

# Project Proposal: Heart Disease Prediction

## Introduction

Heart disease is one of the leading causes of death worldwide. Early detection and prediction of heart disease can significantly prolong patients' life. This project aims to develop a machine learning model to predict the presence of heart disease in individuals based on various health metrics.

## Dataset and Data reading

```
[1]: import pandas as pd

df = pd.read_csv('heart.csv', index_col=False)
```

This data set contains 14 attributes, including the predicted attribute. These 13 variables have been screened and are not affected by multicollinearity. The "target" field refers to the presence of heart disease in the patient.

## Questions

- Can we predict the presence of heart disease based on these variables?
- We will use four models: logistic regression, K-NN, decision tree and SVM to study the data. Which machine learning model performs best in predicting heart disease in this dataset? Which hyperparameters perform best for that model?

## Variables

- Target Variable:
  - **target**: A binary variable where 1 indicates the presence of heart disease and 0 indicates its absence.
- Predictor Variables:
 

<ul style="list-style-type: none"> <li>• <b>age</b></li> <li>• <b>sex</b>: 1 = male, 0 = female</li> <li>• <b>cp</b>: Chest pain type (4 types).</li> <li>• <b>trestbps</b>: Resting blood pressure (in mm Hg).</li> <li>• <b>chol</b>: Serum cholesterol in mg/dl.</li> <li>• <b>fbs</b>: Fasting blood sugar &gt; 120 mg/dl (1 = true; 0 = false).</li> <li>• <b>restecg</b>: Resting electrocardiographic results (values 0, 1, 2).</li> <li>• <b>thalach</b>: Maximum heart rate achieved.</li> </ul>	<ul style="list-style-type: none"> <li>• <b>exang</b>: Exercise-induced angina (1 = yes; 0 = no).</li> <li>• <b>oldpeak</b>: ST depression induced by exercise relative to rest.</li> <li>• <b>slope</b>: The slope of the peak exercise ST segment.</li> <li>• <b>ca</b>: Number of major vessels (0-3) colored by fluoroscopy.</li> <li>• <b>thal</b>: A blood disorder indicator (3 = normal; 6 = fixed defect; 7 = reversible defect)</li> </ul>
---	--

## Methods

- Logistic Regression: Use this as a starting model and interpret the effect of each feature.
- Decision Tree: To explore a non-linear model and identify feature importance.
- K-Nearest Neighbors: To assess a simple distance-based approach.
- Support Vector Machine: To apply a robust classification algorithm for high-dimensional data.