

# Abstract

This dissertation is a comprehensive exploration of the creation and enhancement of predictive models for heart failure prediction, leveraging the heart failure dataset. Heart failure remains a significant global health concern, and the development of accurate predictive models holds immense importance for early intervention and treatment planning. This research investigates four prominent predictive modeling algorithms: Support Vector Machine (SVM), Decision Trees, Random Forest, and XGBoost. Each model's performance is meticulously assessed using a variety of metrics, including Confusion Matrices, ROC Curves, and Precision-Recall Curves, shedding light on their predictive capabilities in the context of heart failure.

In pursuit of model optimization, a detailed examination of hyperparameters is undertaken through Grid Search, which remarkably improves the accuracy of each model. The study extracts key insights from the model evaluations, emphasizing their implications for the broader research objectives. While the results are promising, the dissertation acknowledges potential limitations, such as overfitting, dataset size constraints, and the inherent challenges associated with evaluation metrics.

The dissertation concludes with a call for future research to explore alternative algorithms, emerging technologies, and an expanded dataset scope to further enhance predictive accuracy. It underscores the profound impact of optimization techniques on predictive model performance and their potential to drive innovations in heart failure prediction. This research thus stands as a declaration to the power of machine learning in predictive analytics, providing a roadmap for future endeavors in this critical field of study.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Project Objective . . . . .	1
1.2	Background and Rationale for the Study . . . . .	1
1.3	Methodological Overview . . . . .	2
1.4	Expected Outcomes . . . . .	2
1.5	Structure of the Dissertation . . . . .	3
<b>2</b>	<b>Exploratory Data Analysis (EDA)</b>	<b>5</b>
2.1	Data Sources and Collection . . . . .	5
2.2	Dataset Description . . . . .	5
2.3	Data Preprocessing . . . . .	7
2.4	Basic Statistics and Visualization . . . . .	9
2.4.1	Quantitative Variables: . . . . .	9
2.4.2	Qualitative Variable Insights: . . . . .	10
2.4.3	Inter-Variable Relationships: . . . . .	12
2.5	Key Insights from Exploratory Data Analysis: . . . . .	14
2.6	Data Preparation and Splitting: . . . . .	19
2.7	Standardization: . . . . .	19
<b>3</b>	<b>Predictive Modeling for Heart Failure</b>	<b>21</b>
3.1	Support Vector Machines (SVM): . . . . .	22
3.1.1	Justification for using SVM: . . . . .	22
3.1.2	Mathematical Overview: . . . . .	23
3.2	Decision Trees: . . . . .	24
3.2.1	Justification for using Decision Trees: . . . . .	24
3.2.2	Mathematical Overview: . . . . .	25
3.3	Random Forest: . . . . .	26
3.3.1	Justification for Using Random Forest: . . . . .	26
3.3.2	Mathematical Overview: . . . . .	27
3.4	XGBoost classifier: . . . . .	28
3.4.1	Justification for using XGBoost: . . . . .	28
3.4.2	Mathematical Overview: . . . . .	29
<b>4</b>	<b>Performance Evaluation and Model Optimization</b>	<b>31</b>
4.1	Confusion Matrix: . . . . .	31
4.1.1	Components of a Confusion Matrix: . . . . .	31
4.1.2	Confusion Matrix for SVM: . . . . .	33
4.1.3	Confusion Matrix for Decision Trees: . . . . .	34

4.1.4	Confusion Matrix for Random Forest: . . . . .	34
4.1.5	Confusion Matrix for Xgboost Classifier: . . . . .	35
4.2	ROC Curve Analysis: . . . . .	36
4.2.1	ROC Curve Analysis for SVM and XGBoost: . . . . .	37
4.2.2	ROC Curve Analysis for Decision Trees and Random Forest: . . . . .	37
4.3	Precision-Recall Curve Analysis: . . . . .	39
4.4	Model Selection: Deciding the Best Fit . . . . .	40
4.5	Model Optimization: . . . . .	41
4.5.1	Hyperparameter Tuning Technique: . . . . .	42
4.5.2	SVM: . . . . .	42
4.5.3	Decision Trees: . . . . .	43
4.5.4	Random Forest: . . . . .	43
4.5.5	XGBoost classifier: . . . . .	44
<b>5</b>	<b>Conclusion</b>	<b>47</b>
5.1	Recapitulation . . . . .	47
5.2	Major Takeaways . . . . .	47
5.3	Limitations and Challenges: . . . . .	48
5.4	Future Work . . . . .	49
5.5	Acknowledgments: . . . . .	49

# List of Figures

2.1	Boxplot depicting the distribution of data with outliers . . . . .	8
2.2	Histograms showcasing the distributions of 'Age' (left) with a Gaussian or Normal Distribution and 'Creatinine Phosphokinase' (right) following an exponential distribution in the dataset . . . . .	10
2.3	Histograms depicting the Normal (Gaussian) Distribution patterns for 'Ejection Fraction' (left) and 'Platelets' (right) in the dataset. . . . .	11
2.4	Comparative histograms illustrating the skewed distributions of 'Serum Creatinine' (left) with a positive skew and 'Serum Sodium' (right) exhibiting a negative skew in the dataset . . . . .	11
2.5	Comparative Distribution of Categorical Variables against Survival Outcomes .	12
2.6	Correlation Matrix Heatmap . . . . .	13
2.7	Box plot illustrating the relationship between Age and Mortality Risk . . . . .	14
2.8	Box plot illustrating the relationship between Ejection Fraction and Mortality Risk	15
2.9	Box plot illustrating the relationship between Serum-Creatinine and Mortality Risk . . . . .	16
2.10	Box plot illustrating the relationship between Serum-Sodium and Mortality Risk	16
2.11	Box plot illustrating the relationship between Time and Mortality Risk . . . . .	17
2.12	Distribution of Gender (Male and Female) . . . . .	17
2.13	Distribution of High Blood Pressure . . . . .	18
2.14	Distribution of Anaemia . . . . .	18
2.15	Distribution of Smoking . . . . .	19
4.1	Confusion Matrix for SVM . . . . .	33
4.2	Confusion Matrix for Decision Trees . . . . .	34
4.3	Confusion Matrix for Random Forest . . . . .	35
4.4	Confusion Matrix for Xgboost Classifier . . . . .	36
4.5	Receiver operating Characteristic curves for SVM . . . . .	37
4.6	Receiver operating Characteristic curves for XGBoost . . . . .	38
4.7	Receiver operating Characteristic curves for Decison Trees and Random Forest	38
4.8	Model Performance Evaluation via Precision-Recall Curves . . . . .	40
4.9	Comparative Accuracy of Optimized Machine Learning Models . . . . .	45

# List of Tables

2.1	Description of Features in the Dataset . . . . .	6
2.2	Features, Measurements, and Ranges . . . . .	7
2.3	Summary of Features . . . . .	9
3.1	Performance Metrics of the Support Vector Machine (SVM) Classifier . . . . .	23
3.2	Performance Metrics of the Decision Trees . . . . .	25
3.3	Performance Metrics of the Random Forest . . . . .	27
3.4	Performance Metrics of the XGBoost Classifier . . . . .	29
4.1	Confusion Matrix . . . . .	32

# Chapter 1

## Introduction

### 1.1 Project Objective

Heart health continues to be a major focus on modern medical research in a world where the fine boundaries of biology and technology combine. As one examines the complex workings of the human heart, it becomes clear that understanding its potential weaknesses is not only a learning process but also a critical societal requirement. This dissertation's main goal is to clarify the complex issues related to heart failure mortality. We want to find patterns, correlations, and potential predictive indications using machine learning and predictive modeling to help both patients and healthcare professionals manage their heart health. One would be curious about the motivations behind this search as we set out on this journey. The background and reasoning section of our study explains the context of our research, as well as its urgent importance and the gaps it is aimed at filling.

