

UNIVERSITETI I PRISHTINËS "HASAN PRISHTINA"  
FAKULTETI I SHKENCAVE MATEMATIKO-NATYRORE  
DEPARTMENT OF MATHEMATICS, COMPUTER SCIENCE PROGRAM



**Machine Learning  
Seminar Project  
Data Analysis for Diabetes Prediction Using Health Indicators**

Albana Rexhepi, Eljesa Kqiku, Kaltrina Kuka

May 2025

## CONTENT

# 1 Introduction

## 1.1 Dataset description

This dataset contains detailed health indicators collected from a large population and is designed to support the analysis and prediction of diabetes risk. It includes a variety of columns describing lifestyle factors and demographic data, such as high blood pressure (HighBP), high cholesterol (HighChol), body mass index (BMI), smoking and alcohol use, physical activity, diet, mental and physical health status, as well as gender, age, education level, and income.

The main variable of interest is Diabetes\_binary, which indicates whether an individual has been diagnosed with diabetes (1) or not (0). This dataset can be used to develop classification models aimed at predicting diabetes risk based on known behavioral and health-related factors.

index	feature	description
Diabetes_binary	0 = no diabetes, 1 = diabetes	Diabetes status
HighBP	0 = no high BP, 1 = high BP	High blood pressure
HighChol	0 = no high cholesterol, 1 = high cholesterol	High cholesterol
CholCheck	0 = no cholesterol check in 5 years, 1 = yes	Cholesterol check
BMI	Continuous	Body Mass Index
Smoker	0 = no, 1 = yes	Smoked at least 100 cigarettes in life
Stroke	0 = no, 1 = yes	Ever had a stroke
HeartDiseaseorAttack	0 = no, 1 = yes	Heart disease or heart attack history
PhysActivity	0 = no, 1 = yes	Physical activity in past 30 days (not job-related)
Fruits	0 = no, 1 = yes	Consumes fruit 1+ times per day
Veggies	0 = no, 1 = yes	Consumes vegetables 1+ times per day
HvyAlcoholConsump	0 = no, 1 = yes	Heavy alcohol consumption
AnyHealthcare	0 = no, 1 = yes	Healthcare coverage
NoDocbcCost	0 = no, 1 = yes	Couldn't see a doctor due to cost
GenHlth	1 = excellent, 5 = poor	General health status
MentHlth	1-30 days	Poor mental health days in the past 30 days
PhysHlth	1-30 days	Physical illness/injury days in past 30 days
DiffWalk	0 = no, 1 = yes	Difficulty walking or climbing stairs
Sex	0 = female, 1 = male	Gender
Age	1 = 18-24, 9 = 60-64, 13 = 80 or older	Age category
Education	1-6	Education level
Income	1-8	Income scale

Table 1: List of all attributes in the dataset

## 1.2 Motivation of the study

The motivation of this study is based on three main questions related to understanding the risk factors for diabetes and how these factors can be used to predict the likelihood of developing the disease:

1. **How can machine learning models be used to predict diabetes risk more accurately?** This broader question relates to using advanced techniques to help create a reliable model for predicting diabetes. Machine learning can uncover relationships and patterns that may not be immediately apparent using traditional methods.
2. **Can a subset of factors be used to accurately predict whether an individual has diabetes?** An important question is whether only a few factors, such as BMI and Age, can be used to make an accurate prediction of diabetes, simplifying the prediction process without losing too much accuracy.

Through these questions, the study aims to improve our understanding of the factors influencing diabetes risk and to develop a simple yet accurate method for predicting who may be at risk.

### 1.3 Selection of algorithms for the dataset

For this dataset, where the target variable is `Diabetes_binary` (diabetes status), three well-known machine learning algorithms have been chosen: **MLP (Multilayer Perceptron)**, **Autoencoder**, and **Naïve Bayes**. These algorithms are suitable for this type of problem for various reasons:

- **MLP (Multilayer Perceptron):**  
MLP is a type of neural network that uses multiple processing layers and is excellent for handling classification and regression problems. This algorithm is powerful for capturing complex relationships between different features. For our dataset, it is well-suited to model the connections between health indicators and diabetes risk, as it can effectively analyze factors such as BMI, Age, Physical\_Health, Smoking, and other related variables.
- **Autoencoder:**  
An Autoencoder is a type of neural network used primarily for unsupervised learning and dimensionality reduction. It can help identify hidden patterns in the data and is useful for reducing the number of features while maintaining the underlying structure. For our dataset, it can assist in detecting underlying factors contributing to diabetes risk by learning an efficient representation of the data. Autoencoders are also useful for feature extraction, especially in cases where the data is high-dimensional.
- **Naïve Bayes:**  
Naïve Bayes is a probabilistic supervised learning algorithm based on Bayes Theorem, assuming that all features are independent of one another given the target class (in this case, diabetes status). Although this assumption of independence rarely holds true in real-world data, Naïve Bayes often performs remarkably well, especially in binary classification problems with multiple features. For our dataset, this algorithm is particularly well-suited due to its simplicity, fast training time, and effectiveness even with high-dimensional data.

## 2 Preprocessing

At this part of the project, we begin with the data preprocessing phase. This step is necessary to prepare the dataset for machine learning, ensuring the data is ready for analysis and modeling.

### 2.1 Data cleaning

#### 2.1.1 Removing Duplicates

For this step, we first analyze if there are any duplicate rows in the dataset using the following Python code:

```
1 import pandas as pd
2
3 # Load the dataset
4 df = pd.read_csv('datasets/diabetes_binary_5050split_health_indicators_BRFSS2015.csv')
5
6 # Check for duplicate rows
7 duplicates = df[df.duplicated()]
8
9 # Display the result
10 if duplicates.empty:
11     print("No duplicates found!")
12 else:
13     print(f"Found {duplicates.shape[0]} duplicate rows.")
14     print(duplicates)
```

After running the code, we found that there are 1635 duplicate rows. However, since these duplicates could potentially represent individuals who share the same health indicators (such as age, BMI, or smoking status), and not necessarily represent data entry errors, we decided not to remove them. This is because, in medical datasets, it is common for multiple individuals to have identical data points across various features. Therefore, removing these rows could lead to loss of valid data.

```
1 Found 1635 duplicate rows.
```

As shown, there are 1635 duplicate rows, but they are not removed to ensure we preserve all relevant data for analysis.

