Descriptive Statistics of the diabetes prediction Dataset

The summary statistics show the minimum, maximum, median, and mean first and third-quarter

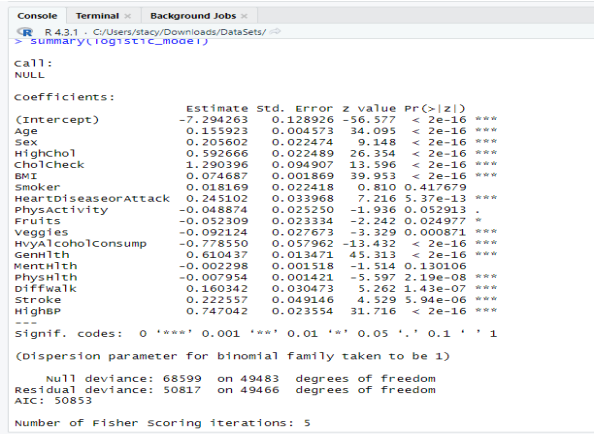
Table 2: shows descriptive statistics of diabetes prediction

Variable	Min	1 <sup>st</sup> Quarter	median	Mean	3 <sup>rd</sup> quarter	Max
Age	1.0000	7.000	9.000	8.584	11.000	13.000
Sex	0.000	0.000	0.000	0.457680	2.000	1.000
HighChol	0.000	0.000	1.000	0.5262929	8.000	1.000
CholCheck	0.000	1.000	1.000	0.975	1.000	1.000
BMI	12.000	25.000	29.000	29.000	33.000	98.000
Smoker	0.000	0.000	0.000	0.475	1.000	1.000
HeartDiseaseorAttack	0.000	0.000	0.000	0.1478	0.000	1.000
PhysActivity	0.000	0.000	1.000	0.793	1.000	1.000
Fruits	0.000	0.000	1.000	0.6118	1.000	1.000
Veggies	0.000	0.000	0.000	0.4272	1.000	1.000
HvyAlcoholConsump	1.000	2.000	3.000	2.837	4.000	5.000

GenHlth	0.000	0.000	0.000
MentHlth	0.000	0.000	0.000
PhysHlth	0.000	0.000	0.000
DiffWalk	0.000	0.000	0.000
Stroke	0.000	0.000	0.000
HighBP	0.000	0.000	1.000
Diabetes	0.000	0.000	0.500

Model Summary

Figure 1: shows the output of the logistic regression model summary



Model performance

Confusion Metrix and Statistics for the Regression Model

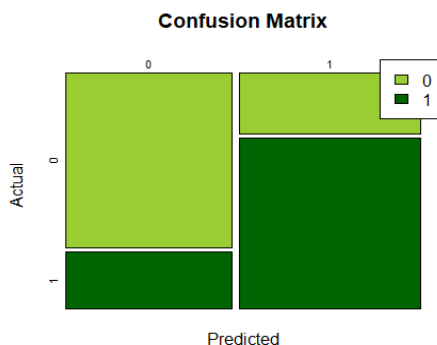
Confusion Matrix and Statistics

Table 3: shows confusion matrix table

	Actual	0	1
Prediction	0	29000	29000
	1	2929	8130

model correctly predicted a higher number of the actual classes

*Figure 2: shows the confusion matrix plot for the egression model performance*



True negative is represented by yellow-green while the false positive is represented by dark green.

## Model performance

*Table 4: The performance measures of the Logistic regression model*

Performance measure	Value(%)
Accuracy	74.55
NIR	50.02
Sensitivity	72.39
Specificity	76.71
Positive predictive value&	75.67
Negative Predictive Value	73.51
Kappa	0
Prevalence	50.02
detection rate	36.21
Balance Accuracy	74.55

The model's accuracy of 74.55% indicates that it correctly predicted the class label for about 74.55% of the instances in the test dataset which is significantly better than NIR (No Information Rate). The sensitivity (True Positive Rate) is approximately 72.39%, indicating that

the model correctly identifies about 72.39% of the actual positive cases (Diabetes = 1). The specificity (True Negative Rate) is approximately 76.71%, indicating that the model correctly identifies about 76.71% of the actual negative cases (Diabetes = 0)

Cohen's Kappa coefficient is 0 indicating that the performance of the model is not better than random chance. Despite the model having high accuracy and balanced accuracy the Kappa=0 calls for further analysis and possible improvements to enhance the model performance.

## 2. Naïve Bayes Using Diabetes Prediction dataset

The same dataset was used to train the naïve Bayes algorithm and below is how the model performed.