### 1.2 Background and Rationale for the Study

Heart issues touch many of our lives. From our grandparents' stories of their peers' heart ailments to our friends mentioning their parents' conditions, heart problems are sadly familiar to many of us. Heart diseases, particularly cardiovascular diseases (CVDs), remain a primary global health concern, causing significant illness and being the leading cause of death, claiming an estimated 17.9 million lives annually (WHO, 2019).

Now, let's take a step back. When we think about heart failures and their devastating impact, it's tempting to view the situation from a purely clinical perspective. But in truth, every piece of data, every percentage, represents a human story. A story of hope, fear, battles won, and sometimes, battles lost. This isn't just a report. It's an effort to bridge the gap between numbers on a screen and the reality of living with heart disease.

So, why this study? Because there's a pressing need to use all tools at our disposal, especially data analytics, to tackle this age-old issue. By taking a deeper dive into the intricacies of heart failures, their causes, and patterns, we might be able to offer new insights, perhaps even

innovative solutions. If data can give us even a small chance to change the ending of someone's story, it's a pursuit worth undertaking.

But, how exactly do we get from data to these life-saving insights? This leads us to the methods employed.

### 1.3 Methodological Overview

A solid scientific approach is crucial in the subject of cardiac disorders (NHS, 2022) because of the many complex and interrelated elements that characterize it. Data, which is often referred to as the "oil" of the modern day, is crucial for understanding, forecasting, and reducing health concerns (The Economist, 2017). The importance of data analysis cannot be overstated, as it allows us to make sense of the complexity that often clouds our understanding. This is achieved by converting raw data into actionable insights, enabling us to effectively apply this newfound knowledge.

This project stands evidence to the power of data. Exploratory Data Analysis (EDA), a key first step that exposes the complexity of our dataset, serves as the starting point of our investigation. EDA not only highlights important trends and spots probable abnormalities, but it also gives researchers a starting point for further research. The groundwork is therefore set for predictive modeling by building on these fundamental discoveries. SVM (Support Vector Machines), Decision Trees, Random Forest, and XGBoost classifier are a few of the machine learning methods we've selected for this application. These models act as our research tools, helping us to peel back the layers and expose the frequently subtle connections between various health variables and their impact on heart health.

However, applying algorithms alone is ineffective. In refining them, the true essence is revealed. We verify that our models' predictions are accurate by thoroughly testing and training them on the dataset. The complex dance of data gathering, analysis, and modeling is what enables us to discover the hidden meanings in the data. Our strategy combines cutting-edge machine learning algorithms with traditional statistical tools. You will see as we go along the breadth and depth of our methodology, which is all focused toward explaining the hidden causes of heart problems. This methodological analysis sets the foundation for our upcoming section, which discusses the expected outcomes.

### 1.4 Expected Outcomes

In the heart of this dissertation's exploration lies the promise of significant transformations in our battle against heart failure. The anticipated outcomes are not mere aspirations; they are the compass guiding our journey.

At the forefront, the spotlight turns to the identification of **key predictors**. Through rigorous model evaluation, this research aims to uncover the variables that utilize the most influence on



heart failure outcomes. These discoveries could provide healthcare professionals with a potent tool to navigate the complex landscape of patient attributes.

But this journey encompasses more than just improvements in prediction. We're venturing into the domains of **generalization and scalability**, where the models we develop should prove adaptable and consistent across diverse datasets and populations. Scalability ensures these models can handle not just today's data but the ever-growing volumes of tomorrow. As we continue, we imagine a substantial enhancement in **predictive accuracy**. The optimized models are poised to excel in forecasting heart failure events, potentially revolutionizing early intervention and patient care. Imagine a world where healthcare providers can identify at-risk individuals with even greater certainty.

As the research unfolds, it promises to offer a wealth of **informative insights**. These insights will illuminate the unique characteristics and capabilities of each predictive model, showcasing their strengths and limitations. These revelations will serve as guiding lights for future explorers in the realm of model selection and optimization. And at its core, this research represents a significant **contribution to heart failure research**. It aspires to join the chorus of voices dedicated to understanding and addressing heart failure. By providing data-driven tools and insights, it seeks to empower healthcare professionals and researchers alike.

In this narrative, the expected outcomes are not merely a remote vision; they are the fuel propelling us toward advancements in heart failure prediction and, by extension, patient care. Through this research, we aim to engrave a meaningful chapter in the ongoing story of human health, one heartbeat at a time.

## 1.5 Structure of the Dissertation

This dissertation is structured into several chapters, each serving a specific purpose in the exploration of predictive modeling for heart failure. Below is a brief overview of the content covered in each chapter:

- **Introduction:** In this introductory chapter, we set the stage for the entire dissertation. We outline the project's objectives, discuss the background and rationale for the study, provide an overview of our chosen methodology, and highlight the expected outcomes. Additionally, we present the structure of the dissertation, giving readers a glimpse of the organization of our work.
- **Exploratory Data Analysis (EDA):** In this chapter, we dive into the practical aspects of our research. We discuss data sources and collection methods, describe the dataset used in our analysis, detail our data preprocessing steps, and present basic statistics and visualizations. The key insights derived from our exploratory data analysis are also highlighted.
- **Predictive Modeling for Heart Failure:** This chapter is dedicated to the core of our research. We explore various predictive modeling techniques, including Support Vector

Machines (SVM), Decision Trees, Random Forest, and XGBoost classifier. For each technique, we provide justifications, mathematical overviews, and their relevance to our study.

- **Performance Evaluation and Model Optimization:** In this chapter we focus on the evaluation of our predictive models. We discuss concepts such as the confusion matrix, ROC curve analysis, and precision-recall curve analysis. We also delve into the critical aspects of model selection and optimization, including hyperparameter tuning techniques.
- **Conclusion:** In the final chapter, we summarize our findings and key takeaways. We revisit the objectives of our research and discuss the major insights gained. Additionally, we acknowledge the limitations and challenges encountered during our study and propose directions for future research.

## Chapter 2

# Exploratory Data Analysis (EDA)

### 2.1 Data Sources and Collection

When studying heart failure, researchers gather data from a variety of places. Many times, they look at records from hospitals or clinical trials because these offer a lot of detail about patients, like their medical history and treatment outcomes. But it's essential to handle this data carefully to protect patient privacy. There are also big health databases in certain countries or areas that give a wider view by pooling information from many hospitals or clinics. As technology advances, we're getting more data from things like smartwatches that track heart rates and activity levels in real-time (Piwek et.al, 2016). This is exciting, but it's vital to make sure this data is accurate and consistent. Surveys where patients share their experiences and universities conducting their own studies provide more information.

We started on an analysis utilizing a dataset comprising medical histories of 299 heart failure patients. These records were gathered from both the Faisalabad Institute of Cardiology and Allied Hospital in Faisalabad, Punjab, Pakistan, spanning the months from April to December 2015. Notably, all these patients were diagnosed with left ventricular systolic dysfunction. They had previously experienced heart failures, categorizing them under the III or IV classes as per the New York Heart Association (NYHA) heart failure stage classification (Ahmad T et.al, 2017). With this, the reader is guided and prepared for an in-depth exploration of the dataset in the upcoming section.

### 2.2 Dataset Description

To fully explore the complicated nature of heart failure, a large and complete dataset is necessary. Our study is based on a dataset that is exceptional both in terms of the quantity of data and the quality of the information it provides. Each entry in this collection, which has a total of 299 records, contains more than just data. It is a narrative, the account of a person who was given a heart failure diagnosis, contained in rows and columns.

Our dataset’s diversity is evident not just in its quantity but in its depth. These records represent individuals ranging from ages 40 to 95, offering an extensive view of how heart failure interacts across different age groups. The gender distribution, featuring 105 females and 194 males, further enriches our understanding by providing a balanced exploration of the impact of this ailment across genders.

However, for a more structured understanding of the data, we’ve organized the information into two tables. The Table 2.1 provides a broad overview of the dataset, presenting a description of each feature. This table ensures that readers unfamiliar with medical jargon or specific terminologies can still grasp the significance of each variable.

