# Data Analysis for Diabetes Prediction Using Health Indicators

Machine Learning

# Introduction

➔ The dataset contains detailed health indicators collected from a large population and is designed to support the analysis and prediction of diabetes risk.

➔ The main variable of interest is Diabetes binary, which indicates whether an individual has been diagnosed with diabetes (1) or not (0).

➔ Motivation:
- 1. How can machine learning models be used to predict diabetes risk more accurately?
- Can a subset of factors be used to accurately predict whether an individual has diabetes?

# Selection of algorithms for the dataset

➔ **MLP (Multilayer Perceptron)**
- ◆ Neural network capable of complex non-linear relationships
- ◆ Requires tuning of hidden layers and neurons
- ◆ Handles high-dimensional data well

➔ **Autoencoder + Gradient Boosting Classifier**
- ◆ Autoencoder reduces dimensionality and extracts features
- ◆ Gradient Boosting builds strong predictors from decision trees
- ◆ Combined approach enhances performance for complex datasets

➔ **Naive Bayes**
- ◆ Probabilistic classifier based on Bayes' theorem
- ◆ Computationally efficient and works well with limited data
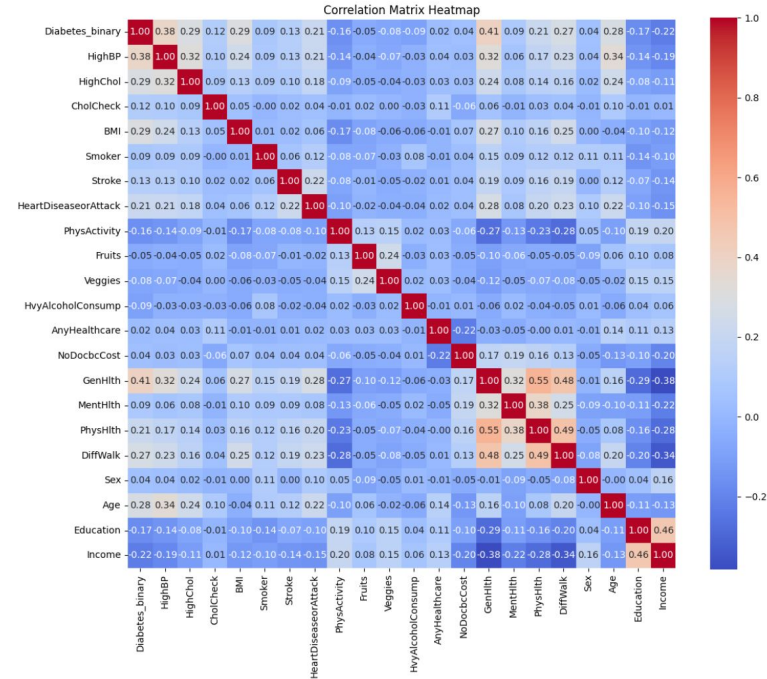- ◆ Assumes feature independence (may be limiting)

# Preprocessing - Data cleaning

➔ Found 1635 **duplicate rows** , but they are not removed to ensure we preserve all relevant data for analysis.

➔ IQR method flagged 40,205 records as containing at least one **outlier** attribute—more than half of the dataset.

➔ Z-score method was far more reasonable giving 1927 outliers. We removed them, that left our dataset with 68,765 remaining records for training and testing

```
1  ####### Summary of the dataframe #######
2                           Min    Max        Mean    Mode
3  Diabetes_binary          0.0    1.0    0.500000     0.0
4  HighBP                   0.0    1.0    0.563458     1.0
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
```
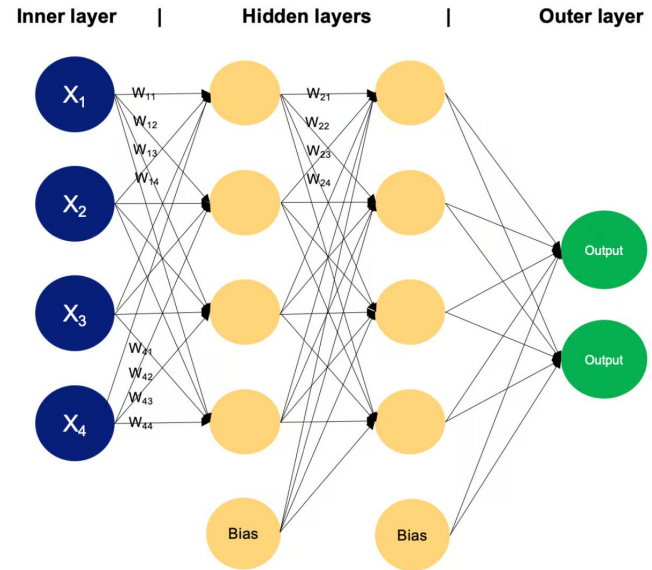
# Preprocessing - Correlation

➔ **Data integration:** all the required data is already combined into a single CSV file.

