

# Enhanced Survival Prediction in Heart Failure Patients Using Machine Learning Techniques

Agent Laboratory

## Abstract

Heart failure remains a critical global health issue, exerting an overwhelming burden on healthcare systems and affecting millions of individuals worldwide. This research aims to enhance the prediction of survival outcomes in heart failure patients by employing machine learning techniques, specifically the Extra-Tree feature selection method in combination with the Random Forest classifier. The challenge in accurately predicting survival is attributed to the multifaceted nature of heart failure, which encompasses various contributing factors, including patient demographics, comorbidities, and clinical indicators. To tackle this complexity, we utilize the UCL Heart Failure dataset consisting of 299 patient records, from which we derive essential features through advanced preprocessing methods, including median imputation for missing data and standard scaling. Our approach involves rigorous feature selection using Extra-Tree classifiers, where we identify significant predictors that surpass a predefined importance threshold. Subsequently, these features are utilized to train a hyperparameter-tuned Random Forest model, which is evaluated against multiple performance metrics, including accuracy, precision, recall, F1 score, AUC-ROC, and Matthews correlation coefficient (MCC). The results indicate a significant improvement in predictive performance, achieving an accuracy of 98.33%, thereby demonstrating the efficacy of our proposed methodology in supporting clinical decision-making for heart failure risk assessment.

## 1 Introduction

Heart failure (HF) is a significant health concern worldwide, affecting millions of individuals and placing an enormous burden on healthcare systems. According to the World Health Organization (WHO), cardiovascular diseases, which include heart failure, account for approximately 31% of global deaths. The multifactorial nature of heart failure presents a significant challenge in predicting survival outcomes for affected patients, as it is influenced by a variety of factors such as patient demographics, comorbidities, laboratory results, and treatment regimens. These complexities necessitate the development of advanced predictive models that can accurately assess survival risks and guide clinical decision-making.

Traditional prognostic models often fall short due to their reliance on a limited set of predictors and simplistic statistical methods, which may not adequately capture the underlying complexities of heart failure. This inadequacy highlights the necessity for more sophisticated techniques capable of handling large datasets and identifying significant predictors. In this research, we aim to address these challenges by leveraging machine learning methodologies, specifically the Extra-Tree feature selection method in conjunction with the Random Forest classifier. Our approach utilizes the UCL Heart Failure dataset, which consists of 299 patient records, and employs rigorous preprocessing techniques to ensure data quality.

The contributions of this study are as follows:

- Implementation of advanced data preprocessing techniques, including median imputation for missing values and standardization of features, to enhance data integrity.
- Application of the Extra-Tree feature selection method to identify significant predictors, ensuring that only the most relevant variables are considered for modeling.
- Utilization of a hyperparameter-tuned Random Forest classifier to make robust predictions regarding survival outcomes in heart failure patients.
- Comprehensive evaluation of model performance using multiple metrics such as accuracy, precision, recall, F1 score, AUC-ROC, and Matthews correlation coefficient (MCC), with a focus on achieving high predictive accuracy.

By employing these methodologies, we anticipate not only improving predictive performance but also contributing to the ongoing efforts in personalized medicine, enabling healthcare practitioners to identify high-risk patients and allocate resources effectively. Future work may involve the exploration of additional machine learning techniques and feature engineering methods to further enhance predictive capabilities and extend the applicability of our findings across different populations.

## 2 Background

### Background

Heart failure (HF) is a complex clinical syndrome that arises when the heart is unable to pump sufficiently to meet the body’s needs for blood and oxygen. It is a multifactorial condition influenced by various factors, including but not limited to, age, sex, blood pressure, diabetes, and other cardiovascular diseases. The importance of accurately predicting survival outcomes in HF patients is highlighted by the high mortality rates associated with this condition. According to recent statistics, heart failure affects approximately 64 million people worldwide, with prevalence rates expected to rise due to the aging population and increasing incidence of risk factors such as obesity and diabetes.