Feature	Description
Age	Represents the patient’s age, measured in years, serving as an essential demographic characteristic.
Anaemia	A binary indicator reflecting the presence (1) or absence (0) of anemia, a condition affecting the blood’s ability to carry oxygen.
CPK	This feature quantifies the Creatinine Phosphokinase (CPK) enzyme level in the blood, with elevated levels possibly signifying heart damage.
Diabetes	This binary variable signifies whether a patient has (1) or doesn’t have (0) diabetes.
Ejection Fraction	Reflects the heart’s efficiency by measuring the percentage of blood pumped out with each contraction.
Hypertension	A binary representation (1 for presence, 0 for absence) of High Blood Pressure (hypertension), a significant cardiovascular risk factor.
Platelets	Denotes the blood’s platelet count, measured, vital for blood clotting processes.
Serum Creatinine	Measured in mg/dL, it denotes the creatinine level in the blood, a key metric for kidney function.
Serum Sodium	The concentration of sodium in the blood, crucial for maintaining the body’s fluid balance.
Sex	A binary variable capturing the patient’s gender (1 for male, 0 for female).
Smoking	A binary variable indicating smoking status (1 for smoker, 0 for non-smoker).
Time	Represents the follow-up period in days from diagnosis to the next follow-up, providing insights into the disease progression.
Death Event	The outcome variable, captured as a binary feature (1 for died, 0 for survived), indicates whether a death event occurred during the follow-up period.

*Table 2.1: Description of Features in the Dataset*

Feature	Measurement	Range
Age	Years	[40, ..., 95]
Anaemia	Boolean	0, 1
High blood pressure	Boolean	0, 1
Creatinine phosphokinase (CPK)	mcg/L	[23, ..., 7861]
Diabetes	Boolean	0, 1
Ejection fraction	Percentage	[14, ..., 80]
Sex	Binary	0, 1
Platelets	kiloplatelets/mL	[25.01, ..., 850.00]
Serum creatinine	mg/dL	[0.50, ..., 9.40]
Serum sodium	mEq/L	[114, ..., 148]
Smoking	Boolean	0, 1
Time	Days	[4,...,285]
Death event (Target)	Boolean	0, 1

*Table 2.2: Features, Measurements, and Ranges*

The Table 2.2 focuses on the measurements and data ranges. It acts as a compass, guiding us through the landscape of numbers, values, and categories, ensuring that our analysis remains grounded in experimental evidence. In essence, while the dataset brings the experiences of 299 heart failure patients to our fingertips, the accompanying tables function as our guidebooks, ensuring clarity, structure, and depth in our exploration. But as with all real-world data, there's an imperative need for refinement and adjustments. Before any meaningful analysis can occur, the raw data must be transformed into a polished, analysis-ready format. This brings us to the crucial phase of data preprocessing, which we'll delve into next.

## 2.3 Data Preprocessing

In our journey through data analytics, imagine a seasoned artist, carefully chipping away at a block of stone. Each strike is precise, slowly revealing the art hidden within. Similarly, data preprocessing is our tool, ensuring our dataset's true story shines through, free of any distortions. Let's uncover the steps taken to make sure our dataset is in its finest form:

1. **Handling Missing Data:** Think of a jigsaw puzzle. Even if one piece is missing, the picture remains incomplete. The same goes for datasets. If any piece of data is missing, it can misguide our analysis. But here's the good news: our dataset is like a perfectly assembled puzzle, without a single piece out of place. No missing values, which means a smooth and uninterrupted analysis journey.
2. **Outlier Detection:** Now, imagine you're listening to a harmonious tune, and suddenly, a few off-key notes play. They catch your attention, don't they? In the world of data, these are outliers. While analyzing, we kept an eagle's eye out for such 'off-key notes'. We

leaned on visual techniques, especially the revealing box plot in Figure 2.1. It highlighted a few features like Creatinine Phosphokinase, Platelets, and Serum Creatinine that had values dancing away from the main group. But here's the twist: given the severity of the conditions of our patients, these weren't missteps in the data. Instead, they were true reflections of their health state. So, instead of dismissing them, we embraced these outliers as part of the real story.

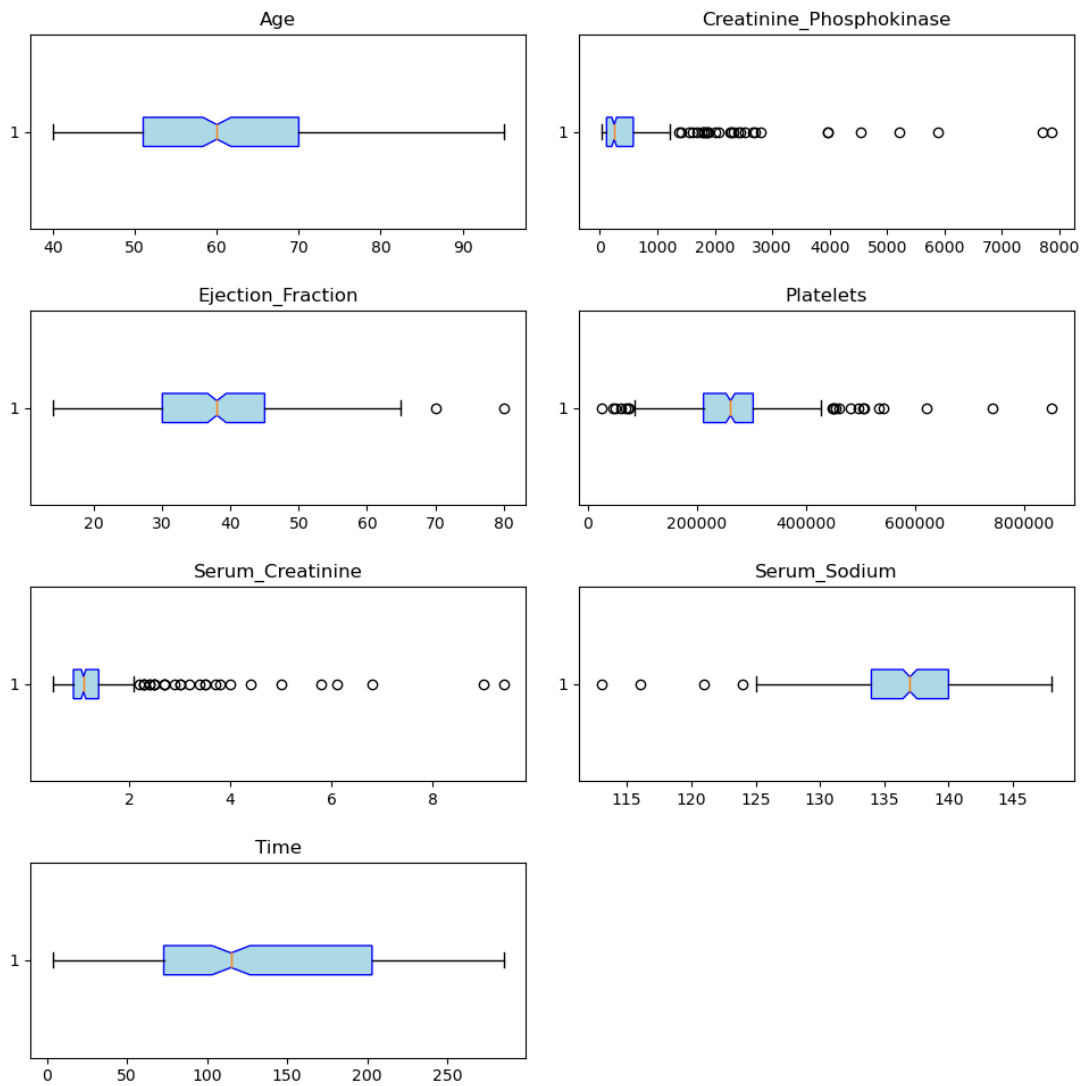


Figure 2.1: Boxplot depicting the distribution of data with outliers

## 2.4 Basic Statistics and Visualization

### 2.4.1 Quantitative Variables:

Journeying through our dataset, one is met with an array of numbers, each echoing the unique health narratives of patients. The age of patients, a crucial variable, spans a significant stretch from a young 40 to a seasoned 95, with the average floating around 61 years. This range underlines the age-related vulnerabilities often associated with cardiovascular health.

