

# Intro to Deep Learning - Final Project

Semester A 2025

Authors: Yigal Meshulam: יגאל משולם  
Yuval Wilf: יובל ווילף

## Title:

Machine & Deep Learning-Based Classification of Heart Disease Using Health Risk Indicators

## About this document:

**Todo:** add about the colab, reference between this document sections and colab cells.

## Keywords:

Deep learning, Heart disease prediction, Neural networks, Classification, Healthcare analytics

## Abstract:

Heart disease remains one of the leading causes of mortality worldwide, necessitating accurate and early prediction for effective prevention. This study employs deep learning techniques to classify heart disease status based on a dataset comprising various health indicators. The [Heart Disease](#) dataset, obtained from Kaggle, includes risk factors such as age, gender, blood pressure, cholesterol levels, smoking habits, and other health-related parameters. We developed and trained a neural network model based on health metrics to classify individuals as having or not having heart disease. The model's performance is evaluated using standard classification metrics such as accuracy, precision, recall, and F1-score. The results demonstrate the effectiveness of deep learning in identifying heart disease risk, providing insights for healthcare applications. The code implementation is available in a Colab notebook on GitHub: [Deep Learning Project](#).

## 1. Introduction

Heart disease is a major public health issue and a leading cause of death worldwide. Traditional risk assessment methods rely on statistical analysis and expert-driven evaluations, which may not fully capture complex interactions among risk factors. Deep learning has emerged as a powerful tool in healthcare analytics, offering improved prediction capabilities through automated feature extraction. In this study, we apply a deep learning model to classify heart disease based on a dataset containing multiple health risk indicators.