### 2.1.2 Handling missing values

In this step, we aim to check whether the dataset contains any missing values. Missing data can impact the accuracy of machine learning models, so identifying and addressing it is an important part of the data preprocessing process. Below is the Python code used for this task:

Finally, we used the z-score method, which was more suitable for our dataset. This method found 1,927 outlier records, which is only 2.7% of the total dataset. Since this number was small, we thought it was reasonable to remove these records without further analysis, as we still had enough data left to work with.

```
1 import pandas as pd
2
3 # Load the dataset
4 df = pd.read_csv('datasets/diabetes_binary_5050split_health_indicators_BRFSS2015.csv')
5
6 # Check for missing values
7 missing_values = df.isnull().sum()
8
9 # Display missing values per column
10 print(missing_values)
```

The output after executing the script was:

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

As we can see, all columns in the dataset have a count of 0 missing values. This means the dataset is complete in terms of data presence, and no additional cleaning or imputation for missing values is necessary.

### 2.1.3 Handling outliers

The first step we used to analyze the outlier or noise values was checking the min, max, average and mode of each attribute. From these data we can conclude that there are no noises since all the min-max fields are within the range declared on the metadata.

Next, we applied the interquartile range (IQR) method to identify outliers. However, this approach was not very effective, as it flagged 40,205 records as containing at least one outlier attribute—more than half of the dataset. Because this didn't seem like a logical result, we decided to ignore this method and try something else.

The next test was the z-score method and this one was far more reasonable giving 1927 outlier lines. Since that would be less than 3% of the dataset, it is acceptable to remove the outlier lines completely. That left our dataset with 68,765 remainin records, a considerable number of rows.

```
1 import pandas as pd
2 import numpy as np
3 from scipy.stats import zscore
4
5 df = pd.read_csv('datasets/diabetes_binary_5050split_health_indicators_BRFSS2015.csv')
6
7 ##### Making a summary of the fields #####
8 min_vals = df.min()
9 max_vals = df.max()
10 avg_vals = df.mean()
11 mode_vals = df.mode().iloc[0]
12 summary_df = pd.DataFrame({
13     'Min': min_vals,
14     'Max': max_vals,
15     'Mean': avg_vals,
16     'Mode': mode_vals
17 })
18 print("##### Summary of the dataframe #####")
19 print(summary_df)
20
21
22 print("\n##### Checking for outliers with Interquartile Range (IQR) method #####")
23 Q1 = df.quantile(0.25)
24 Q3 = df.quantile(0.75)
25 IQR = Q3 - Q1
26 outliers = (df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))
27 outlier_counts = outliers.sum(axis=1)
28 print("Rows with >=1 outliers:", (outlier_counts >= 1).sum())
29
30
31 print("\n##### Checking for outliers with Z-Score Method (Assumes Normal Distribution) #####")
32 z_scores = df.select_dtypes(include='number').apply(zscore)
33 outliers = (z_scores.abs() > 5)
34 print("Outlier rows (Z-score):", outliers.any(axis=1).sum())
35
36
37 ##### Removing the outliers #####
38 rows_to_remove = outliers.any(axis=1)
39 cleaned_df = df[~rows_to_remove]
40 cleaned_df.to_csv('datasets/dataset_without_outliers.csv', index=False)
41 print("Cleaned dataset saved as 'dataset_without_outliers.csv'")
```

```
1 ##### Summary of the dataframe #####
2
3 Diabetes_binary      Min    Max      Mean    Mode
4 HighBP              0.0    1.0    0.500000    0.0
5 HighChol            0.0    1.0    0.525703    1.0
```

```

6 CholCheck      0.0    1.0    0.975259    1.0
7 BMI            12.0   98.0   29.856985   27.0
8 Smoker         0.0    1.0    0.475273    0.0
9 Stroke         0.0    1.0    0.062171    0.0
10 HeartDiseaseorAttack 0.0    1.0    0.147810    0.0
11 PhysActivity  0.0    1.0    0.703036    1.0
12 Fruits        0.0    1.0    0.611795    1.0
13 Veggies       0.0    1.0    0.788774    1.0
14 HvyAlcoholConsump 0.0    1.0    0.042721    0.0
15 AnyHealthcare 0.0    1.0    0.954960    1.0
16 NoDocbcCost   0.0    1.0    0.093914    0.0
17 GenHlth       1.0    5.0    2.837082    3.0
18 MentHlth      0.0   30.0    3.752037    0.0
19 PhysHlth      0.0   30.0    5.810417    0.0
20 DiffWalk      0.0    1.0    0.252730    0.0
21 Sex           0.0    1.0    0.456997    0.0
22 Age           1.0   13.0    8.584055   10.0
23 Education     1.0    6.0    4.920953    6.0
24 Income        1.0    8.0    5.698311    8.0
25
26 ##### Checking for outliers with Interquartile Range (IQR) method #####
27 Rows with >=1 outliers: 40205
28
29 ##### Checking for outliers with Z-Score Method (Assumes Normal Distribution) #####
30 Outlier rows (Z-score): 1927
31 Cleaned dataset saved as 'dataset_without_outliers.csv'

```

## 2.2 Data integration

In this project, all the required data is already combined into a single CSV file. Therefore, no additional integration from multiple sources is necessary. The dataset is self-contained and ready for the next steps of data preprocessing.

## 2.3 Data transformation

In this step, we apply feature scaling to the dataset in order to standardize the features. Scaling is necessary because some machine learning algorithms, such as the ones we have chosen for this analysis, are sensitive to the scale of the data. These algorithms work better when all features have a similar range.

### 2.3.1 Data scaling

In this step, we apply feature scaling to the dataset in order to standardize the features. Scaling ensures that all features have a similar range. The `StandardScaler` from the `sklearn.preprocessing` library is used to perform scaling.

The following Python code demonstrates the process of scaling the features of the dataset:

```

1 from sklearn.preprocessing import StandardScaler
2 import pandas as pd
3
4 # Load the dataset
5 df = pd.read_csv('datasets/dataset_without_outliers.csv')
6
7 # Separate features and target
8 X = df.drop(columns=['Diabetes_binary']) # Features
9 y = df['Diabetes_binary']               # Target
10
11 # Initialize and apply the scaler
12 scaler = StandardScaler()
13 X_scaled = scaler.fit_transform(X)
14

```

```

15 # Convert scaled features back to DataFrame
16 X_scaled_df = pd.DataFrame(X_scaled, columns=X.columns)
17
18 # Add the target column back
19 X_scaled_df['Diabetes_binary'] = y.values
20
21 # Save to a new CSV file
22 X_scaled_df.to_csv('datasets/diabetes_scaled.csv', index=False)
23
24 print("Scaled dataset saved as 'datasets/diabetes_scaled.csv'.")