Creatinine\_Phosphokinase, instrumental in decoding heart health, presents a broad spectrum – from a mere 23, climbing to a staggering 7861, with its median stationed at 250. This variable offers insights into the varied extents of heart damages prevalent in our group.

Ejection\_Fraction gives another perspective. With its median pinned at 38%, it hints at the cardiac challenges these patients struggle with, emphasizing the significance of this variable in cardiac diagnostics.

Variables like platelet counts, Serum\_Creatinine, Serum\_Sodium, and the span denoted by Time, together sketch the broader health landscape. The latter, denoting the follow-up period, varies from as brief as 4 days to as long as 285 days.

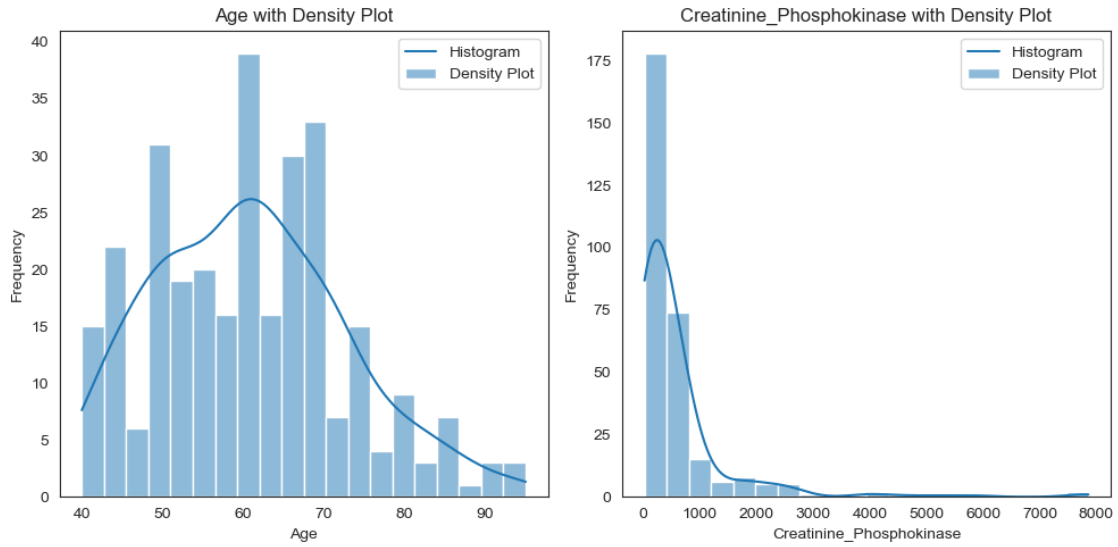
For those seeking a more compact glimpse, a comprehensive statistical summary is available in Table 2.3.

Feature	Count	Mean	Std	Min	25%	50%	75%	Max
Age	299	60.8339	11.8948	40	51	60	70	95
CPK	299	581.8395	970.2879	23	116.5	250	582	7861
Ejection_Fraction	299	38.0836	11.8348	14	30	38	45	80
Platelets	299	263358.0293	97804.2369	25100	212500	262000	303500	850000
Serum_Creatinine	299	1.3939	1.0345	0.5	0.9	1.1	1.4	9.4
Serum_Sodium	299	136.6254	4.4125	113	134	137	140	148
Time	299	130.2609	77.6142	4	73	115	203	285

Table 2.3: Summary of Features

Every variable in our dataset has its own tale to tell, reflecting unique patterns and behaviors. The 'Age' variable, for instance, moves like the steady beat of a heart. It's balanced and consistent, showcasing a familiar pattern we often see in statistics: the Gaussian or Normal Distribution. Most of the subjects' ages gather around a central age, highlighting its symmetrical nature. On the other hand, when we look at the 'Creatinine Phosphokinase' data, the story changes. Picture standing atop a hill, watching water rush down quickly at first and then slowing as it reaches the bottom. This imagery captures the Exponential Distribution of this data. A lot of patients show lower levels, but there's a rapid drop off, pointing to only a few with higher levels. Figure 2.2, the histogram helps visualize this, with a sharp drop at the beginning that gradually tapers off, much like the water slowing down after its initial rush.

The histogram of Ejection Fraction, as represented in Figure 2.3 on the left, adheres to the



*Figure 2.2:* Histograms showcasing the distributions of 'Age' (left) with a Gaussian or Normal Distribution and 'Creatinine Phosphokinase' (right) following an exponential distribution in the dataset

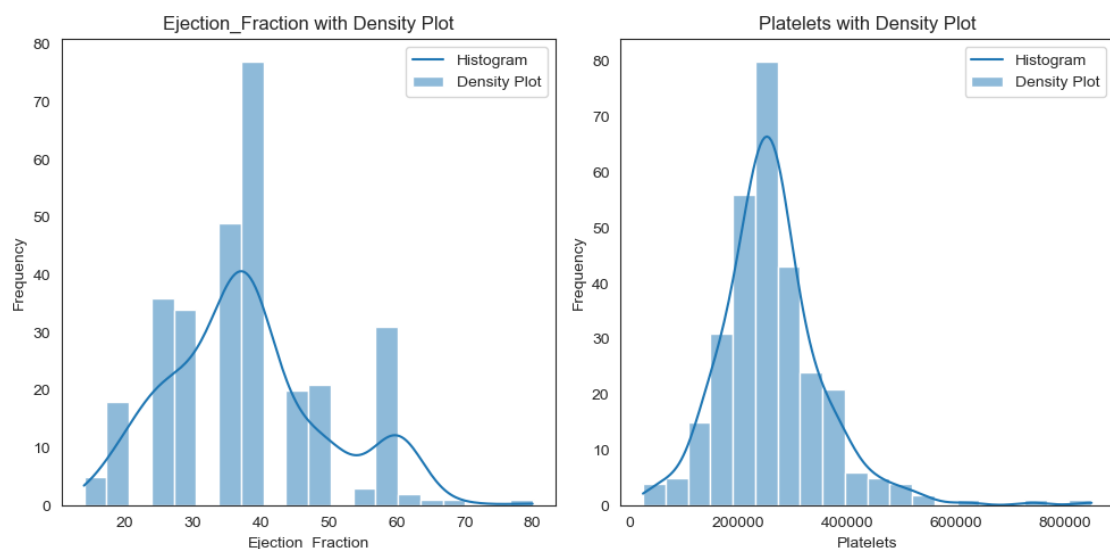
Normal Distribution, indicating a balanced representation of this feature in our dataset. This suggests that most patients' ejection fraction values are centered around the mean, with a symmetrical spread on both sides. Similarly, the distribution of 'Platelets,' as shown on the right of Figure 2.3, also conforms to the Gaussian or Normal Distribution. A majority of patients have platelet counts that cluster around the average, ensuring an even distribution on either side of the mean. This showcases a typical representation of platelet counts in the studied group.

The variable 'Serum Creatinine' shows a Positive Skew in its distribution, left graph in Figure 2.4. In this pattern, most of the data points are concentrated on the lower side, with the tail extending towards the higher values. Essentially, a majority of patients have lower serum creatinine levels, while fewer have higher levels. Conversely, 'Serum Sodium' exhibits a Negative Skew, right graph in Figure 2.4. Here, most of the data points are gathered towards the higher values, with the tail pointing towards the lower side. This means most patients have higher serum sodium levels, but there's a decline with fewer individuals having lower levels. Yet, this numerical tale is merely one facet of our dataset. Let's pivot to the next section, Qualitative Variable Insights, to explore the categorical dimensions, unraveling patterns of habits and lifestyles that further enrich our comprehension.