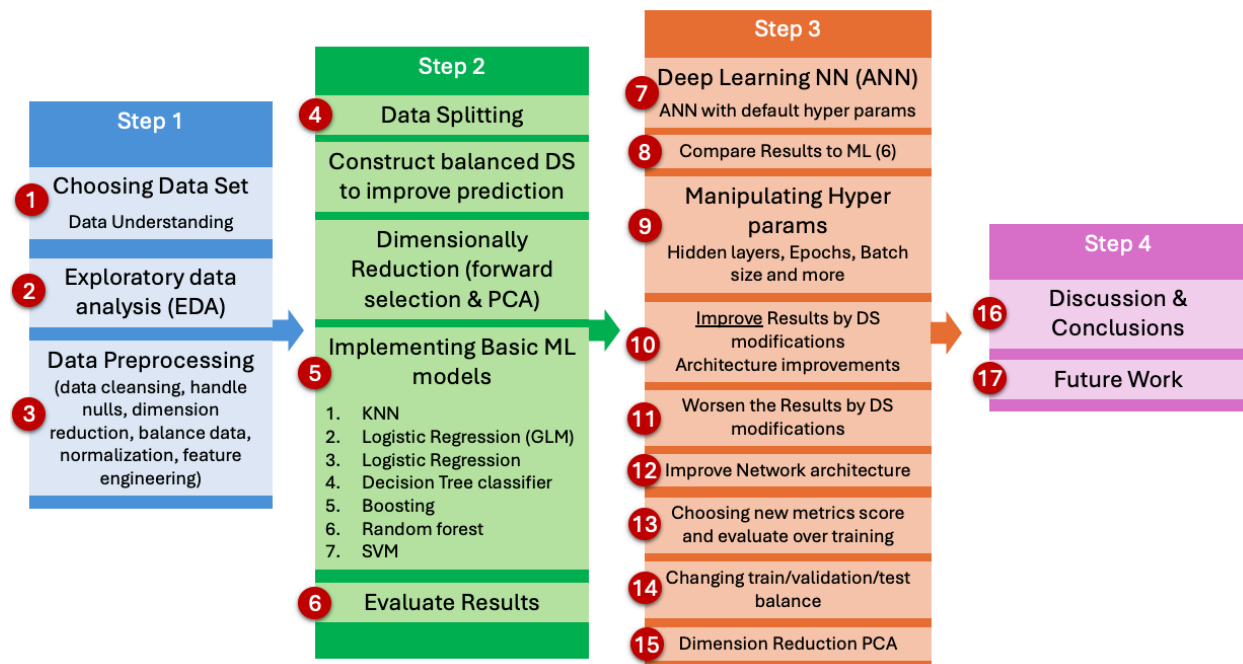
## 2. Literature Review

Recent advances in artificial intelligence (AI) and machine learning (ML) have significantly improved medical diagnosis and prediction systems. Various studies have explored logistic regression, decision trees, support vector machines (SVM), and ensemble models for heart disease prediction. However, deep learning models, particularly neural networks, have shown superior performance due to their ability to learn hierarchical patterns in complex datasets.

For example, Rajkomar et al. (2018) demonstrated the effectiveness of deep learning in medical diagnosis, highlighting its ability to detect complex interactions in patient data. Similarly, Esteva et al. (2017) applied deep learning in dermatology, showing improved classification performance compared to traditional methods. In the cardiovascular domain, Krittanawong et al. (2019) reviewed AI applications in cardiology, emphasizing the role of neural networks in predicting cardiovascular risk. Another study by Al'Aref et al. (2020) showcased the use of deep learning in coronary artery disease assessment, demonstrating high diagnostic accuracy. These studies reinforce the potential of deep learning in healthcare applications, particularly in heart disease classification.

## 3. Methodology

This project comprises four steps to determine whether we can predict whether a person will suffer from a heart disease. The steps are described in the following chart:



The following subsection will describe each step and discuss the results.

### 3.1. **Step 1** - Data Understanding and Processing

#### 3.1.1. Choosing a data set - data set description

3.1.1.1. The labeled [heart disease](#) dataset used in this study is sourced from Kaggle and contains multiple health indicators associated with heart disease risk. It includes demographic features (age, gender), physiological parameters (blood pressure, cholesterol, BMI), behavioral factors (smoking, alcohol consumption, exercise habits), and metabolic indicators (fasting blood sugar, triglyceride level, CRP level, homocysteine level). The target variable is "Heart Disease Status," which is labeled as "Yes" or "No."

3.1.1.2. Data Set Analysis - to start analyzing the dataset, we have to check the values of our features. First of all - our dataset contains both categorical and numeric values. For example, Exercise Habits values are [high, medium, low], Gender [Male, Female], while Blood pressure and BMI are numeric float numbers.

3.1.1.3. In the dataset, we found out that some values, such as "Alcohol consumption," contain "None". We need to read the CSV files into a data frame in such a way that those None values will be regarded as None and not NaN, or Null - the way to do this is using this code:

```
df = pd.read_csv('heart_disease.csv', na_values=["", "NA", "null"], keep_default_na=False)
```

### 3.1.2. Exploratory Data Analysis (EDA)

#### 3.1.2.1. Statistics: shape, numeric & categorical distribution

3.1.2.1.1. The purpose here is to understand the fundamental characteristics of each feature (including measures of central dependency and spread) - that will help identify potential anomalies or skew in both numeric and categorical variables.

```
Heart Disease dataset contain 10000 rows and 21 columns.

Number of numeric columns: 9
Number of categorical columns: 12
```

3.1.2.1.2. Shape and features & Target variable balance: we see imbalance in the target variable - 80% are labeled false (no heart disease), 20% true.

```
Heart Disease Status

No      0.8
Yes     0.2
```