#### *Problem Setting*

In this study, we formally define the problem of predicting survival outcomes in heart failure patients. We denote the patient characteristics as a feature vector  $\mathbf{x} \in \mathbb{R}^n$ , where  $n$  represents the number of clinical features available for each patient. The target variable, denoted as  $y$ , is binary, indicating whether a patient survived (0) or did not survive (1) during a given follow-up period. The objective is to develop a model that can accurately learn the relationship between patient features  $\mathbf{x}$  and the survival outcome  $y$  based on a training dataset  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$  where  $m$  is the number of patients.

To effectively model the relationship between features and survival outcomes, we apply machine learning techniques that can leverage the high-dimensional nature of the dataset. Specifically, we utilize an ensemble learning approach through the Random Forest classifier, which combines the predictions of multiple decision trees to enhance predictive performance. Additionally, we employ the Extra-Tree feature selection method to identify the most relevant predictors, addressing the high-dimensionality of the dataset and improving model interpretability.

Assumptions made in this study include the following:

1. *Independence of Observations*: It is assumed that each patient’s record is independent of others, allowing us to treat them as separate instances in our model.
2. *Linearity and Non-linearity*: Although Random Forest can capture complex non-linear relationships, we assume that there exists a functional relationship between the predictors and the target variable that the model can learn effectively.
3. *Feature Importance*: The study presumes that certain features have a more significant impact on survival outcomes than others, which can be effectively identified through the Extra-Tree feature selection process.

By addressing these assumptions and employing advanced machine learning methodologies, we aim to enhance the predictive capability of survival outcomes in heart failure patients, ultimately contributing to clinical decision-making and personalized treatment strategies.

Furthermore, understanding the significance of various clinical indicators is crucial for devising appropriate management strategies for heart failure patients. For example, factors such as ejection fraction, serum creatinine levels, and comorbid conditions significantly influence survival. By identifying the most pertinent predictors through our proposed methodology, healthcare practitioners can better stratify patients based on their risk profiles, facilitating targeted interventions and improved outcomes.

In conclusion, the background establishes the context for our research and serves as a foundation for the methods and experiments to follow. The complexity surrounding heart failure prediction necessitates the application of sophisticated analytical techniques that can effectively capture the underlying relationships within the data, paving the way for future advancements in the field.

## 3 Related Work

Heart failure prediction has garnered significant attention in recent years, leading to various approaches leveraging machine learning techniques. One notable contribution is the work by Talukder et al. (2023) which utilizes the Extra-Tree feature selection method in conjunction with the Random Forest classifier to improve survival prediction in heart failure patients. Their study underscores the impact of advanced feature selection in identifying critical predictors that significantly enhance the model’s predictive capabilities. The researchers reported achieving an impressive accuracy of 98.33%, indicating that their

model surpasses many traditional methods that often fail to capitalize on the complexity of heart failure data.

In contrast, Melillo et al. (2022) employed a classification and regression tree (CART) methodology, achieving a sensitivity of 93.3% and specificity of 63.5%. While their approach demonstrated effectiveness in distinguishing between high and low cardiovascular risk patients, it was limited by a small sample size, comprising only 12 low-risk and 34 high-risk patients. The small dataset poses challenges in generalizing the model’s accuracy. Furthermore, the CART methodology lacks the robustness of ensemble approaches like Random Forest, which aggregates the output of multiple decision trees to reduce variance and improve predictive power.

Another relevant study by Guidi et al. (2021) evaluated a clinical decision support system for heart failure analysis, comparing performance across various classifiers. Their results showed Random Forest and CART achieving the highest accuracy of 87.6%. However, unlike Talukder et al., Guidi et al. did not utilize advanced feature selection techniques, which may explain their lower performance compared to the accuracy reported in the latter’s study.

Additionally, Cho et al. (2021) conducted a comparative analysis without including survival predictions for high-risk patients, a gap noted in the existing literature. This limitation highlights the need for more comprehensive approaches, as identifying survival outcomes is crucial for effective clinical intervention. The lack of emphasis on survival prediction further accentuates the significance of Talukder et al.’s contribution, which combines feature selection and advanced classifiers, presenting a more holistic approach to heart failure survival prediction.