➔ **Data transformation:** we apply feature scaling to the dataset in order to standardize the features.
  ◆ StandardScaler, which standardizes the features by removing the mean and scaling to unit variance.

➔ **Correlation between features:**
  ◆ Removal of weakly correlated features
  ◆ Mitigation of multicollinearity
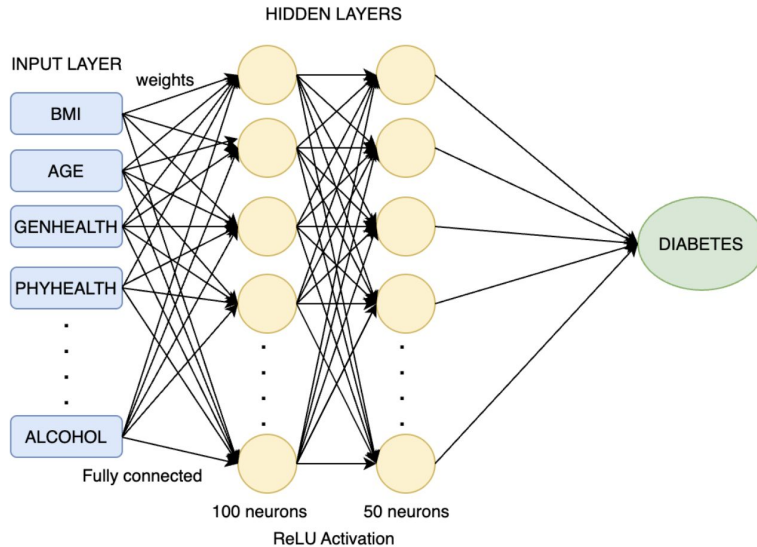


Correlation Matrix Heatmap

# Multilayer Perceptron - Theory

➔ MLP consists of three primary layers: the input layer, one or more hidden layers, and the output layer.
➔ 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.