3.1.2.1.3. Numeric features statistics - using the describe function, we retrieved statistics about the dataset's numeric values (age, blood pressure, cholesterol level, etc.)

```
Numeric Summary:

count  Age      Blood Pressure      Cholesterol Level      BMI \
mean   49.296259  149.757740      225.425577      29.077269
std    18.193970  17.572969      43.575809      6.307098
min    18.000000  120.000000      150.000000      18.002837
25%    34.000000  134.000000      187.000000      23.658075
50%    49.000000  150.000000      226.000000      29.079402
75%    65.000000  165.000000      263.000000      34.520015
max    80.000000  180.000000      300.000000      39.996954

count  Sleep Hours      Triglyceride Level      Fasting Blood Sugar      CRP Level \
mean   6.991329      9974.000000      9978.000000      9974.000000
std    1.753195      87.067226      23.584011      4.340248
min    4.000005      100.000000      80.000000      0.003647
25%    5.449866      176.000000      99.000000      3.674126
50%    7.003252      250.000000      120.000000      7.472164
75%    8.531577      326.000000      141.000000      11.255592
max    9.999952      400.000000      160.000000      14.997087

count  Homocysteine Level
mean   12.456271
std    4.323426
min    5.000236
25%    8.723334
50%    12.409395
75%    16.140564
max    19.999037
```

3.1.2.1.4. Categorical features statistics - we did the same for categorical features.

3.1.2.1.5. Summary and conclusions from dataset analysis: We can see that many categorical features are balanced around 50% besides the target variable. We can see that some fields are missing (do not sum up to 10000 rows) we need to take care of those rows. Interestingly, the mean of age is around 50. We have a maximum age of 80 even. The mean values of all numeric

```
Categorical Summary:

count  Gender      Exercise Habits      Smoking      Family Heart      Disease      Diabetes \
unique 2          3          2          2          2          2
top    Male      High      Yes      Yes      No      No
freq   5003      3372      5123      5004      5018

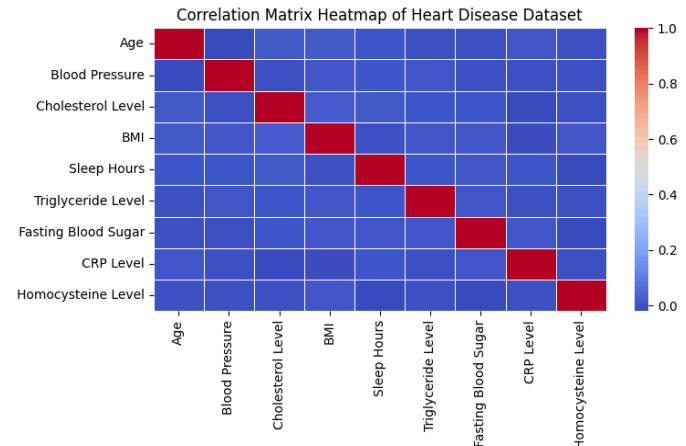
count  High Blood Pressure      Low HDL Cholesterol      High LDL Cholesterol \
unique 2          2          2
top    Yes      Yes      No
freq   5022      5000      5036

count  Alcohol Consumption      Stress Level      Sugar Consumption      Heart Disease      Status
unique 4          3          3          3          2
top    None      Medium      Low      No
freq   2554      3387      3390      8000
```

values have valuable information for this research.