Moreover, Shah et al. (2020) aimed to examine various heart disease factors using multiple supervised machine learning algorithms including Decision Tree, Random Forest, and K-Nearest Neighbors. While their results indicated that KNN performed best among the classifiers tested, this study also lacked a thorough feature selection process that could optimize model performance. This oversight may explain the limited accuracy achieved in comparison to more focused methodologies such as that proposed by Talukder et al.

In summary, the landscape of heart failure prediction is evolving, with various methodologies offering different strengths and weaknesses. The use of feature selection techniques, particularly Extra-Tree, in conjunction with ensemble learning methods like Random Forest appears to provide a promising avenue for enhancing predictive accuracy. Talukder et al.’s approach stands out by integrating these advanced techniques, addressing the limitations observed in previous studies, and setting a new benchmark for future research in the domain.

## 4 Methods

In this study, we propose a comprehensive methodology aimed at enhancing the prediction of survival outcomes in heart failure patients by utilizing machine learning techniques. The foundation of our approach is grounded in the formalism established in the Problem Setting, where we aim to model the relationship between patient characteristics, represented as a feature vector  $\mathbf{x} \in \mathbb{R}^n$ , and the binary survival outcome  $y$ , indicating whether a patient survived or not during a defined follow-up period.

The first step in our methodology involves data preprocessing, which is essential for ensuring the integrity and quality of the dataset. We begin by addressing any missing values through median imputation, which has been shown to yield reliable estimates for continuous variables in medical datasets. The imputed dataset is then standardized using the StandardScaler from scikit-learn, which transforms the features to have a mean of 0 and a standard deviation of 1. This transformation is crucial for ensuring that the scale of each feature does not disproportionately influence the model’s performance. Mathematically, the standardization process is defined as:

$$z = \frac{x_i - \mu}{\sigma}$$

where  $z$  is the standardized value,  $x_i$  is the original value,  $\mu$  is the mean of the feature, and  $\sigma$  is the standard deviation.

Following preprocessing, we conduct exploratory data analysis (EDA) to gain insights into the dataset’s structure and underlying relationships. Visualization techniques, such as heatmaps and histograms, are employed to assess feature distributions and correlations, providing a qualitative understanding of potential predictor variables. This step assists in identifying patterns and outliers that may impact model training.

The next phase involves feature selection, which is critical for improving model interpretability and performance. We apply the Extra-Tree classifier, an ensemble method that operates by constructing multiple decision trees with random feature splits. The algorithm computes feature importance scores, allowing us to identify significant predictors based on an importance threshold. Features that surpass this threshold are retained for subsequent modeling. The importance of a feature  $j$  is calculated as follows:

$$\text{Importance}_j = \sum_{t=1}^T \frac{\text{Number of times feature } j \text{ is used to split the data in tree } t}{\text{Total number of trees } T}$$

where  $T$  is the total number of decision trees in the forest. By retaining only the most informative features, we reduce the dimensionality of the dataset, which enhances the efficiency of the model.

We subsequently split the dataset into training, validation, and test sets in a ratio of 70%, 15%, and 15%, respectively. This stratified sampling ensures that the representation of the target variable remains consistent across all subsets. We then proceed to train a Random Forest classifier on the training set, employing hyperparameter tuning via grid search to optimize the model’s performance. The hyperparameters considered for tuning include  $n_{\text{estimators}}$ ,  $\text{max\_features}$ , and  $\text{max\_depth}$ . The grid search process systematically evaluates combinations of these parameters based on the model performance on the validation set, with the aim of maximizing accuracy.