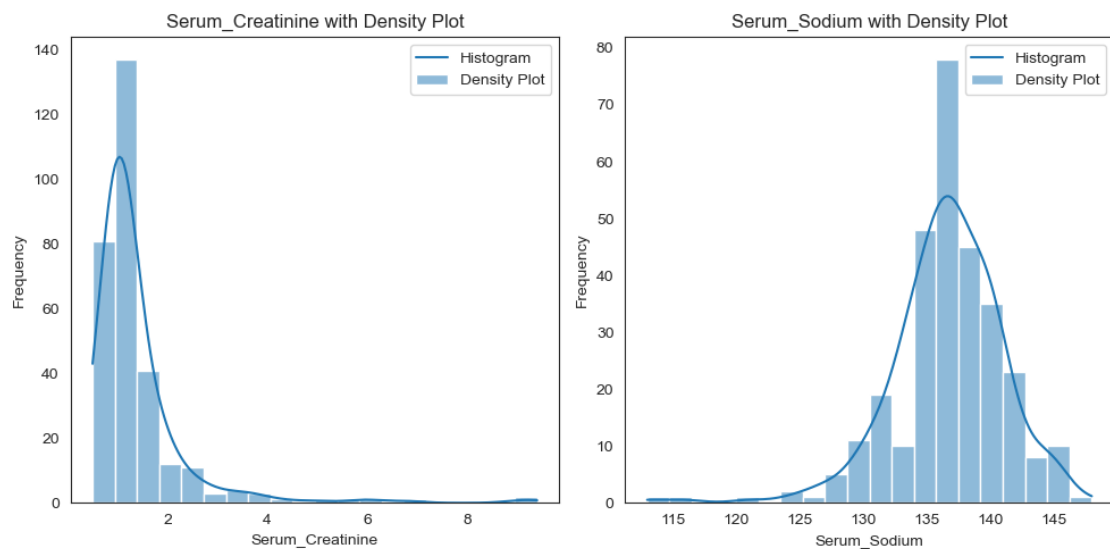
## 2.4.2 Qualitative Variable Insights:

Moving from the quantifiable aspects of our data, we transition into the realm of categorical features, the qualitative variables. These variables, unlike their numerical counterparts, are often described by characteristics or categories rather than numbers. They provide a snapshot of the various groups within our dataset, painting a vibrant picture of diversity and pattern.





*Figure 2.3:* Histograms depicting the Normal (Gaussian) Distribution patterns for 'Ejection Fraction' (left) and 'Platelets' (right) in the dataset.



*Figure 2.4:* Comparative histograms illustrating the skewed distributions of 'Serum Creatinine' (left) with a positive skew and 'Serum Sodium' (right) exhibiting a negative skew in the dataset

Take, for instance, variables like 'Smoking', 'Sex', 'Anaemia', 'Diabetes', and 'High Blood Pressure'. They don't resonate in terms of a continuous scale but rather in clearly defined categories: 'smoker' or 'non-smoker'; 'male' or 'female'; 'yes' or 'no' for the presence of anaemia, diabetes, or high blood pressure. Understanding the distribution of these categories not only gives us insights into the diversity of our dataset but also sketches the broader health profile of the individuals it represents. To further illuminate these insights, We incorporated a bar graph, refer Figure 2.5, contrast each category against the outcomes of 'survived' and 'died'. This visual representation provides a quick comparative glance, showcasing the interplay between these categorical attributes and the survival outcomes. As we delve deeper, we'll scrutinize each of these variables, offering a detailed perspective on their distribution and the patterns they reveal in relation to survival outcomes.

Although knowing each element individually gives us a foundation, it is the relationships between them that reveal deeper stories and patterns. In order to gain additional insights from the interactions they represent, we are prepared to explore the complex network of inter-variable relationships in the section that follows.

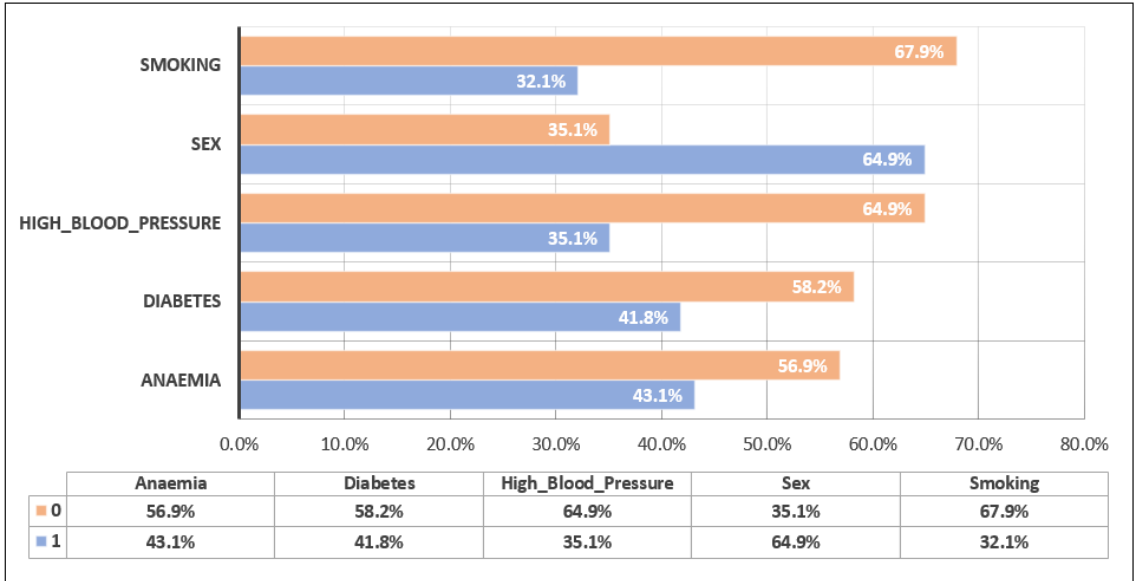


Figure 2.5: Comparative Distribution of Categorical Variables against Survival Outcomes

### 2.4.3 Inter-Variable Relationships:

Looking more closely at the data reveals that not just the individual traits but also how they interact and relate to one another are important. Variables frequently reveal a more complex tale when they are combined than when they are studied alone. Whether direct or inverse, this link between the variables can offer vital insights into the bigger picture and help forecast events more precisely. We use a correlation matrix, a visual representation of the linear correlations between various characteristics, to unravel these relationships, refer Figure 2.6. We can identify

patterns and relationships between the features of our dataset by taking a closer look at this heatmap-based matrix. Some key observations:

- **Positive correlations** indicate that as one variable increases, the other does too. For instance, there might be a positive correlation between age and certain health conditions, signifying that as patients get older, the likelihood of those conditions may also increase (thebmj, 1997).
- **Negative correlations** suggest that as one variable goes up, the other goes down (thebmj, 1997). An example could be a negative correlation between exercise frequency and the risk of certain heart ailments.
- **Zero or near-zero correlation** indicates little to no linear relationship between variables (thebmj, 1997).

However, as this section of our study comes to a close, we stand on the edge of going even further and moving from comprehending relationships to making significant judgments. This viewpoint from above has equipped us for the more focused journey that awaits us as we continue into the crucial section that follows: Key Insights from Exploratory Data Analysis (EDA). There, we will distill our observations, pinpointing the most prominent discoveries that our data has to offer.