```

This code scales all the features of the dataset using the `StandardScaler`, which standardizes the features by removing the mean and scaling to unit variance. The scaled features are then saved in a new CSV file called `diabetes_scaled.csv`.

## 3 Prepare for analysis

### 3.1 Correlation between features

Understanding how features interact within the dataset is a key step before building predictive models. One effective way to explore these interactions is through correlation analysis. This technique reveals the strength and direction of linear relationships between variables, helping us detect redundant features, highlight those most relevant to the target variable, and optimize the feature set for better model efficiency and accuracy. By refining the input space early on, we set a strong foundation for the performance of our machine learning models.

One important aspect of correlation analysis is examining how strongly each feature relates to the target variable, in this case, `Diabetes_binary`. Identifying variables that show a higher correlation with the target can guide us in selecting the most predictive features for our model. On the other hand, features that exhibit very weak relationships with the target may contribute little to the predictive power and can be excluded to simplify the model. Additionally, correlation analysis helps reveal whether certain features are highly related to each other. When two features show a strong correlation such as above 0.5 it may indicate redundancy. Retaining both could introduce multicollinearity, which not only complicates model interpretation but can also impair the algorithm's performance. In such cases, it's often beneficial to remove one of the correlated features to streamline the dataset and enhance model robustness.

To calculate the correlation between features, we used the `corr()` function from the `pandas` library. We visualized the results using a Heatmap. A heatmap is a graphical representation that uses color gradients to indicate the strength of relationships in a two-dimensional matrix. This is especially useful for identifying patterns and connections between features.

The implementation can be seen in the code below, and the visual results are presented in Figure 1:

```

1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 df = pd.read_csv('datasets/diabetes_binary_5050split_health_indicators_BRFSS2015.csv')
6
7 # Compute correlation matrix
8 corr_matrix = df.corr()
9
10 # Plot the heatmap
11 plt.figure(figsize=(12, 10))
12 sns.heatmap(corr_matrix, annot=True, fmt=".2f", cmap="coolwarm", square=True)
13 plt.title("Correlation Matrix Heatmap")
14 plt.tight_layout()
15 plt.savefig("report/images/diabetes_correlation_matrix.png")
16 plt.close()

```



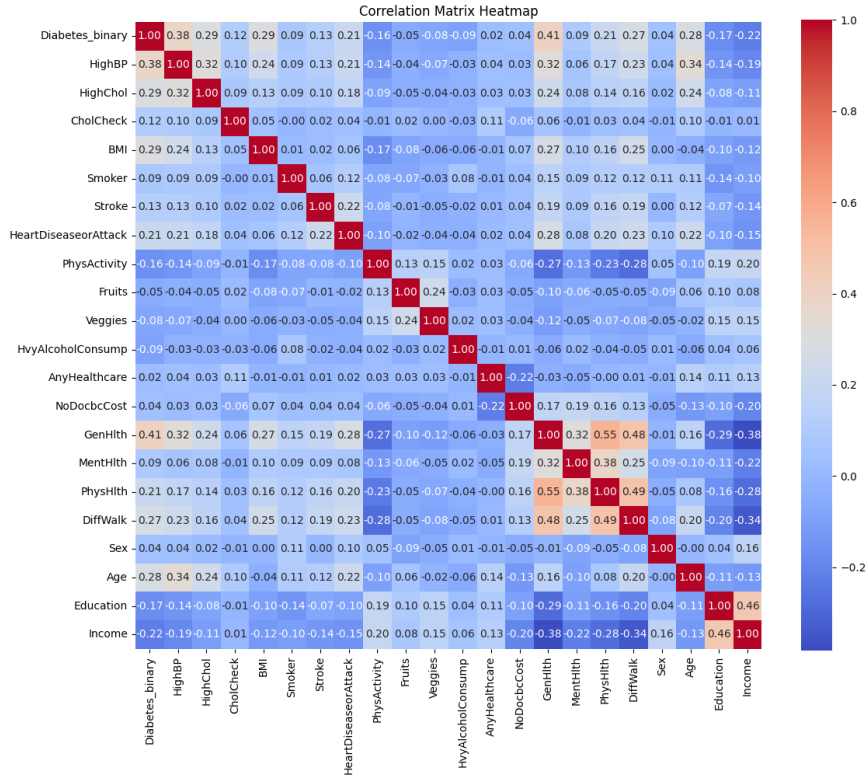


Figure 1: Correlation matrix between features.

From the heatmap visualization, we observe that some features are more strongly correlated with each other. We then filtered the feature pairs that have an absolute correlation higher than 0.45. The implementation of this filtering process in Python is shown below:

```
1 import pandas as pd
2 import numpy as np
3
4 df = pd.read_csv('datasets/diabetes_binary_5050split_health_indicators_BRFSS2015.csv')
5 corr_matrix = df.corr()
6
7 # Unstack and filter correlations
8 high_corr = (
9     corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k=1).astype(bool))
10     .stack()
11     .reset_index()
12 )
13 high_corr.columns = ['Feature 1', 'Feature 2', 'Correlation']
14 filtered_corr = high_corr[high_corr['Correlation'].abs() > 0.45]
15 print(filtered_corr.sort_values(by='Correlation', ascending=False))
```

The results are displayed in Table 2:

## Feature Selection Strategy

To optimize the feature set for model performance, we implemented a two-step strategy:

Feature 1	Feature 2	Correlation
GenHlth	PhysHlth	0.552757
PhysHlth	DiffWalk	0.487976
GenHlth	DiffWalk	0.476639
Education	Income	0.460565

Table 2: Feature pairs with correlation  $c > |0.45|$

## 1. Removal of weakly correlated features

To address weak correlation with the target variable, we will remove features exhibiting an absolute Pearson correlation coefficient of less than 0.1 with the 'Diabetes\_binary' column. This step aims to eliminate features with a minimal linear relationship to the outcome we are trying to predict, potentially reducing noise and simplifying the model.