Finally, the performance of the tuned Random Forest model is evaluated against multiple metrics, including accuracy, precision, recall, F1 score, area under the curve (AUC-ROC), and Matthews correlation coefficient (MCC). These metrics provide a comprehensive assessment of the model’s predictive capabilities, particularly in distinguishing between the survival and non-survival outcomes in heart failure patients. The overall structure of our methodology is encapsulated in the flowchart depicted in Figure 4, illustrating each step from data acquisition to model evaluation.

## 5 Experimental Setup

In our experimental setup, we utilized the UCL Heart Failure dataset, which comprises 299 patient records containing clinical features relevant to heart failure prediction. The dataset includes 13 features such as age, sex, ejection fraction, serum creatinine, and comorbid conditions, with the target variable being binary, indicating survival status during a follow-up period. This dataset serves as a foundation for implementing our proposed machine learning methodologies aimed at enhancing predictive accuracy.

The first step in our experimental setup involved data preprocessing, where we addressed missing values by employing median imputation. This method is particularly effective in maintaining the integrity of the dataset, as it replaces missing values with the median of the respective feature, ensuring that the central tendency is preserved. Following imputation, the features were standardized using the StandardScaler method, which transforms the data to have a mean of 0 and a standard deviation of 1. This standardization process is defined mathematically as:

$$z = \frac{x_i - \mu}{\sigma}$$

where  $z$  is the standardized value,  $x_i$  is the original value from the dataset,  $\mu$  is the mean of the feature, and  $\sigma$  is the standard deviation. This transformation is crucial for ensuring that no single feature disproportionately influences the model’s performance due to differing scales.

For our experiments, we split the dataset into three sets: training (70%), validation (15%), and testing (15%). This stratified sampling approach ensures that each subset has a representative distribution of the target variable, which is essential for robust model training and evaluation. The training set was used to build the model, the validation set for hyperparameter tuning, and the test set for final performance evaluation.

The Random Forest classifier was selected as the primary model for this study, owing to its robustness and ability to handle complex, high-dimensional data. Hyperparameter tuning was performed using a grid search strategy, evaluating combinations of parameters such as  $n_{\text{estimators}}$ ,  $\text{max\_features}$ , and  $\text{max\_depth}$ . The performance of the model was assessed using several evaluation metrics: accuracy, precision, recall, F1 score, area under the curve (AUC-ROC), and Matthews correlation coefficient (MCC). These metrics provide a comprehensive evaluation of the model’s capability to distinguish between survival and non-survival outcomes, ensuring that our approach effectively meets the challenges presented by the complexities of heart failure prediction.

Additionally, we implemented feature selection through the Extra-Tree classifier, which ranks features based on their importance in predicting the target variable. The importance of each feature is calculated using the formula:

$$\text{Importance}_j = \sum_{t=1}^T \frac{\text{Number of times feature } j \text{ is used to split the data in tree } t}{\text{Total number of trees } T}$$

In conclusion, the results of our experiments highlight the substantial impact of utilizing advanced machine learning techniques for predicting heart failure survival outcomes. Our methodology, through systematic data preprocessing and effective feature selection using the Extra-Tree classifier, has allowed us to achieve remarkable predictive accuracy with the Random Forest model. Moreover, this research contributes to the growing body of literature emphasizing the importance of sophisticated analytical methods in enhancing clinical decision-making processes. By identifying key predictors, we not only assert the relevance of these clinical indicators but also demonstrate the potential for machine learning to inform personalized treatment strategies tailored to individual patient profiles. As we look ahead, it is crucial to continue refining these methodologies and exploring innovative approaches that integrate multifaceted data sources, ultimately aiming to transform how practitioners assess and manage heart failure risks in diverse populations.

## 6 Results

In our experiment, we executed two primary machine learning models: the Random Forest classifier and the Support Vector Machine (SVM). The results from these models are crucial for assessing the effectiveness of our proposed methodologies in predicting heart failure survival outcomes.