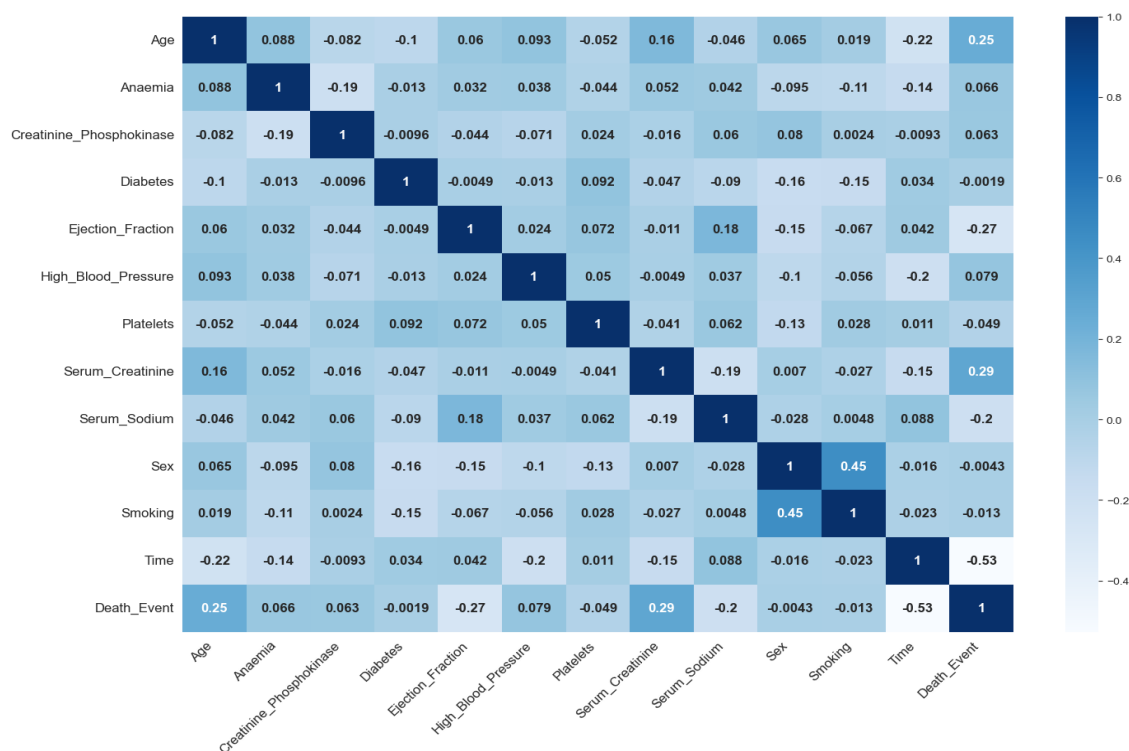


Figure 2.6: Correlation Matrix Heatmap

## 2.5 Key Insights from Exploratory Data Analysis:

In this section, we delve into the key insights gained from our Exploratory Data Analysis (EDA). By examining the relationships and patterns within our dataset, we uncover valuable information that sheds light on the factors influencing mortality due to heart failure. These insights pave the way for a more in-depth exploration of individual subsections, each providing a comprehensive examination of specific variables and their impact on heart failure-related mortality. Here, we offer an overview of the comprehensive findings from our EDA, setting the stage for a deeper dive into the details.

### 1. Age & Mortality Risk:

- **Descriptive Insight:** Cases of mortality incidence predominantly begin after the age of 40, with noticeable spikes at ages 45, 50, 60, 65, 70, 75, 80, and 90. Additionally, the mortality incidence ratio demonstrates a marked increase post the age of 70, refer Figure 2.7.
- **Quantitative Insight:** The correlation between Age and mortality incidence is 0.25, indicating a moderate positive relationship. As age advances, the risk of mortality due to heart failure escalates, refer Figure 2.6.

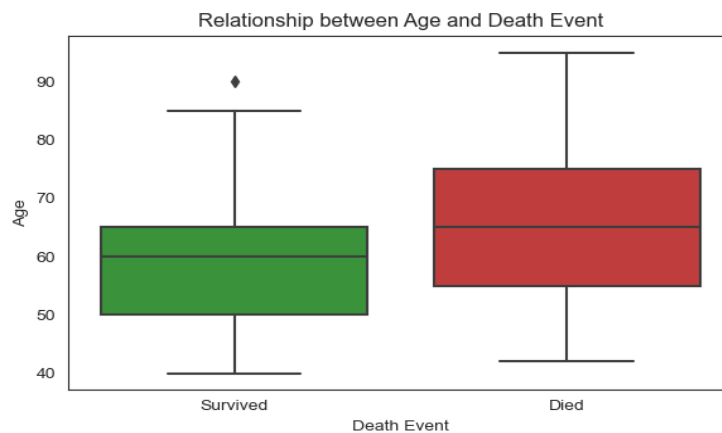


Figure 2.7: Box plot illustrating the relationship between Age and Mortality Risk

### 2. Heart Health Metrics:

- **Descriptive Insight:** The data shows that individuals with **ejection fraction** values below 40% are more likely to experience mortality related to heart failure. This means that as the ejection fraction, a measure of the heart's pumping efficiency, decreases below this threshold, the risk of mortality increases. This observation is visually represented in Figure 2.8, which provides a clear visual representation of this relationship. Another critical factor is **serum creatinine** levels. When serum

creatinine values exceed 1.2, it indicates a heightened susceptibility to mortality due to heart failure. In other words, higher levels of serum creatinine are associated with an increased risk of heart failure-related mortality. Figure 2.9, visually presents this relationship, making it easier to grasp the significance of serum creatinine in predicting mortality risk.

Similarly, the dataset shows that individuals with **serum sodium** levels falling within the range of 127-145 also face an elevated risk of heart failure-related mortality. This suggests that deviations from this range, whether too low or too high, are associated with an increased likelihood of mortality. Figure 2.10, provides a visual representation of how serum sodium values within or outside this range correlate with mortality risk.

- **Quantitative Insight:**

- **Ejection\_Fraction** and mortality incidence correlate at -0.27, underscoring that diminished ejection fraction elevates mortality risk, refer Figure 2.6.
- **Serum\_Creatinine** registers a positive correlation of 0.29 with Mortality Incidence. Elevated values signal heightened heart failure risk, refer Figure 2.6.
- **Serum\_Sodium** correlates negatively with Mortality Incidence at -0.20, suggesting reduced levels as potential risk indicators, refer Figure 2.6.

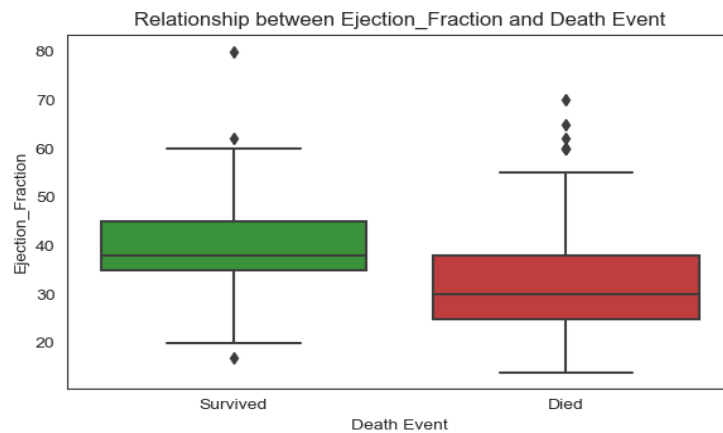


Figure 2.8: Box plot illustrating the relationship between Ejection Fraction and Mortality Risk

### 3. Follow-up Duration:

- **Descriptive Insight:** The time variable, potentially indicating the duration of medical follow-up, exhibits a strong relationship with Mortality Incidence in the range of 25-100 days, refer Figure 2.11.
- **Quantitative Insight:** A strong negative correlation of -0.53 between the duration of follow-up after heart failure detection (referred to as "Time") and the occurrence

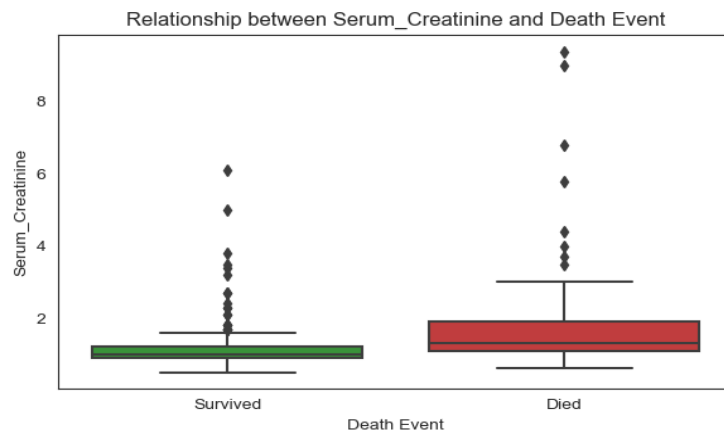


Figure 2.9: Box plot illustrating the relationship between Serum-Creatinine and Mortality Risk