### 3.1.2.1.6. **Correlations:** numeric vs. numeric - predictor vs. target

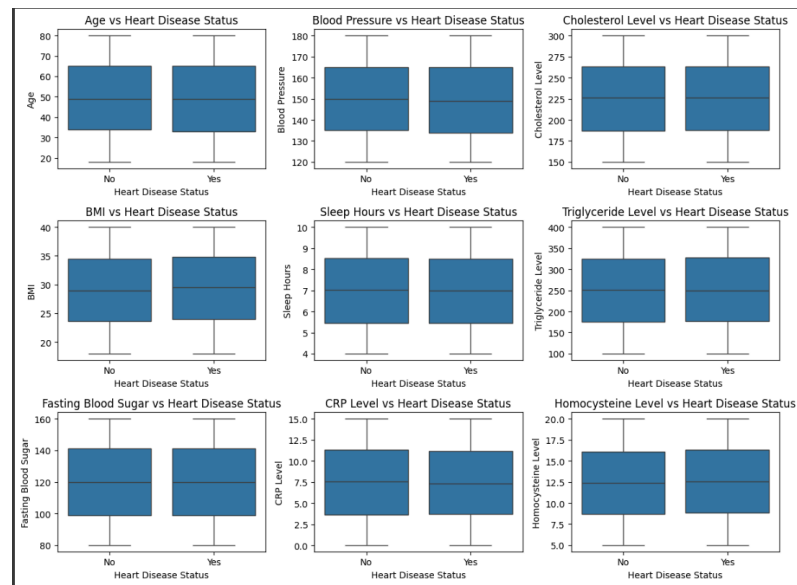
3.1.2.1.6.1. In this section, we examine how numeric predictors relate to each other and the target variable, indicating potential multicollinearity or strong relationships that can inform further analysis or feature engineering. We must determine whether categorical variables are significantly associated with heart disease status, providing a statistical basis for deciding if these features may be helpful to the predictors.



3.1.2.1.6.2. Findings: the graph shows that cross correlations between numeric features are close to 0. That means the multicollinearity is zero. There is no overlapping.

3.1.2.1.6.3. Checking correlation between numeric features and the target variable Heart disease status. We will draw boxplot graphs of each numeric value and its target variable.

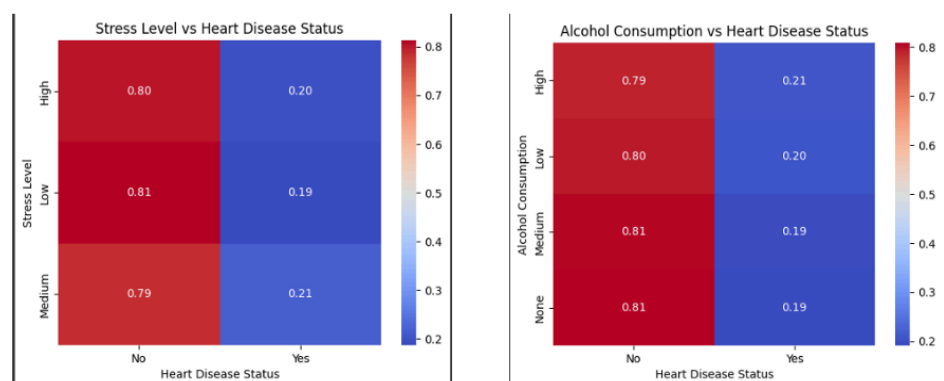
3.1.2.1.6.4. Findings: As seen from the charts, there are no correlations. [Todo: Why not use a heat map? Yigal? Hey : ) We are not looking for correlation in a sense of changing together. Instead, we try to find a clear "range of separation" between them both.]



### 3.1.2.1.7. Correlations: Categorical columns vs. Heart Disease Status - using (Chi-squared)

3.1.2.1.7.1. We examine a table of each predictor and its chi-square value and p-value. From the table, we can see that our *H<sub>0</sub>* hypothesis (there is no difference or correlation between each categorical variable and the target variable) is rejected for stress level, which means that stress level is the only categorical variable that influences the target variable.

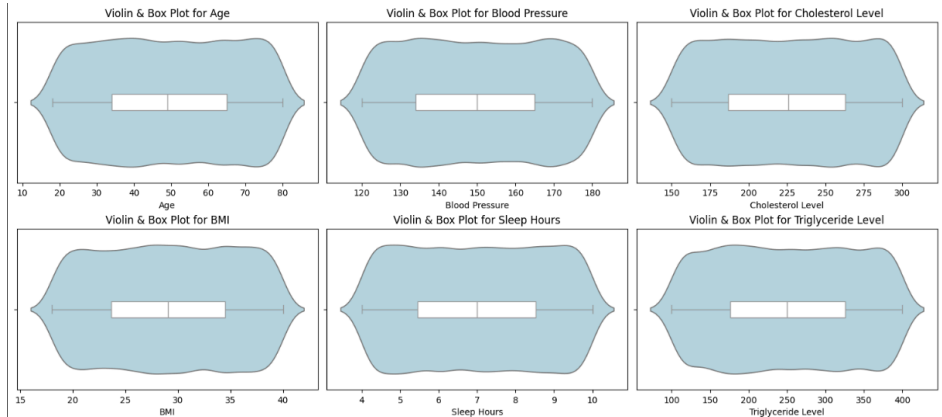
3.1.2.1.7.2. We plot the stress level chi-squared and the target variable and get the following:



3.1.2.1.7.3. **Conclusion:** the correlation between different stress levels [High, Medium, Low] and the target variable Heart disease is relatively low (for all types ~20%), the same goes for Alcohol consumption.

### 3.1.2.2. Features distribution

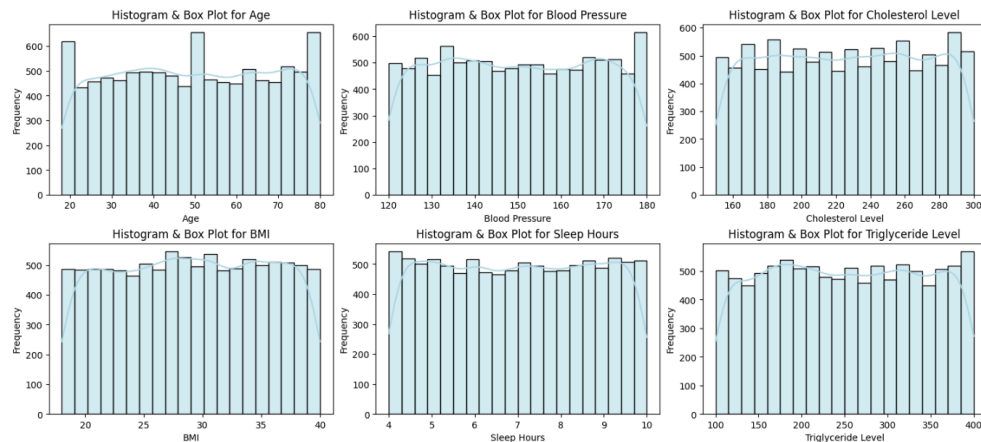
visualization: Violin, Boxplots - to visually inspect the distribution, spread and potential outliers of each numeric feature, providing an intuitive understanding of each variable's behaviour and detecting anomalies or skewed data.



#### 3.1.2.2.1. The graph (full

graph can be found in the colab notebook attached to this project) shows that all violins are similar in different features, showing the spread of data and outliers. If we take age, for example, the median is around 50, and the interquartile range (box) is between 35 and 65. We do not see potential outliers (dots outside the whiskers). The kernel density (the violin part) shape is more or less evenly spread. We do not see a normal distribution.

3.1.2.3. Features distribution - Histograms - it shows the frequency of individuals in different age ranges, indicating how age is distributed within the dataset. We can see here that data is evenly spread between different ranges.



### 3.1.3. Data Quality and Preparations: (Conclusions and actions from EDA process)

3.1.3.1. **Missing Values:** Both numeric and categorical predictors have missing values in varying amounts. Missingness must be addressed through imputation (mean/median for numeric, mode or "Missing" category for categorical) or exclusion, depending on the severity and distribution of missing data. Example: **Alcohol Consumption** has significant missing data (2,586 missing values out of 10,000), which could bias the analysis if not handled carefully.