## 2. Mitigation of multicollinearity

Subsequently, to mitigate the issue of multicollinearity among the remaining features, we will identify pairs of features with a high absolute Pearson correlation coefficient (greater than 0.45) as shown in table 2. For each such pair, one of the features will be removed. The decision on which feature to remove will be based on its correlation with the target variable; the feature with the weaker absolute correlation to 'Diabetes\_binary' will be prioritized for removal. This strategy aims to reduce redundancy in the feature set, which can lead to instability and interpretability issues in some models, without significantly sacrificing predictive power related to the target.

The implementation process in Python is shown below:

```

1 import pandas as pd
2 import numpy as np
3
4 # Load the dataset
5 df = pd.read_csv('datasets/diabetes_binary_5050split_health_indicators_BRFSS2015.csv')
6 target = 'Diabetes_binary'
7
8 # Step 1: Remove features with very low correlation to the target
9 abs_corr_target = df.corr()[target].abs()
10 low_corr_features = abs_corr_target[(abs_corr_target < 0.1) & (abs_corr_target.index != target)].
    index.tolist()
11 df_dropped_low = df.drop(columns=low_corr_features)
12
13 # Step 2: Identify and handle multicollinearity (high inter-feature correlation)
14 corr_remaining = df_dropped_low.corr()
15 upper_triangle = np.triu(np.ones(corr_remaining.shape), k=1).astype(bool)
16
17 # Extract high-correlation feature pairs (|correlation| > 0.45)
18 high_corr_pairs = (
19     corr_remaining.where(upper_triangle)
20     .stack()
21     .reset_index()
22 )
23 high_corr_pairs.columns = ['Feature 1', 'Feature 2', 'Correlation']
24 filtered_high_corr = high_corr_pairs[high_corr_pairs['Correlation'].abs() > 0.45]
25
26 # Drop the less important feature (based on correlation with target) from each high-correlation
    pair
27 drop_multicollinearity = set()
28 abs_corr_remaining_target = df_dropped_low.corr()[target].abs()
29

```

```

30 for _, row in filtered_high_corr.iterrows():
31     f1, f2 = row['Feature 1'], row['Feature 2']
32     if f1 != target and f2 != target:
33         corr_f1_target = abs_corr_remaining_target.get(f1, 0)
34         corr_f2_target = abs_corr_remaining_target.get(f2, 0)
35         if corr_f1_target < corr_f2_target:
36             drop_multicollinearity.add(f1)
37         else:
38             drop_multicollinearity.add(f2)
39
40 # --- Step 3: Final dataset after feature selection ---
41 df_final = df_dropped_low.drop(columns=list(drop_multicollinearity))
42 df_final.to_csv('datasets/datasetF.csv', index=False)
43
44 # Output summary
45 print("Low correlation features:", low_corr_features)
46 print("\nHigh correlation pairs (> 0.45):")
47 print(filtered_high_corr.sort_values(by='Correlation', ascending=False))
48 print("\nFeatures dropped for multicollinearity:", list(drop_multicollinearity))
49 print("\nFinal features:", df_final.columns.tolist())

```

## 4 Algorithm implementation

Next, we will implement the three chosen algorithms: MLP (Multilayer Perceptron), Autoencoder, and [Third Algorithm]. These algorithms have been selected due to their suitability for classification problems and their ability to handle structured and complex datasets. First, we will provide a brief introduction to each algorithm, highlighting their key characteristics. After that, we will proceed with the implementation and performance analysis for our dataset.

### 4.1 Multilayer Perceptron (MLP)

The **Multilayer Perceptron (MLP)** is a core model in the field of machine learning, used for both classification and regression tasks. MLP consists of three primary layers: the **input layer**, one or more **hidden layers**, and the **output layer**. Each layer contains multiple neurons, and neurons from one layer are fully connected to neurons in the next layer. This dense connectivity allows MLP to capture complex relationships in data, particularly non-linear ones, which are common in many real-world applications [1].

In the context of our dataset, which includes health-related features such as BMI, Age, Physical\_Health, and Sleep\_Time, MLP is well-suited to identify hidden patterns and make predictions about whether an individual has diabetes or not. The network learns from the data through a process called *backpropagation*, where errors from the output are propagated back through the network to adjust the weights of the connections between neurons. This process is combined with *gradient descent*, an optimization technique that minimizes the prediction error by adjusting the weights during each iteration of training. Over time, the model learns the optimal weights, improving its ability to make accurate predictions [2].

The basic structure of an MLP is illustrated in the figure below. The input features, such as BMI and Age, are fed into the input layer. From there, they pass through one or more hidden layers, where neurons transform the data by applying learned weights and activation functions. The transformed data then flows to the output layer, which produces the final prediction, such as whether an individual has diabetes.

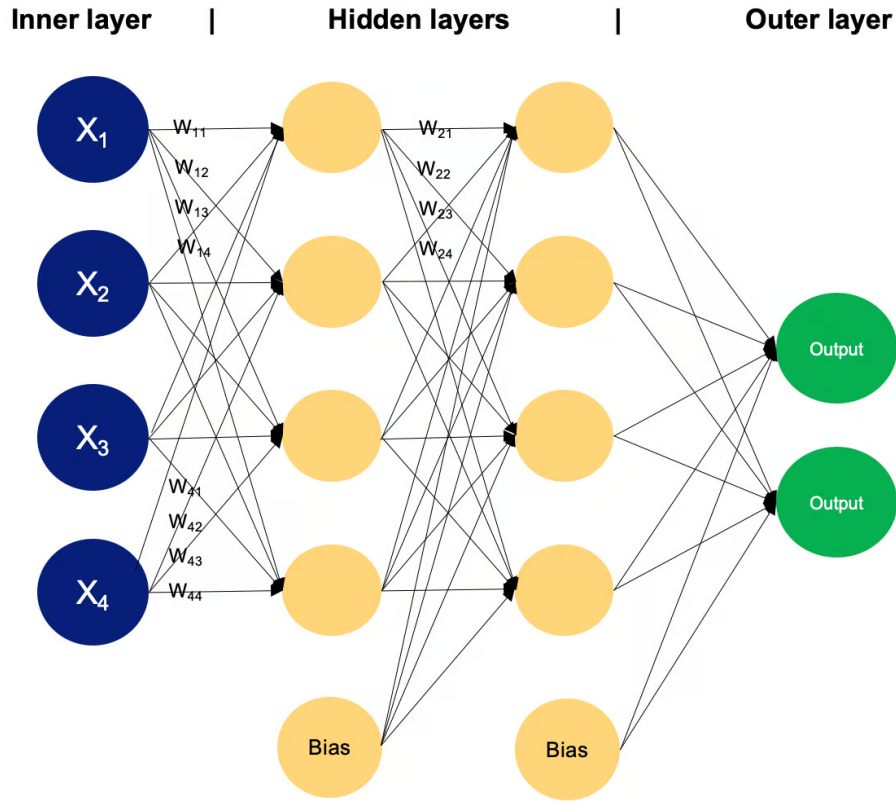


Figure 2: Basic structure of a Multilayer Perceptron (MLP) [3].