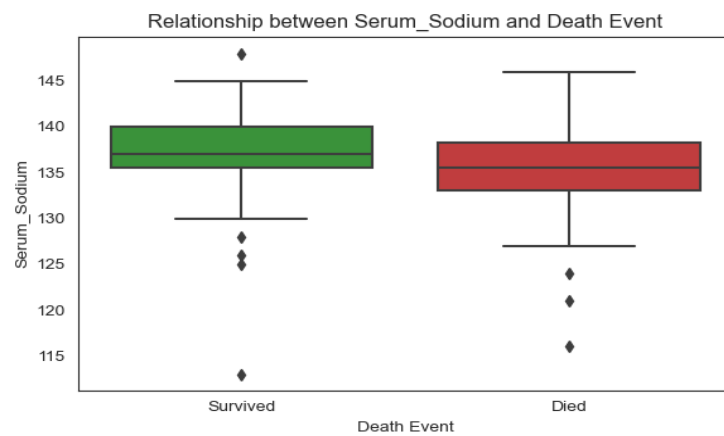


Figure 2.10: Box plot illustrating the relationship between Serum-Sodium and Mortality Risk

of mortality highlights that extended monitoring periods substantially reduce the risk of death, refer Figure 2.6.

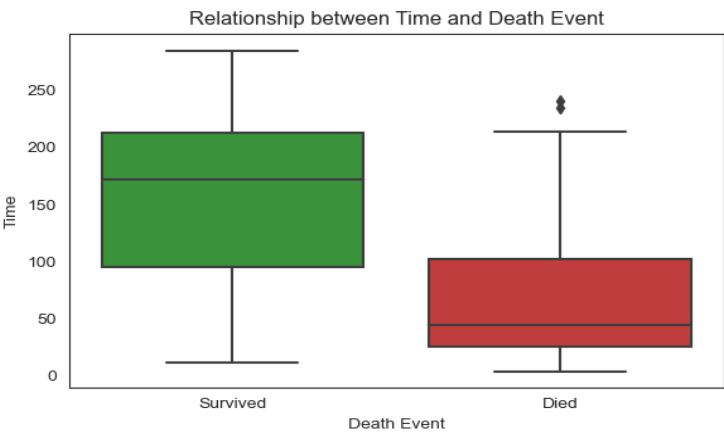


Figure 2.11: Box plot illustrating the relationship between Time and Mortality Risk

4. Gender Disparities:

- **Descriptive Insight:** The male population displays a greater tendency for heart failure-related Mortality Incidence compared to the female population, refer Figure 2.12.
- **Quantitative Insight:** The weak negative correlation of -0.0043 between Gender and Mortality Incidence provides marginal quantitative evidence to this observation, refer Figure 2.6.

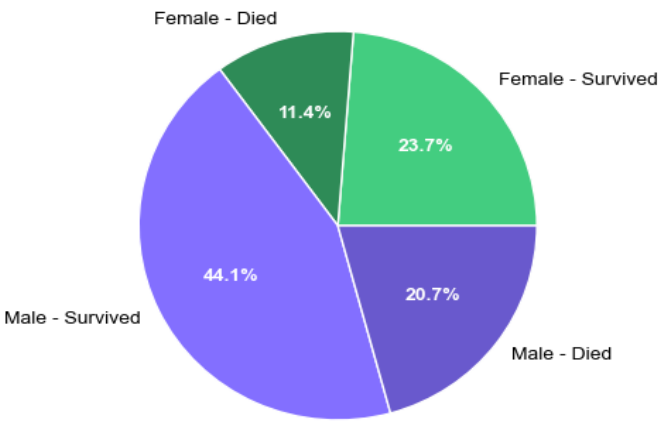


Figure 2.12: Distribution of Gender (Male and Female)

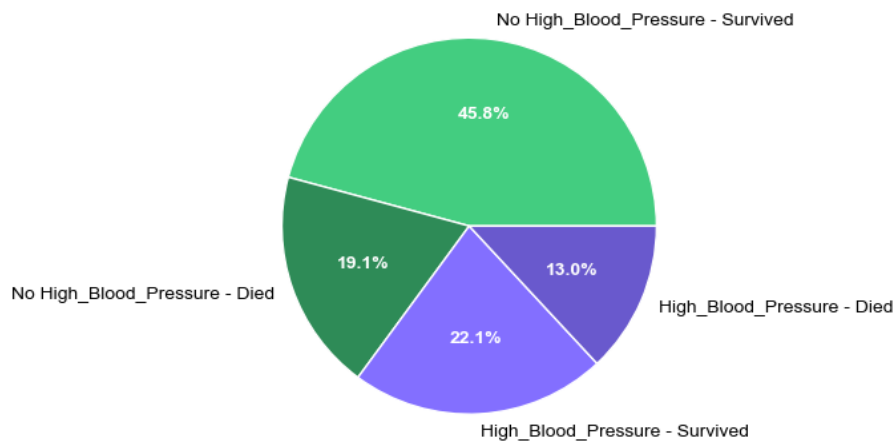
5. Underlying Health Conditions & Lifestyle Indicators:

- **Descriptive Insight:** Individuals with high blood pressure demonstrate a higher susceptibility to heart failure. Conversely, in this dataset, there's insufficient evidence

to link diabetes, anemia, and smoking directly to heart failure, despite established medical knowledge (NHS, 2022).

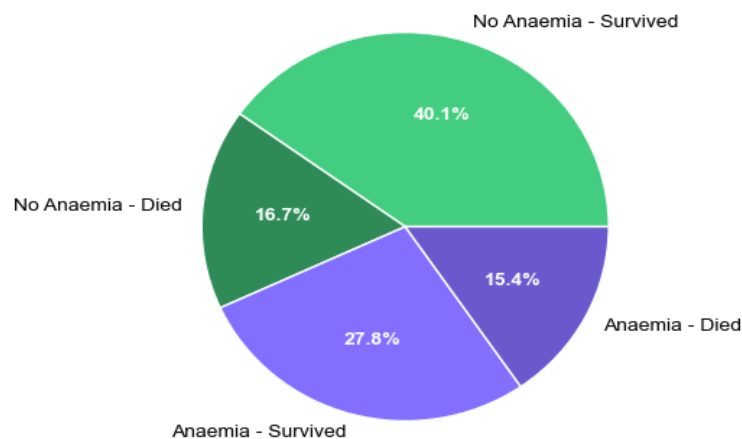
- **Quantitative Insight:**

- A positive correlation of 0.079 between **High\_Blood\_Pressure** and Mortality Incidence is evident, refer Figure 2.6.



*Figure 2.13: Distribution of High Blood Pressure*

- **Anaemia** shows a mild positive correlation with Mortality Incidence at 0.066, refer Figure 2.6.



*Figure 2.14: Distribution of Anaemia*

- **Smoking**, while a known cardiovascular risk factor (NHS, 2022), only exhibits a weak correlation of -0.013 in this dataset, refer Figure 2.6, possibly due to dataset limitations.