- 3.1.3.2. **Outliers and Scaling:** Some numeric predictors, such as **CRP Level** and **Triglyceride Level**, exhibit wide ranges or extreme values that may act as outliers. These should be analyzed and, if needed, transformed or capped.
- 3.1.3.3. **Standardization or normalization** of numeric predictors will be necessary to ensure consistent scaling across features, especially for algorithms sensitive to magnitude differences (e.g., logistic regression)
- 3.1.3.4. **Class Imbalance** - The target variable, **Heart Disease Status**, has a significant imbalance (**No: 80%, Yes: 20%**). This imbalance must be addressed using techniques like resampling (oversampling or undersampling) or cost-sensitive learning.
- 3.1.3.5. **Predictor Insights and Potential Relationships with the target:**
  - 3.1.3.5.1. **Age:** Older age is likely positively correlated with heart disease risk.
  - 3.1.3.5.2. **Blood Pressure, Cholesterol Level, BMI:** These are known risk factors for cardiovascular health and show a wide range of values that could differentiate between individuals with and without heart disease.
  - 3.1.3.5.3. **CRP Level and Triglycerides:** Indicators of inflammation and lipid levels, respectively, are likely significant predictors.
  - 3.1.3.5.4. **Smoking, Exercise Habits, and Family Heart Disease:** Strongly tied to lifestyle and genetic risk factors, these are highly likely to predict heart disease.
  - 3.1.3.5.5. **High Blood Pressure, Diabetes, and Cholesterol Levels (HDL/LDL):** Directly linked to cardiovascular health, these are expected to be important predictors.
  - 3.1.3.5.6. Predictors like **Stress Level**, **Sugar Consumption**, and **Alcohol Consumption** may indirectly contribute to heart disease risk and offer actionable insights for prevention.
- 3.1.3.6. **Range & Distribution:** Predictors like **Sleep Hours** and **Fasting Blood Sugar** may require further analysis to understand their exact influence, as they seem to have narrower ranges and may exhibit non-linear effects.
- 3.1.3.7. **Features Engineering and Selection:**
  - 3.1.3.7.1. **Transformations** - Numeric predictors with non-linear relationships to the target (e.g., **Sleep Hours**, **BMI**) may require feature transformations (e.g., polynomial terms, logarithmic transformations). Ordinal categorical predictors (e.g., **Stress Level**, **Sugar Consumption**) should be encoded with order preserved to capture their progression.
  - 3.1.3.7.2. **Interactions** - Interaction terms between predictors (e.g., **Age × Exercise Habits** or **BMI × Smoking**) might reveal additional insights and improve predictive performance.
- 3.1.3.8. **Key Hypotheses:**
  - 3.1.3.8.1. Older individuals with higher blood pressure, cholesterol levels, and CRP levels are at a higher risk of heart disease.
  - 3.1.3.8.2. Smoking, low exercise, and high stress levels are likely to have strong associations with heart disease.
  - 3.1.3.8.3. Gender and family history of heart disease may introduce inherent differences in risk profiles.
  - 3.1.3.8.4. Lifestyle choices (e.g., alcohol consumption, sugar intake) influence heart disease risk, though they may be weaker predictors than clinical or genetic factors.



3.1.4.**Data Pre-Processing** (see exact methods in the colab notebook).

3.1.4.1. In the third section of step 1, based on the conclusion of EDA, we determined the extent and pattern of missing data across features. We assessed potential biases in the target variable. This guides whether dropping missing values or applying imputation techniques is most appropriate.

3.1.4.2. We prepare categorical variables for machine learning algorithms by converting them into suitable numeric representations, facilitating proper model training and improved predictive performance.

3.1.4.3. Handling Empty cells & Nulls (Drop and Imputation)

3.1.4.4. We examined the null distribution vs the target variable and dropping nulls.

3.1.4.5. We did not need to do any data imputation procedures.

3.1.4.6. Data transformations - We used transformations to convert categorical variables into numerical representations using one-hot encoding, ensuring compatibility with machine learning models. The drop\_first=True parameter prevents redundancy and multicollinearity while converting all values to float and maintaining data consistency for efficient computation. This step enhances model interpretability and performance by making the dataset fully numeric. We implemented another way to transform categorical features to numeric values using dummies.

