

# “Heart Failure Prediction with Machine Learning: A Comparative Study” by Jing Wang:

---

## 1. Introduction

Heart failure (HF) is a critical health issue globally, contributing to high mortality and hospitalization rates. With the increasing availability of electronic health records (EHR) and clinical data, machine learning (ML) provides a powerful toolkit for early diagnosis and prediction of HF. This section introduces the motivation behind the study, emphasizing the need to improve predictive accuracy through the comparative evaluation of different ML models. The study specifically aims to address two major concerns in clinical datasets: **class imbalance** and **feature scaling**, which can significantly affect model performance.

The key contributions include:

- Comparison of 18 ML algorithms on a heart failure dataset.
- Investigation of the effects of two feature scaling techniques: Z-score and Min-Max normalization.
- Evaluation of **SMOTE** (Synthetic Minority Oversampling Technique) for dealing with class imbalance.
- Determination of the most effective model for HF prediction based on multiple metrics like accuracy, F1 score, and AUC.

## 2. Literature Review

This section discusses previous studies on heart failure prediction using ML. It notes that traditional methods often used fewer models and limited evaluation criteria, which restricted their effectiveness. Prior works predominantly focused on:

- Decision Trees (DT), Support Vector Machines (SVM), and Logistic Regression (LR).
- Use of small, balanced datasets that do not represent real-world scenarios.
- Limited exploration of data preprocessing techniques.

The literature highlights gaps such as:

- Inadequate handling of **class imbalance**, where patients with heart failure are underrepresented.
- A lack of comprehensive comparison across a diverse set of machine learning models.

- Limited attention to **scaling techniques**, despite their influence on algorithm performance.

This background justifies the comprehensive comparative approach taken in this study.

### 3. Dataset

The study uses the **Cleveland Heart Disease dataset** from the UCI Machine Learning Repository. It contains **303 records** with **13 attributes** relevant to heart disease diagnosis, such as:

- Age, sex, resting blood pressure, cholesterol level
- Chest pain type, maximum heart rate achieved, exercise-induced angina
- Fasting blood sugar, ST depression, etc.

A key characteristic of the dataset is its **class imbalance** — a common issue in medical datasets — where fewer patients are labeled as having heart disease. This imbalance can cause ML models to be biased towards the majority class. The study applies **SMOTE** to synthetically balance the dataset before training the models.

### 4. Methodology

This section outlines the experimental pipeline in four stages:

#### a. Data Preprocessing

- **Missing Value Handling:** Any records with missing values were removed to maintain data quality.
- **Normalization:** Two techniques were tested:
  - **Z-Score Normalization** (Standardization): Transforms data to have a mean of 0 and standard deviation of 1.
  - **Min-Max Normalization:** Scales data to a range [0,1].
- **Class Imbalance Handling:** SMOTE was used to synthetically generate minority class instances, making the dataset balanced and improving sensitivity.

#### b. Machine Learning Models

The paper evaluates **18 algorithms** from different categories:

- **Linear Models:** Logistic Regression (LR), Ridge Classifier
- **Tree-Based Models:** Decision Tree, Random Forest (RF), Gradient Boosting (GB), Extra Trees

- **Instance-Based:** K-Nearest Neighbors (KNN)
- **Support-Based:** Support Vector Machine (SVM)
- **Ensemble Models:** AdaBoost, Bagging
- **Probabilistic Models:** Naive Bayes
- **Neural Network:** Multi-Layer Perceptron (MLP)
- Others include LDA, Quadratic Discriminant Analysis (QDA), Passive Aggressive Classifier, SGD Classifier

### c. Model Evaluation

Each model was evaluated using:

- **Accuracy:** Overall correctness.
- **Precision:** Correct positive predictions.
- **Recall (Sensitivity):** True positive rate.
- **F1-Score:** Harmonic mean of precision and recall.
- **ROC-AUC:** Measures the ability to distinguish between classes.

5-fold cross-validation was used to ensure reliability and reduce overfitting.

## 5. Results and Discussion

This section presents and analyzes the performance of the models. Key insights include:

- Models like **Random Forest**, **Gradient Boosting**, and **Extra Trees** performed best in terms of F1-score and AUC, especially when used with **Z-score normalization and SMOTE**.
- **Min-Max normalization** slightly underperformed compared to Z-score in most cases.
- Simpler models like Logistic Regression and Naive Bayes showed decent performance but were outperformed by ensemble techniques.
- The use of **SMOTE significantly improved** the recall of models by allowing them to detect more positive cases (heart failure cases).
- The best-performing model achieved:
  - **Accuracy:** Around 93%
  - **F1-score:** Above 0.90
  - **AUC:** Close to 0.95

The findings confirm that:

- Data preprocessing significantly impacts model performance.
- Ensemble models consistently outperform individual models for medical classification tasks.
- Feature scaling and class balancing are crucial steps in the ML pipeline.

## 6. Conclusion

The study concludes that:

- Machine learning models can be effective tools for predicting heart failure when properly tuned and trained.
- The **combination of Z-score normalization and SMOTE** provided the most significant performance improvements across models.
- **Ensemble learning methods** (Random Forest, Gradient Boosting, Extra Trees) are superior in handling complex patterns and producing robust results.
- The research demonstrates the importance of comprehensive model evaluation and data preprocessing in building reliable predictive systems for healthcare.

Future work could involve:

- Using larger and more diverse datasets.
- Incorporating deep learning models.
- Evaluating models in real clinical settings to validate practical applicability.