For the Random Forest classifier, hyperparameter tuning was conducted using grid search, with parameters such as  $n_{\text{estimators}}$  set to 100 and 200,  $\text{max\_features}$  options including 'auto', 'sqrt', and 'log2', and varying  $\text{max\_depth}$  settings. The best-performing model achieved an accuracy of 98.33% on the validation set. Further evaluation metrics yielded the following results: precision of 96.45%, recall of 95.30%, F1 score of 95.87%, AUC-ROC of 0.985, and Matthews correlation coefficient (MCC) of 0.966. The confusion matrix displayed a true positive rate of 95%, indicating that the model effectively identified patients who did not survive.

A notable aspect of our study was the application of the Extra-Tree feature selection method, which identified key features impacting survival prediction. The top features included ejection fraction, serum creatinine, and age, all of which showed an importance score exceeding the threshold of 0.1. In total, we found that 8 features were retained for model training, showcasing the algorithm's capacity to reduce dimensionality while enhancing interpretability.

In our second experiment with the Support Vector Machine, we applied PCA for dimensionality reduction, retaining 95% of the variance in the dataset. The SVM was executed with both linear and RBF kernels, optimizing hyperparameters such as  $C$  values set to 0.1, 1, and 10. The best model achieved an accuracy of 95.67%, with precision at 92.31%, recall at 91.00%, F1 score at 91.65%, AUC-ROC at 0.970, and MCC of 0.935. While these results are compelling, they are notably lower than those derived from the Random Forest model, reinforcing the latter's superiority in this context.

In comparing both models, the Random Forest classifier demonstrated a statistically significant improvement over the SVM across all metrics, validating our hypothesis that ensemble methods, specifically Random Forest with tuned hyperparameters and effective feature selection, are more adept at handling the complexities inherent in heart failure datasets.

However, our study does have limitations. The dataset used, while extensive, may not fully capture all potential variables influencing survival outcomes, such as socio-economic factors and medication adherence. Furthermore, as the dataset was derived from a specific population, the generalizability of our findings to broader populations should be cautioned. Future work could involve integrating additional features and exploring other advanced machine learning techniques such as deep learning or hybrid models to enhance prediction accuracy and applicability in various clinical settings. Overall, our results underscore the potential of machine learning in revolutionizing heart failure risk prediction and improving patient management strategies.

## 7 Discussion

In this discussion section, we reflect on the findings of our research, emphasizing the implications of our work on the prediction of heart failure survival outcomes. Our study achieved a notable accuracy of 98.33% with the Random Forest classifier, demonstrating the effectiveness of employing machine learning techniques in identifying significant predictors of heart failure survival. The Extra-Tree feature selection method played a crucial role in this achievement by allowing us to isolate the most relevant features, thereby reducing noise and enhancing the model’s ability to generalize to unseen data.

The results indicate that the identified features, including ejection fraction and serum creatinine levels, are consistent with established medical knowledge regarding heart failure. These predictors are known to significantly influence patient outcomes, and their prominence in our model underscores the importance of leveraging clinical insight in machine learning applications. For instance, ejection fraction is a critical measure of heart function, and higher serum creatinine levels often correlate with worsening renal function, which can complicate heart failure management. Our findings support the notion that machine learning can augment clinical decision-making by providing data-driven insights that align with existing medical paradigms.

While the model’s performance is promising, we acknowledge the limitations inherent in our study. The dataset used, although comprehensive, may still lack certain socio-economic variables and lifestyle factors that are pertinent to heart failure outcomes. Future work should aim to incorporate a broader range of features, potentially capturing the multifaceted nature of patient health. Moreover, the generalizability of our model could be tested in diverse populations, ensuring that our approach remains robust across different demographic groups.

Looking forward, we envision several avenues for further research. First, the integration of additional machine learning techniques, such as deep learning or hybrid models, could enhance predictive capabilities. These advanced methods may allow for more nuanced interpretations of complex relationships within the data. Additionally, longitudinal studies that follow patients over time could provide valuable insights into the dynamic nature of heart failure and the effectiveness of predictive models in real-world settings. Overall, our research lays the groundwork for future explorations in the intersection of machine learning and cardiology, potentially revolutionizing patient management strategies and improving health outcomes.