3.2. **Step 2** - Implementing Basic ML models

3.2.1. Data Splitting

3.2.2. Construct balanced DS to improve predictions

3.2.3. Dimensional Reduction (forward selection & PCA)

3.2.4. Implementing Basic ML Models

3.2.4.1. KNN

3.2.4.2. Logistic Regression (GLM)

3.2.4.3. Logistic Regression

3.2.4.4. Decision Tree Classifier

3.2.4.5. Boosting

3.2.4.6. Random Forest

3.2.4.7. SVM

3.2.5. Evaluate Results

3.3. Step 3 - Deep Learning Neural Network Model (ANN)

3.4.A

3.5.A

4. Results

5. Discussion and Conclusions

6. Future Work

7. References

=====

### 3.2 Data Preprocessing

- **Handling Missing Values:** Any missing or inconsistent data points were addressed using imputation techniques.
- **Feature Encoding:** Categorical variables (e.g., smoking, gender, stress levels) were converted into numerical representations using one-hot encoding.
- **Normalization:** Continuous variables (e.g., cholesterol levels, BMI) were normalized to ensure uniform feature scaling.
- **Train-Test Split:** The dataset was split into training (80%) and testing (20%) sets to evaluate model performance.

### 3.3 Model Architecture

A deep neural network (DNN) was constructed with the following configuration:

- **Input Layer:** Accepts all selected features.
- **Hidden Layers:** Three fully connected layers with ReLU activation functions.
- **Dropout Layer:** Applied to prevent overfitting.
- **Output Layer:** A single neuron with a sigmoid activation function for binary classification.
- **Loss Function:** Binary cross-entropy.
- **Optimizer:** Adam optimizer.

### 3.4 Model Training

The model was trained using backpropagation with stochastic gradient descent (SGD). Hyperparameter tuning was performed to optimize learning rate, batch size, and number of epochs.

## 4. Results

The performance of the deep learning model was assessed using:

- **Accuracy:** Measures the proportion of correctly classified instances.
- **Precision & Recall:** Evaluate model reliability in predicting positive cases.
- **F1-Score:** Provides a balance between precision and recall.
- **Confusion Matrix:** Visual representation of true positives, false positives, true negatives, and false negatives.



The model achieved high accuracy in distinguishing individuals with and without heart disease, outperforming traditional machine learning models. Feature importance analysis revealed that factors such as cholesterol levels, blood pressure, and smoking habits had the highest impact on predictions.

## 5. Discussion and Conclusions

This study demonstrates the potential of deep learning in predicting heart disease based on multiple health indicators. The proposed neural network model effectively learns patterns from the dataset and provides reliable classifications. The high performance of the model suggests that deep learning can serve as a valuable tool in clinical decision-making, potentially aiding healthcare professionals in early diagnosis and prevention strategies.

## 6. Future Work

While the proposed model achieves high accuracy, further research can be conducted to enhance its effectiveness. Future directions include:

- **Integrating Additional Clinical Data:** Incorporating medical imaging or genomic data may improve model performance.
- **Exploring Alternative Architectures:** Investigating convolutional neural networks (CNNs) or transformer-based models for enhanced feature extraction.
- **Improving Interpretability:** Implementing explainable AI (XAI) techniques to make the model's decisions more transparent to healthcare professionals.
- **Expanding the Dataset:** Utilizing larger and more diverse datasets to increase model generalizability.

## 7. References

- Rajkomar, A., et al. (2018). "Scalable and accurate deep learning with electronic health records." NPJ Digital Medicine, 1(1), 18.
- Esteva, A., et al. (2017). "Dermatologist-level classification of skin cancer with deep neural networks." Nature, 542(7639), 115-118.
- Krittanawong, C., et al. (2019). "Artificial intelligence in precision cardiovascular medicine." Journal of the American College of Cardiology, 74(10), 1307-1321.
- Al'Aref, S. J., et al. (2020). "Machine learning in cardiovascular medicine: are we there yet?" Journal of the American College of Cardiology, 75(5), 573-584.