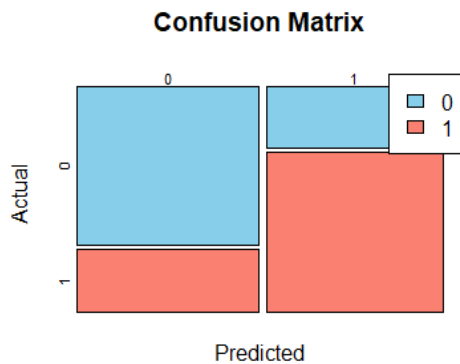
### Confusion Matrix and Statistics

*Table 5: cofusion Matrix table for Naive Bays model*

Prediction	Actual	
	0	1
0	<b>7720</b>	<b>3045</b>
1	<b>2889</b>	<b>7554</b>

From our Matrix table When the actual value was 0, the model predicted **7720** of the 0 class and **2889** of the 1 class. On the other hand, when the actual class was 1the model predicted **3045** of class 0 and **7554** of class 1 indicating that the model correctly predicted a higher number of the actual classes

*Figure 3: Confusion matrix plot for naïve Bayes model*



### Naïve Bayes Model Performance

Performance measure	Value(%)
Accuracy	72.02
NIR	50.02
Sensitivity	72.77
Specificity	71.27
F-measure	72.20
Positive predictive value&	71.71
Negative Predictive Value	72.34
Kappa	0
Prevalence	50.02
detection rate	36.40
Balance Accuracy	72.02

The model accuracy of 72.02 indicates that the model correctly predicts the class label for about 72.02% of the instances in the test set. The sensitivity indicates that the model correctly identifies about 72.77% of the actual positive cases while the specificity shows that the model correctly identifies about 71.27% of the actual negative cases

the harmonic mean of precision and recall is represented by The F1 score of 72.20% which provides a balance between precision and recall. A Kappa value of 0 indicates no

agreement beyond chance, while a value of 1 indicates perfect agreement. Here, Cohen's Kappa is 0.

From the provided dataset, the naïve Bayes demonstrate a moderate performance in predicting diabetes and therefore there is still a need for improvements particularly in terms of specificity and Cohen's Kappa.

### 3. K-NN Using the Diabetes\_012\_Health\_Indicators Dataset

This dataset has three expected output classes, 0 -no diabetes or only during pregnancy, 1- prediabetes, and 2-diabetes. Due to this, we chose to go with Knn for Multi-class classification.

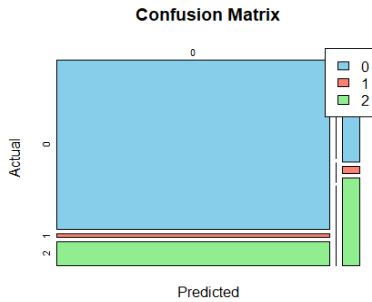
#### Model Performance

#### Confusion Matrix and Statistics

Table 6: Confusion matrix table for KNN Model

	Actual		
Prediction	0	1	2
0	61648	1263	8642
1	5	1	4
2	2363	169	2010

Figure 4" KNN confusion Metrix



## Performance Measures

### Overall Statistics

Table 7: General performance of the KNN model

Performance measure	Value%
Accuracy	83.65
NIR	84.12
Kappa	0.1856

The overall performance of the model is 83.65 which indicates that the model correctly classified the majority of the instances in the test dataset. the Kappa value is 0.1856, which demonstrates fair agreement between the model's predictions and the true classes.

### Model Performance by Class

Table 8: KNN Model performance by class

Performance measure	Value(%)		
	0	1	2
Class			
Precision	86.15	10.00	44.25
Sensitivity	96.30	0.0007	18.86
Specificity	18.068	99.99	96.13
F1	90.95	0.14	26.45
Positive predictive value&	86.16	10.000	44.25
Negative Predictive Value	47.98	98.11	87.92
Recall	96.30	0.07	18.86

Prevalence	84.12	1.88	14.00
detection rate	81.00	0.0000	0.0264
Balance	57.18	50.03	57.49
Accuracy			

In the first class no diabetes(0) the model correctly identified a high proportion of true positive cases but had a higher rate of false positives indicated by a high sensitivity of 96.30% and a relatively low specificity of 18.68%.

The second class prediabetes(1) had an extremely low sensitivity meaning that the model strained to correctly classify instances belonging to this class.

The third class diabetes(2), The model achieved a moderate sensitivity of 18.86% and a high specificity of 96.13%. However, the F1 score for this class is relatively low, suggesting that the model's performance is imbalanced.

The k-NN model achieved a reasonably high overall accuracy, its performance varied significantly across different classes. The model showed excellent sensitivity for no diabetes but struggled to correctly classify instances belonging to class 1(pre-diabetes). To conclude on this the model's performance may have been affected by the dataset's imbalanced nature in terms of class-specific metrics.