**Inner layer** | **Hidden layers** | **Outer layer**

$X_1$ $X_2$ $X_3$ $X_4$

$W_{11}$ $W_{12}$ $W_{13}$ $W_{14}$

$W_{21}$ $W_{22}$ $W_{23}$ $W_{24}$

$W_{41}$ $W_{42}$ $W_{43}$ $W_{44}$

Output
Output

Bias  Bias

# Multilayer Perceptron - Implementation



```python
# 4. Define MLP model
mlp = MLPClassifier(
    hidden_layer_sizes=(100, 50),
    activation='relu',
    solver='sgd',
    alpha=0.001,
    learning_rate='adaptive',
    max_iter=500,
    random_state=42
)

# 5. Train model
mlp.fit(X_train, y_train)

# 6. Predict on test set
y_pred = mlp.predict(X_test)

# Plot loss curve for full feature set
plt.figure(figsize=(8, 4))
plt.plot(mlp.loss_curve_)
plt.title("Loss Curve (Full Feature Set)")
plt.xlabel("Iterations")
plt.ylabel("Loss")
plt.grid(True)
```
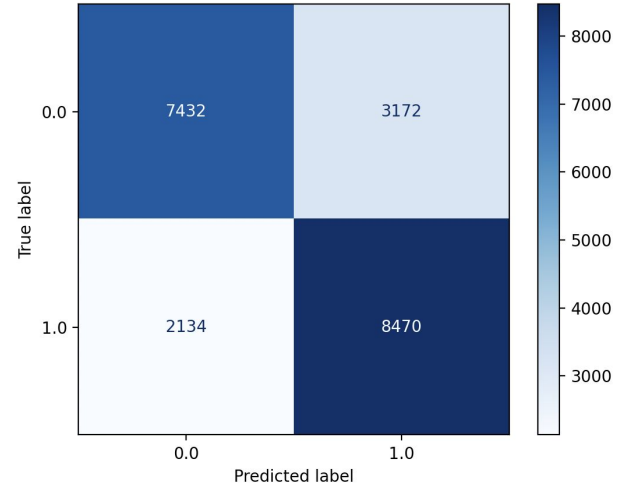
# Multilayer Perceptron - Results

```
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
```
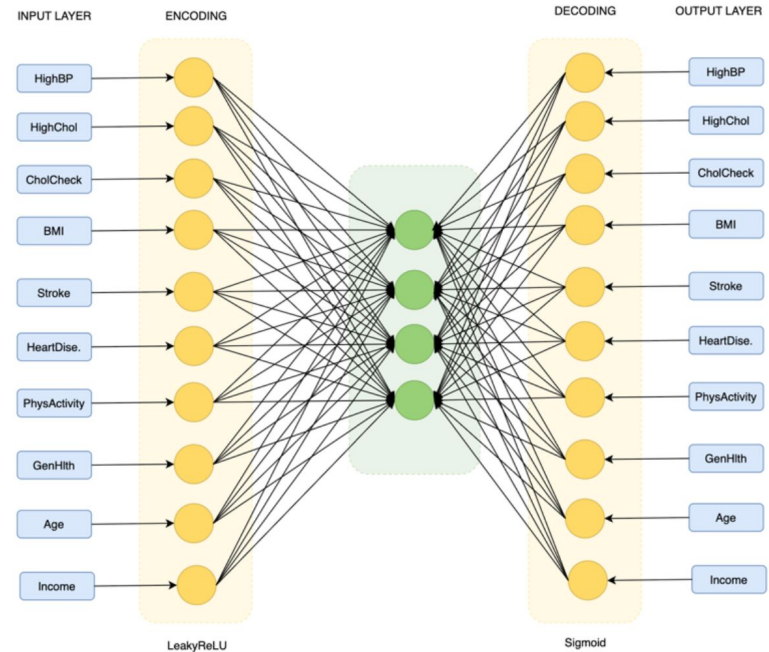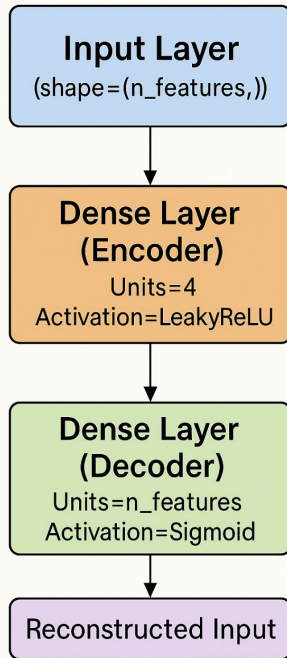


Confusion Matrix (Full Feature Set)

# Autoencoder - Theory

➔ Autoencoders are neural networks designed to learn efficient representations of data by compressing input through an encoder, capturing essential features in a compact bottleneck layer, and reconstructing the original data using a decoder.

# Autoencoder - Implementation



```python
class Autoencoder:
    def __init__(self, encoding_dim, activation1, activation2, optimizer1, optimizer2):
        self.input_dim = X_train.shape[1]
        self.encoding_dim = encoding_dim

        input_layer = keras.layers.Input(shape=(self.input_dim,))
        encoder = keras.layers.Dense(encoding_dim, activation1)(input_layer)
        decoder = keras.layers.LeakyReLU(encoder)

        self.autoencoder = keras.Model(inputs=input_layer, outputs=decoder)
        self.encoder_model = keras.Model(inputs=input_layer, outputs=encoder)

        self.autoencoder.compile(optimizer=optimizer1, loss=optimizer2)
        self.autoencoder.summary()

    def train(self, X_train, X_test, epochs, batch_size=32):
        self.autoencoder.fit(
            X_train, X_train,
            epochs=epochs,
            batch_size=batch_size,
            shuffle=True,
            validation_data=(X_test, X_test)
        )

    def encode(self, X):
        return self.encoder_model.predict(X)
```
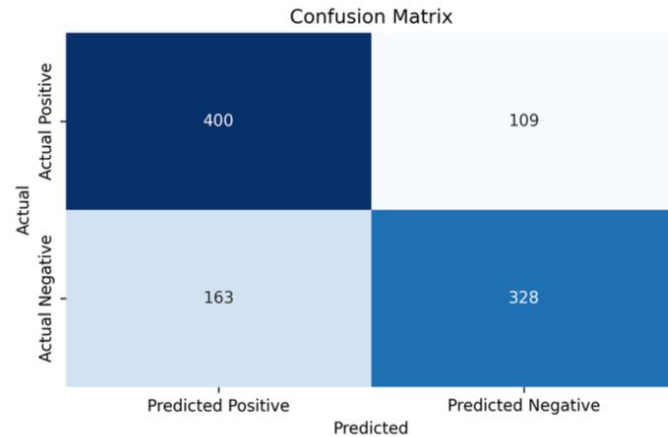
# Autoencoder - Results

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0.0 | 0.77 | 0.68 | 0.72 | 10604 |
| 1.0 | 0.71 | 0.79 | 0.75 | 10604 |
| **Accuracy** | | 0.74 | | 21208 |
| **Macro Avg** | 0.74 | 0.74 | 0.74 | 21208 |
| **Weighted Avg** | 0.74 | 0.74 | 0.74 | 21208 |



Confusion Matrix

# Naive Bayes - Theory

➤ The Naive Bayes algorithm offers a straightforward yet often surprisingly effective approach to supervised learning, particularly for classification tasks like predicting diabetes status.