In the image, you can observe the flow of information from the input layer through the hidden layers to the output layer. Each layer plays a critical role in transforming the input data to make sense of the patterns and relationships. The neurons in the hidden layers learn features in the data, and as the network adjusts its weights through training, it becomes more accurate in predicting the target outcome.

MLPs are particularly effective in tasks where the relationships between input features and the target variable are non-linear, as is the case in health-related predictions like diabetes classification [1]. This model's ability to learn from complex and large datasets allows it to handle a wide variety of problems, from image recognition to medical diagnoses.

What makes MLPs particularly well-suited for our dataset is their capacity to model intricate relationships between multiple features simultaneously. By learning these relationships, MLPs can make accurate predictions, even when the data contains complex, interdependent factors. This makes MLP an ideal choice for predicting the likelihood of diabetes based on a variety of health indicators [3].

#### 4.1.1 Implementation and Evaluation

The MLP model was implemented using the `scikit-learn` library, leveraging two hidden layers with sizes 100 and 50 respectively. The activation function used was ReLU, and stochastic gradient descent was chosen as the optimizer with adaptive learning rate and a maximum of 500 iterations.

Before finalizing the model, we performed hyperparameter tuning using `GridSearchCV` to identify the best combination of parameters. This ensures that the model generalizes well to unseen data.

The best parameters found were:

```
{'mlp_activation': 'relu', 'mlp_alpha': 0.001, 'mlp_hidden_layer_sizes': (100, 50),
 'mlp_learning_rate': 'adaptive', 'mlp_max_iter': 500, 'mlp_solver': 'sgd'}
```

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 from sklearn.model_selection import train_test_split, cross_validate
6 from sklearn.neural_network import MLPClassifier
7 from sklearn.metrics import (
8     accuracy_score, precision_score, recall_score,
9     f1_score, confusion_matrix, ConfusionMatrixDisplay,
10    classification_report
11 )
12
13 # 1. Load data
14 df = pd.read_csv("datasets/diabetes_scaled.csv")
15
16 # 2. Features & Labels
17 X = df.drop("Diabetes_binary", axis=1)
18 y = df["Diabetes_binary"]
19
20 # 3. Train/Test
21 X_train = pd.read_csv("datasets/X_train.csv")
22 X_test = pd.read_csv("datasets/X_test.csv")
23 y_train = pd.read_csv("datasets/y_train.csv")
24 y_test = pd.read_csv("datasets/y_test.csv")
25
26 # 4. Define MLP model
27 mlp = MLPClassifier(
28     hidden_layer_sizes=(100, 50),
29     activation='relu',
30     solver='sgd',
31     alpha=0.001,
32     learning_rate='adaptive',
33     max_iter=500,
34     random_state=42
35 )
36
37 # 5. Train model
38 mlp.fit(X_train, y_train)
39
40 # 6. Predict on test set
41 y_pred = mlp.predict(X_test)
42
43 # Plot loss curve for full feature set
44 plt.figure(figsize=(8, 4))
45 plt.plot(mlp.loss_curve_)
46 plt.title("Loss Curve (Full Feature Set)")
47 plt.xlabel("Iterations")
48 plt.ylabel("Loss")
49 plt.grid(True)
50 plt.show()
51
52 # 7. Evaluate metrics
53 print("Performance on Full Feature Set:")
54 print(f"Accuracy: {accuracy_score(y_test, y_pred):.4f}")
55 print(f"Precision: {precision_score(y_test, y_pred):.4f}")
56 print(f"Recall: {recall_score(y_test, y_pred):.4f}")
57 print(f"F1 Score: {f1_score(y_test, y_pred):.4f}")
58
59 # 8. Confusion matrix visualization
60 cm = confusion_matrix(y_test, y_pred)
61 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=mlp.classes_)
62 disp.plot(cmap=plt.cm.Blues)

```

```

63 plt.title("Confusion Matrix (Full Feature Set)")
64 plt.show()
65
66 # 9. Cross-validation without plot
67 cv_results = cross_validate(
68     mlp, X_train, y_train,
69     cv=5,
70     scoring=["accuracy", "precision", "recall", "f1"],
71     return_train_score=True
72 )
73
74 print("\nCross-Validation Mean Scores (5-Fold):")
75 for metric in ["accuracy", "precision", "recall", "f1"]:
76     test_score = np.mean(cv_results[f'test_{metric}'])
77     train_score = np.mean(cv_results[f'train_{metric}'])
78     print(f"{metric.capitalize()}: Train = {train_score:.4f}, Test = {test_score:.4f}")
79
80 print("\nDetailed Classification Report (Full Feature Set):")
81 print(classification_report(y_test, y_pred, target_names=["No Diabetes", "Diabetes"]))
82
83 # === Performance on Selected Features + Interaction Terms ===
84
85 # Select features of interest + target
86 df_interact = df[["GenHlth", "PhysHlth", "DiffWalk", "Education", "Income", "Diabetes_binary"]].
87     copy()
88
89 # Add interaction features
90 df_interact["GenHlth_x_PhysHlth"] = df_interact["GenHlth"] * df_interact["PhysHlth"]
91 df_interact["PhysHlth_x_DiffWalk"] = df_interact["PhysHlth"] * df_interact["DiffWalk"]
92 df_interact["GenHlth_x_DiffWalk"] = df_interact["GenHlth"] * df_interact["DiffWalk"]
93 df_interact["Education_x_Income"] = df_interact["Education"] * df_interact["Income"]
94
95 # Prepare features and labels
96 X_interact = df_interact.drop("Diabetes_binary", axis=1)
97 y_interact = df_interact["Diabetes_binary"]
98
99 # Train/Test split
100 X_train_i, X_test_i, y_train_i, y_test_i = train_test_split(
101     X_interact, y_interact, test_size=0.2, stratify=y_interact, random_state=42
102 )
103
104 # Define MLP model
105 mlp_interact = MLPClassifier(
106     hidden_layer_sizes=(100, 50),
107     activation='relu',
108     solver='sgd',
109     alpha=0.001,
110     learning_rate='adaptive',
111     max_iter=500,
112     random_state=42
113 )
114
115 # Train model
116 mlp_interact.fit(X_train_i, y_train_i)
117
118 # Predict on test set
119 y_pred_i = mlp_interact.predict(X_test_i)
120
121 # Plot loss curve for selected features + interaction terms
122 plt.figure(figsize=(8, 4))
123 plt.plot(mlp_interact.loss_curve_)
124 plt.title("Loss Curve (Interactions Model)")

```