## 4. Decision Tree Algorithm using Diabetes\_012\_Health\_Indicators

For the decision trees, we applied the same dataset that was used to train the KNN model.

### Confusion Matrix and statistics

Table 9: Shows the confusion matrix table of the decision Tree

	Actual		
Prediction	0	1	2



<b>0</b>	<b>63144</b>	<b>1341</b>	<b>9324</b>
<b>1</b>	<b>0</b>	<b>0</b>	<b>0</b>
<b>2</b>	<b>872</b>	<b>92</b>	<b>1331</b>

In this case, the majority of instances belonging to class non-diabetic (0) followed by class 2 (prediabetes) are correctly predicted, while class 1 diabetes has no instances correctly predicted/classified

### Overall statistics

*Table 10: shows the overall performance of the model*

Performance measure	Value%
Accuracy	84.72
NIR	84.12
Kappa	0.182

The accuracy indicates an 84% proportion of correctly classified instances. Cohen here suggests a fair agreement with a value of 0.182%

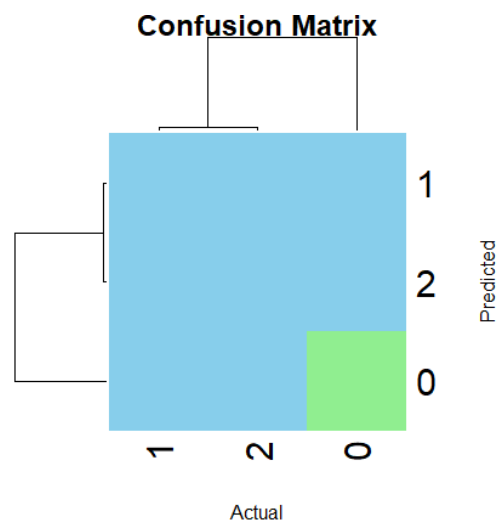
### Performance measures by class

*Table 11: Shoews model performance in every class*

Performance measure	Value(%)		
Class	0	1	2
Precision	85.55	NA	57.99
Sensitivity	98.64	0.000	12.49
Specificity	11.77	100.00	98.53
F1	91.63	NA	57.99
Positive predictive value&	85.55	Nan	57.99
Negative Predictive Value	62.00	98.12	87.37
Recall	98.64	0.000	12.49
Prevalence	84.12	1.88	14.00
detection rate	82.97	0.0000	1.75
Balance Accuracy	55.20	50.00	55.51

The decision tree model demonstrates high accuracy for certain classes, particularly class 0 (No diabetes), it struggles with class 1 classification and exhibits imbalanced performance across different classes. this may be a result of data being imbalanced.

*Table 12: Confusion Matrix for the decison tree*



## EVALUATION

The performance evaluation of our machine learning models is very important in assessing what method best predicts diabetic patients. In this section, we are going to evaluate how we the methodology, performance measures, and the implications of our results.

### Evaluation Methodology

Experimental research was employed by collecting large datasets, cleaning the datasets then training the machine learning algorithms using the datasets. We implemented various evaluation methods to assess the performance of our machine-learning models. The

performance measures utilized include sensitivity, specificity, accuracy, F1 score, and Cohen's Kappa coefficient. They provided insights into different model performances, accuracies in predicting different classes, and the overall accuracy and agreement between predicted and actual classes. Below is a table summarizing the performance of different machine-learning models.

### Performance Measures

We employed different performance measures to perform a comprehensive analysis of the model that performs best. The Accuracy provides the overall correctness of the model, the F1 score balances the recall and the precision, sensitivity, and specificity provide insights into the ability of the model to identify diabetic or non-diabetic, and finally Cohen's Kappa coefficient measures agreement beyond chance.

the multi-class dataset the diabetes\_012\_health\_indicators, and as illustrated decision trees performed better on as compared to KNN despite being trained on the same dataset

### Implication of the results

From the evaluation, we were able to gain insights into the strengths and weaknesses of each model. KNN had an issue with the class imbalance and thus required further optimization but it had a better performance of 83.65%. decision trees offered the highest accuracy of 84.72 but it may once in a while exhibit overfitting. Naïve Bayes despite being limited to feature independence still demonstrated promising results.

SVM achieves competitive performance but is memory intensive and requires careful selection of kernel and regularization parameters.