➤ Variants:
   ◆ Gaussian Naive Bayes
   ◆ Multinomial Naive Bayes
   ◆ Bernoulli Naive Bayes

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

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

# Naive Bayes - Implementation

```python
# Load training and test sets
X_train = pd.read_csv("datasets/X_train.csv")
X_test = pd.read_csv("datasets/X_test.csv")
y_train = pd.read_csv("datasets/y_train.csv").values.ravel()
y_test = pd.read_csv("datasets/y_test.csv").values.ravel()

# Define and train the Gaussian Naive Bayes model
gnb = GaussianNB()
gnb.fit(X_train, y_train)

# Predict on the test set
y_pred = gnb.predict(X_test)
```

```python
# Evaluate performance metrics
print("Performance on Full Feature Set (Naive Bayes):")
print(f"Accuracy: {accuracy_score(y_test, y_pred):.4f}")
print(f"Precision: {precision_score(y_test, y_pred):.4f}")
print(f"Recall: {recall_score(y_test, y_pred):.4f}")
print(f"F1 Score: {f1_score(y_test, y_pred):.4f}")

# Confusion matrix visualization
cm = confusion_matrix(y_test, y_pred)
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=gnb.classes_)
disp.plot(cmap=plt.cm.Oranges)
plt.title("Confusion Matrix (Naive Bayes - Full Feature Set)")
plt.show()

# 5-Fold Cross-validation
cv_results = cross_validate(
    gnb, X_train, y_train,
    cv=5,
    scoring=["accuracy", "precision", "recall", "f1"],
    return_train_score=True
)
```
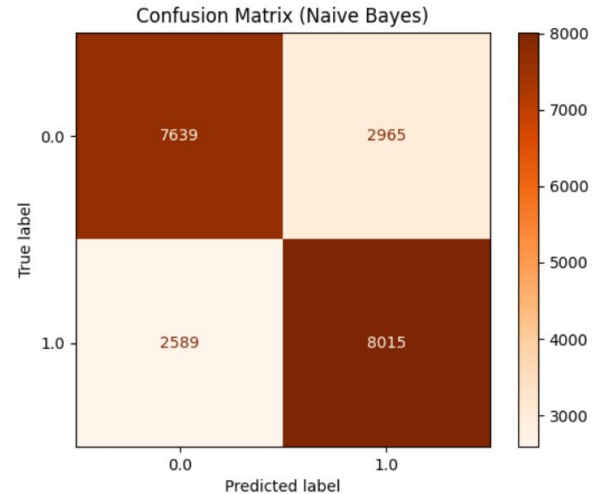
# Naive Bayes - Results

**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   |

Confusion Matrix (Naive Bayes)

|            | Predicted 0.0 | Predicted 1.0 |
|------------|---------------|---------------|
| True 0.0   | 7639          | 2965          |
| True 1.0   | 2589          | 8015          |

# Results & Reasoning

➔ We compare the three applied algorithms—Multilayer Perceptron (MLP), Autoencoder, and Gaussian Naive Bayes—using four key performance metrics: Accuracy, Precision, Recall, and F1 Score.

➔ In a medical setting, missing a true diabetic case (false negative) is far more critical than issuing a false alarm. Thus, we prioritize models with higher **recall**.

➔ Among the three models evaluated, the Multilayer Perceptron (MLP) demonstrates the best overall performance

.

| Algorithm | Parameters | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|---|
| MLP | hidden_layers=(100,50), activation=ReLU, solver=SGD | 0.7526 | 0.7362 | 0.8011 | 0.7673 |
| Autoencoder | encoder=(32,16), decoder=(16,32), activation=ReLU | 0.7400 | 0.7100 | 0.7900 | 0.7500 |
| Gaussian NB | default (GaussianNB) | 0.7381 | 0.7300 | 0.7558 | 0.7427 |

# Thank you!

## Questions?

Albana Rexhepi
Eljesa Kqiku
Kaltrina Kuka