```

124 plt.xlabel("Iterations")
125 plt.ylabel("Loss")
126 plt.grid(True)
127 plt.show()
128
129 # Evaluate metrics
130 print("\nPerformance on Selected Features + Interactions:")
131 print(f"Accuracy: {accuracy_score(y_test_i, y_pred_i):.4f}")
132 print(f"Precision: {precision_score(y_test_i, y_pred_i):.4f}")
133 print(f"Recall: {recall_score(y_test_i, y_pred_i):.4f}")
134 print(f"F1 Score: {f1_score(y_test_i, y_pred_i):.4f}")
135
136 # Confusion matrix visualization
137 cm_i = confusion_matrix(y_test_i, y_pred_i)
138 disp_i = ConfusionMatrixDisplay(confusion_matrix=cm_i, display_labels=mlp_interact.classes_)
139 disp_i.plot(cmap=plt.cm.Greens)
140 plt.title("Confusion Matrix (Interactions Model)")
141 plt.show()
142
143 print("\nDetailed Classification Report (Interactions Model):")
144 print(classification_report(y_test_i, y_pred_i, target_names=["No Diabetes", "Diabetes"]))

```

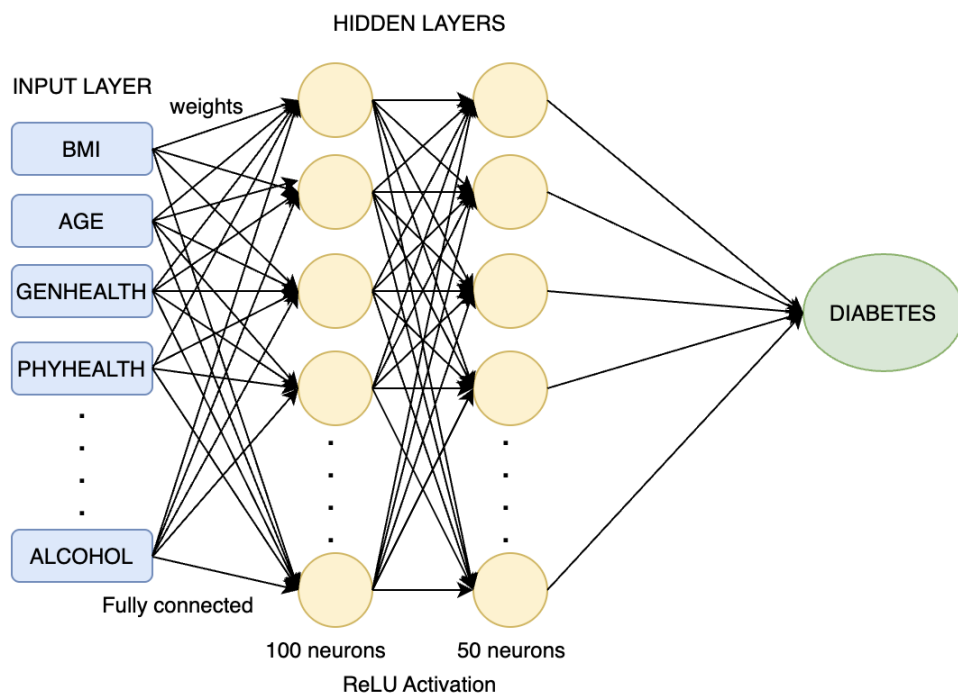


Figure 3: Our structure of a Multilayer Perceptron (MLP).

Figure 3 illustrates the architecture of the Multilayer Perceptron used in our implementation. It includes an input layer that receives features such as BMI, Age, and GenHealth, followed by two hidden layers with 100 and 50 neurons respectively. These layers apply ReLU activation and are fully connected. The final layer outputs a binary prediction indicating whether the individual is diabetic or not.

### 4.1.2 Results on Full Feature Set

Performance on Full Feature Set:

Accuracy: 0.7526  
Precision: 0.7362  
Recall: 0.8011  
F1 Score: 0.7673

This baseline performance shows that the MLP model handles the full feature set well, achieving over 75% accuracy and strong recall (80.11%). Precision and F1-score are also relatively balanced, which means the model is both identifying and distinguishing cases effectively.

Next, we look at cross-validation results to assess the model's ability to generalize across different data splits:

Cross-Validation Mean Scores (5-Fold):

Accuracy: Train = 0.7558, Test = 0.7484  
Precision: Train = 0.7359, Test = 0.7295  
Recall: Train = 0.8113, Test = 0.8040  
F1: Train = 0.7718, Test = 0.7649

These results confirm that the model generalizes well, as train and test scores remain consistent across folds. Importantly, test recall remains high (80.40%), reinforcing the model's reliability in correctly identifying diabetic cases in unseen data. This cross-validation step validates the single-run results and supports the model's robustness.

Now we present the full classification report to explore class-wise performance in more detail:

Detailed Classification Report (Full Feature Set):

	precision	recall	f1-score	support
No Diabetes	0.77	0.70	0.74	6753
Diabetes	0.74	0.80	0.77	7000
accuracy			0.75	13753
macro avg	0.75	0.75	0.75	13753
weighted avg	0.75	0.75	0.75	13753

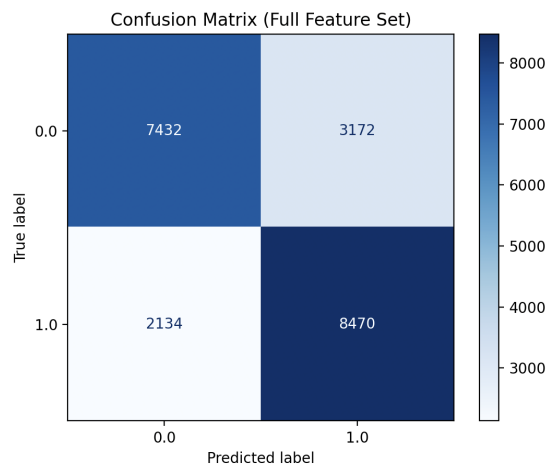


Figure 4: Confusion Matrix for MLP on Full Feature Set



From this report, it's evident that the model performs slightly better at detecting positive (diabetic) cases than negative ones, with a higher recall for the "Diabetes" class (80%). This is desirable in our context: it is far more critical to minimize false negatives (undiagnosed diabetic patients) than false positives. Therefore, recall is the metric we prioritize most.

### 4.1.3 Results on Selected Features with Interaction Terms

Performance on Selected Features + Interactions:

Accuracy: 0.6917

Precision: 0.6633

Recall: 0.8007

F1 Score: 0.7256

While accuracy and precision have dropped in this simplified model, recall has remained nearly identical (80.07%) to the full-feature model. This suggests that the reduced feature set is still effective in identifying diabetic patients, even if it leads to more false positives (lower precision).

We now show the detailed classification report for this model:

Detailed Classification Report (Interactions Model):				
	precision	recall	f1-score	support
No Diabetes	0.74	0.58	0.65	6753
Diabetes	0.66	0.80	0.73	7000
accuracy			0.69	13753
macro avg	0.70	0.69	0.69	13753
weighted avg	0.70	0.69	0.69	13753

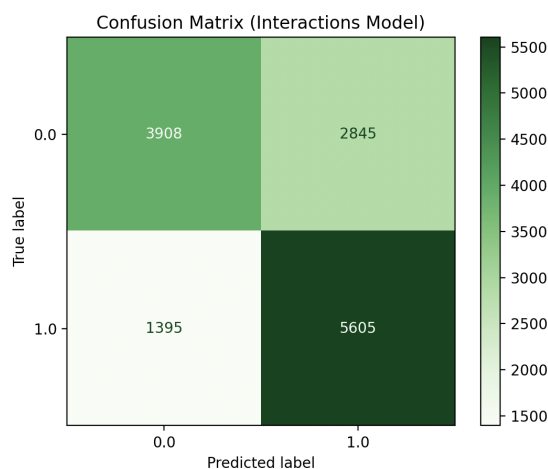


Figure 5: Confusion Matrix for MLP on Selected Features + Interaction Terms

Here again, recall for diabetic cases remains at 80%, but performance for non-diabetic predictions drops noticeably (recall of 58%). In practice, this means more people without diabetes are misclassified as diabetic. While this may raise false alarms, it is an acceptable trade-off in medical diagnostics where missing true cases is more dangerous than flagging potential ones.

In both models, recall is prioritized above all, aligning with our goal of reducing undiagnosed diabetic cases. The full-feature model is preferred for balanced performance, but the interaction-based model still holds value in more constrained environments.

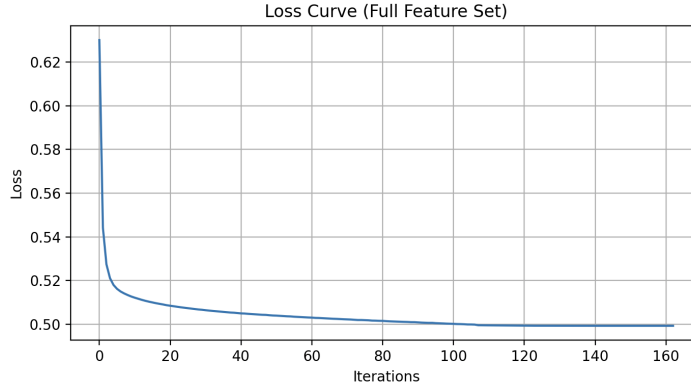


Figure 6: Training loss curve for the MLP model over 500 iterations.

The loss curve for the diabetes dataset shows a rapid initial decrease followed by convergence, indicating stable training and successful optimization.

## 4.2 Autoencoder

## 4.3 Naive Bayes

The Naive Bayes algorithm offers a straightforward yet often surprisingly effective approach to supervised learning, particularly for classification tasks like predicting diabetes status. At its core, it leverages Bayes' Theorem, a fundamental concept in probability that allows us to update the probability of an event based on new evidence. What makes Naive Bayes "naive" is its key simplifying assumption: it treats every feature in the dataset as if it were completely independent of all other features, given the class label (in our case, whether an individual has diabetes or not). The algorithm works by learning the conditional probability of each feature value for each class. For instance, it calculates the probability of having a certain BMI range given that the individual has diabetes, and the probability of having the same BMI range given that they do not have diabetes. It does this for all features in our dataset. When presented with a new individual's health profile, Naive Bayes uses these learned probabilities, along with the overall prevalence of diabetes in the training data (the prior probability), to calculate the probability of that individual belonging to each class (diabetes or no diabetes). The class with the highest calculated probability is then assigned as the prediction. The mathematical foundation lies in Bayes' Theorem:

$$P(\text{Class}|\text{Features}) = \frac{P(\text{Class}) \times P(\text{Features}|\text{Class})}{P(\text{Features})}$$

The "naive" independence assumption simplifies the term  $P(\text{Features}|\text{Class})$  into the product of the individual feature probabilities:

$$P(x_1, x_2, \dots, x_n|\text{Class}) = P(x_1|\text{Class}) \times P(x_2|\text{Class}) \times \dots \times P(x_n|\text{Class})$$

For our diabetes dataset, Naive Bayes presents a compelling choice due to its speed, ease of implementation, and its ability to effectively handle datasets with numerous features. Although the health indicators in the dataset are unlikely to be completely independent—a key assumption of Naive Bayes—the algorithm has demonstrated strong empirical success across various domains. This suggests that it can still provide valuable insights and deliver reliable predictive performance in identifying individuals at risk of diabetes.

There are several variants of Naive Bayes classifiers, each suited to different types of data. Gaussian Naive Bayes is commonly used when features are continuous and assumed to follow a normal distribution. Multinomial Naive Bayes typically applies to discrete features, while Bernoulli Naive Bayes is designed for binary or boolean features. Given that our dataset contains both continuous variables such as BMI and Age, as well as binary indicators like HighBP and Smoker status, Gaussian Naive Bayes is the most appropriate choice. This model assumes that the continuous features are normally distributed within each class, which aligns well with the nature of our data.

Overall, the simplicity of Naive Bayes makes it an excellent baseline model against which to compare more complex algorithms like Multilayer Perceptron (MLP) and Autoencoder. While it may not capture intricate interactions between features as effectively as MLP, Naive Bayes remains a solid and interpretable choice for predicting diabetes risk, offering a good balance between performance and computational efficiency.

### 4.3.1 Implementation and Evaluation

For the implementation of the Naive Bayes classifier, we utilized the Gaussian Naive Bayes variant due to the presence of both continuous and binary features in the dataset. The model was trained on the preprocessed diabetes dataset, where continuous features such as BMI and Age were assumed to follow a Gaussian distribution within each class. Binary features, including indicators like HighBP and Smoker, were also incorporated directly without additional transformation.

To ensure consistency in evaluation, we trained the model using the same training and testing sets applied in all algorithms we used. Additionally, we used 5-fold cross-validation to assess the model's generalization ability.

To evaluate the performance of the Naive Bayes classifier, several metrics were employed, including accuracy, precision, recall, and the F1-score. Given the medical context and the importance of correctly identifying positive cases, special emphasis was placed on recall and F1-score to assess the model's ability to detect individuals with diabetes effectively.

Below is the Python code used for training and evaluating the Gaussian Naive Bayes classifier:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 from sklearn.naive_bayes import GaussianNB
6 from sklearn.model_selection import cross_validate
7 from sklearn.metrics import (
8     accuracy_score, precision_score, recall_score,
9     f1_score, confusion_matrix, ConfusionMatrixDisplay,
10    classification_report
11 )
12
13 # Load training and test sets
14 X_train = pd.read_csv("datasets/X_train.csv")
15 X_test = pd.read_csv("datasets/X_test.csv")
16 y_train = pd.read_csv("datasets/y_train.csv").values.ravel()
17 y_test = pd.read_csv("datasets/y_test.csv").values.ravel()
18
19 # Define and train the Gaussian Naive Bayes model
20 gnb = GaussianNB()
21 gnb.fit(X_train, y_train)
22
23 # Predict on the test set
24 y_pred = gnb.predict(X_test)
25
26 # Evaluate performance metrics
27 print("Performance on Full Feature Set (Naive Bayes):")
28 print(f"Accuracy: {accuracy_score(y_test, y_pred):.4f}")
29 print(f"Precision: {precision_score(y_test, y_pred):.4f}")
30 print(f"Recall: {recall_score(y_test, y_pred):.4f}")
31 print(f"F1 Score: {f1_score(y_test, y_pred):.4f}")
```

```

32
33 # Confusion matrix visualization
34 cm = confusion_matrix(y_test, y_pred)
35 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=gnb.classes_)
36 disp.plot(cmap=plt.cm.Oranges)
37 plt.title("Confusion Matrix (Naive Bayes - Full Feature Set)")
38 plt.show()
39
40 # 5-Fold Cross-validation
41 cv_results = cross_validate(
42     gnb, X_train, y_train,
43     cv=5,
44     scoring=["accuracy", "precision", "recall", "f1"],
45     return_train_score=True
46 )
47
48 print("\nCross-Validation Mean Scores (10-Fold):")
49 for metric in ["accuracy", "precision", "recall", "f1"]:
50     test_score = np.mean(cv_results[f'test_{metric}'])
51     train_score = np.mean(cv_results[f'train_{metric}'])
52     print(f"{metric.capitalize()}: Train = {train_score:.4f}, Test = {test_score:.4f}")
53
54 print("\nDetailed Classification Report (Naive Bayes - Full Feature Set):")
55 print(classification_report(y_test, y_pred, target_names=["No Diabetes", "Diabetes"]))