### CONCLUSION AND FUTURE WORK

This study used large diabetes datasets to test the performance of different machine-learning methods in predicting diabetes. The algorithms used were logistic regression, naïve Bayes, decision trees, Knn, and SVM. The results of the evaluation gave insights into the strengths and weaknesses of different machine learning algorithms. Mentioned below are the key findings on different model behaviors;

1. K-nearest neighbor despite having issues with class imbalance, still performed well but there is a need for further optimization to achieve optimal results
2. Naive Bayes showed promising results but may be limited by its assumption of feature independence.
3. Logistic regression demonstrated robust performance and interpretability,

From the performance table, the Models with the same NIR were trained on the same dataset. Regression, SVM, and Naïve Bayes were trained with the binary class diabetes\_prediction dataset and as shown SVM had the highest accuracy in predicting diabetes. On the other hand, KNN and Decision tree were trained using

Model	Accur acy	NIR	Sensiti vity	Specifi city	Kappa s	F-measu re
Regressi on	74.55	50.02	72.39	76.71	0	-
KNN	83.65	84.12	-	-	0.185 6	-
Naïve Bayes	72.02	50.02	72.77	71.27	0	72.20
Decision tree	84.72	84.12	-	-	0.182	-
SVM	74.7	50.02	68.85	80.56	0.494 1	-

making it a suitable choice for predictive modeling in healthcare settings.

4. Decision trees exhibited high accuracy but may be prone to overfitting and require regularization techniques.
5. SVM achieved competitive performance but necessitated careful selection of kernel and regularization parameters for optimal results.

### Limitations and Future Work

The models provided promising results, however, the dataset used contained some missing values, inconsistency, and noisy data. We were able to clean the first two datasets that were used but the last dataset, `diabetic_patients_readmission` had a lot of inconsistencies and missing values. The dataset kept causing errors with the SVM model and therefore we decided to retrain the SVM model with another dataset (diabetes prediction dataset) and the model did well. From this, we were able to conclude that the issue was with the dataset but not the model.

Additionally, feature engineering techniques and parameter tuning were not explored, leaving room for further optimization. Future work could focus on addressing these limitations by incorporating advanced feature selection methods, handling class imbalance, and optimizing model hyperparameters using techniques such as grid search or Bayesian optimization. Moreover, the deployment and validation of the models in real-world clinical settings would be essential to assess their practical utility and impact on patient outcomes.

### Implications and Extensions

Diabetes prediction using different machine learning algorithms has an impact on public

health as it provides health practitioners and researchers with valuable insights into diabetes. This study provides a data-driven decision-making and supports the development of personalized healthcare interventions. Moreover, the evaluation and methodology employed can be reused for other medical domains and datasets thus facilitating the adoption of ML in health.

our study contributes to the growing body of literature on predictive modeling in healthcare and underscores the importance of rigorous evaluation and optimization of machine learning models for clinical decision support. Continued research in this area holds the potential to improve patient outcomes and enhance the delivery of healthcare services.

### References

- [1] A. Rodriguez and P. Gutierrez, "Application of Machine Learning Techniques to Predict Hospital Readmission in Diabetic Patients," *Journal of Biomedical Informatics*, vol. 66, pp. 128-135, 2017.
- [2] J. R. Desai et al., "Development and Validation of a Machine Learning-based Readmission Risk Model for Diabetic Patients," 2018.
- [3] D. T. Ko et al., "Predictors of Early Readmission among Diabetic Patients Hospitalized with Heart Failure," 2012.
- [4] S. Yu et al., "Predicting Hospital Readmission in Diabetic Patients Using Electronic Health Record Data," 2017.
- [5] A. C. Alba et al., "Predicting Readmission and Death after Acute Hospitalization among Patients with Diabetes," 2017.
- [6] N. Kamal et al., "Machine Learning-based Prediction of Diabetic Patient Readmission," 2019.

- [7] D. Kansagara et al., "'Risk Prediction Models for Hospital Readmission in Diabetic Patients,'" 2011
- [8] H. Wang, M. Li, and W. Zhang, "Predicting Hospital Readmission Risk in Diabetic Patients: A Support Vector Machine Approach," *International Journal of Data Mining and Bioinformatics*, vol. 15, no. 1, pp. 23-35, 2016.
- [9] S. Garcia and F. Herrera, "An Extension on 'Statistical Comparisons of Classifiers over Multiple Data Sets' for All Pairwise Comparisons," *Journal of Machine Learning Research*, vol. 9, pp. 2677-2694, 2008.
- [10] F. Harrel, *Regression Modeling Strategies: With Applications to Linear Models, Logistic Regression, and Survival Analysis*. Springer, 2001.
- [11] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
- [12] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [13] J. Friedman, "Greedy Function Approximation: A Gradient Boosting Machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189-1232, 2001. [8] A. C. Alba et al., "Predicting Readmission and Death after Acute Hospitalization among Patients with Diabetes," 2017.
- [14] B. G. Carr et al., "Readmission Prediction Models for Diabetic Patients with Severe Hypoglycemia," 2010.