```

### 4.3.2 Results on Full Feature Set (Naive Bayes)

Accuracy: 0.7381  
 Precision: 0.7300  
 Recall: 0.7558  
 F1 Score: 0.7427

This initial evaluation of the Gaussian Naive Bayes model on the full feature set yields an accuracy of 0.7381. The model demonstrates a precision of 0.7300, indicating that when it predicts "Diabetes," it is correct approximately 73% of the time. The recall of 0.7558 suggests that the model successfully identifies about 75.6% of the actual diabetic cases. The F1-score, which balances precision and recall, is 0.7427.

#### Detailed Classification Report (Naive Bayes):

	precision	recall	f1-score	support
No Diabetes	0.76	0.66	0.71	6753
Diabetes	0.71	0.81	0.76	7000
accuracy			0.73	13753
macro avg	0.74	0.74	0.73	13753
weighted avg	0.73	0.73	0.73	13753

The classification report offers a more granular view of the model's performance for each class. For the "No Diabetes" class, the precision is 0.76 and the recall is 0.66. This means that 76% of the predictions for "No Diabetes" were correct, and 66% of all actual "No Diabetes" cases were identified.

For the "Diabetes" class, the precision is 0.71, and the recall is 0.81. This indicates that 71% of the "Diabetes" predictions were correct, and 81% of all actual "Diabetes" cases were identified.

The overall accuracy is 0.73, consistent with the initial performance evaluation. The macro and weighted averages for precision, recall, and F1-score provide a consolidated view of the model's effectiveness across both classes. The

higher recall for the "Diabetes" class suggests that the model is better at identifying positive cases, which is often a crucial consideration in medical applications.

## 5 Result interpretation

### References

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. Available: <https://www.deeplearningbook.org>
- [2] Michael Nielsen. *Neural Networks and Deep Learning*. 2015. Available: <http://neuralnetworksanddeeplearning.com>
- [3] *Multilayer Perceptrons in Machine Learning*. <https://www.datacamp.com/tutorial/multilayer-perceptrons-in-machine-learning>
- [4] <https://hex.tech/blog/autoencoders-for-feature-selection/>