6. **Constraints & Future Directions:** The study highlights the intricate interplay between different parameters and the likelihood of mortality due to heart failure. Some acknowl-



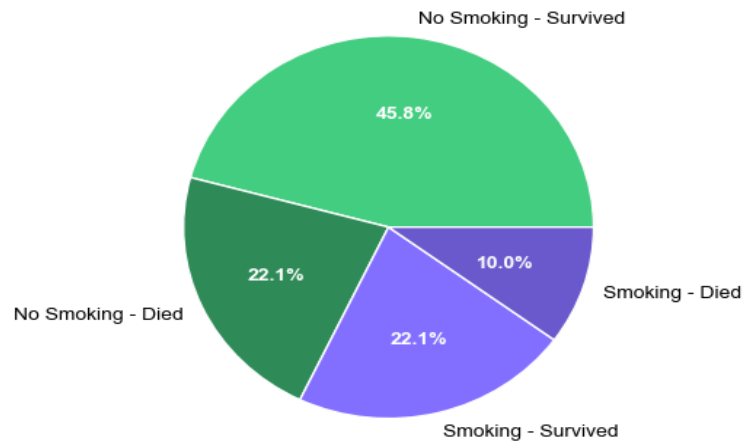


Figure 2.15: Distribution of Smoking

edged risk factors, such as diabetes, show weak or almost negligible correlations, suggesting limitations in the dataset or the need for more comprehensive analytical methods. A more encompassing model, considering interrelations and non-linear patterns among the variables, could offer a deeper understanding of the underlying causes of heart failure.

## 2.6 Data Preparation and Splitting:

In the journey towards predictive modeling, meticulous data preparation is our initial act. We commence by isolating the vital components of our dataset: the input variables that hold the key to understanding heart failure. Then, with a magician's precision, we split our dataset into two distinct groups – our guiding star, the training set encompassing 80% of our data, and the unseen audience, the testing set comprising 20%. In this narrative, the training set serves as the tutor, conveying knowledge to our models, teaching them to discern patterns and predict outcomes. Meanwhile, the testing set remains hidden, a benchmark against which our models' predictions will be assessed. This division, incorporated with the element of randomness and sealed with a random state of 1, sets the stage for our quest to decode the mysteries of heart failure prediction.

## 2.7 Standardization:

In the symphony of data analysis, every feature or variable in a dataset is like an instrument. For the entire ensemble to produce a harmonious sound, each instrument must be finely tuned. Similarly, in a dataset, if one feature has a vast numeric range, it might overpower or unduly influence the outcome when modeling, merely because of its scale, even if it isn't necessarily more important.

This is where standardization steps in, ensuring that every feature in our dataset plays its part without overshadowing the others. Specifically, features like 'Creatinine Phosphokinase',

'Ejection Fraction', 'Platelets', 'Serum Creatinine', and 'Serum Sodium' had varying scales, which could introduce inadvertent biases.

To align these features and enable them to contribute fairly to our machine learning models, we employed the method of standard scaling (Jain, A et.al, 2005). It's a process that transforms our data in such a way that its distribution has a mean of 0 and a standard deviation of 1. In essence:

$$Z = \frac{X - \mu}{\sigma} \quad (2.1)$$

where:

$Z$  : Standardized value

$X$  : Original value

$\mu$  : Mean of the feature

$\sigma$  : Standard deviation of the feature

So, with these steps, our dataset was transformed from a raw piece of stone into a polished masterpiece, ready to narrate its tale of heart health. Onward we go, deeper into the story.

As we delve into the next chapter Predictive Modeling for Heart Failure, we transition from understanding the intricate relationships among variables to harnessing this knowledge for predictive modeling.

## Chapter 3

# Predictive Modeling for Heart Failure

The essence of medicine lies not just in curing but also in predicting. Predictive modeling, at its core, is a statistical technique that uses algorithms and data to predict future outcomes. These models process current and historical data to forecast a future event, ideally with significant accuracy. In the realm of heart health, these predictions can mean the difference between timely intervention and a missed opportunity. The advent of machine learning has supercharged the capabilities of predictive modeling. Instead of relying solely on historical patterns, these models consider numerous variables and intricate patterns in the data, often revealing insights that might remain hidden to traditional analysis.

Given the vast landscape of algorithms available for predictive modeling, it's crucial to be strategic in our selections. Our choices have been influenced heavily by the insights derived from the Exploratory Data Analysis (EDA). Specifically, our decision to implement certain modeling techniques was based on the following considerations:

- **Binary Outcome:** We observed that the outcome variable, pertaining to the mortality due to heart failure, presents a binary distinction.
- **Data Imbalance:** Digging deeper, we discerned an imbalance in the distribution of our outcome classes. Some predictive models, by default, might struggle with imbalanced datasets, thereby emphasizing the need for techniques either inherently adept at handling such disparities or those that can be fine-tuned accordingly.
- **Correlational Insights:** Certain attributes, especially serum creatinine and age, unveiled a marked correlation with mortality outcomes. This implied a need for models that can duly weigh these influential features in their decision-making process.
- **Complex Relationships:** While some variables presented linear relationships with the outcome, others, such as ejection fraction and time, demonstrated non-linear dynamics. Such observations underscored the importance of considering models capable of capturing these intricate patterns.

- **Interactions Among Variables:** Deeper analysis revealed potential cumulative outcomes. For example, the combined effect of age, serum creatinine, and ejection fraction may have greater significance than either factor acting alone. This interaction of variables showed how useful models that can account for such interacting effects could be.

With these insights from our EDA at hand, we chose four powerful techniques—Support Vector Machines (SVM), Decision Trees & Random Forest, and Xgboost classifier—to help us navigate the path to predictive modeling. These models each have particular advantages that fit our dataset's particularities well. In the upcoming sections, we'll examine each model's justification for use and examine how it fits into the story that our data tells. Let's start this investigation now and examine each predictive model separately.

### 3.1 Support Vector Machines (SVM):

Support Vector Machines (SVM) is a popular supervised machine learning algorithm primarily designed for binary classification tasks, but it can also handle multiclass problems. SVM aims to find the optimal hyperplane that best divides a dataset into classes. The strength of SVM lies in its ability to manage high-dimensional data and its versatility in modeling both linear and non-linear relationships (Cortes et.al, 1995). Given the binary outcome observed in our dataset and the potential high-dimensional nature of the feature space, SVM can be an excellent fit. Its ability to maximize the margin between the classes can offer a distinct advantage in scenarios where the distinction between the outcome classes is not Clearly evident.

#### 3.1.1 Justification for using SVM:

1. **Binary Outcome:** Our dataset is primarily focused on a binary outcome, making SVM a natural choice.
2. **Correlational Insights:** Given the significant correlations observed, especially with variables like serum creatinine and age, SVM can help in capturing these linear (or non-linear, with kernels) relationships efficiently.
3. **Complex Relationships:** With the help of kernel tricks, SVM can model the intricate relationships in our data, such as those exhibited by attributes like ejection fraction.
4. **High-Dimensional Feature Space:** If interactions among features are considered, the dimensionality can grow, making SVM a suitable choice due to its strength in high-dimensional spaces.

Support Vector Machine (SVM) requires the data to be clean. This means that the input data should be devoid of noise and there shouldn't be any missing values. It's essential for the model's optimal performance. By its very design, the standard SVM algorithm is tailored for

binary classification tasks. So, it's inherently adapted to problems with two outcomes. That said, there are extensions of SVM that can handle multiclass problems. One of the prerequisites for using SVM is feature scaling. This is because SVM is sensitive to the scale of the data. As such, it's typically expected that the data undergo preprocessing steps like standardization before being fed into the algorithm. Lastly, SVM is equipped to handle data with a large number of features. It's adept at performing in scenarios where the dimensionality is high, even when there are more features than there are data points (Cortes et.al, 1995).

### 3.1.2 Mathematical Overview:

The main objective of SVM is to find the best hyperplane that separates the data into its respective classes. For a linearly separable dataset, the optimal hyperplane is the one that has the maximum margin between the two classes. Mathematically, the hyperplane is defined as:

The equation (Shuzhan Fan, 2018) of the hyperplane is given by:

$$w \cdot x + b = 0$$

Where:

$w$  : the weight vector perpendicular to the hyperplane

$b$  : the bias

The decision function is defined as:

$$f(x) = w \cdot x + b$$

After having set the stage with the theoretical understanding of SVM, we applied the algorithm to our dataset. Leveraging Jupyter Notebook, a widely accepted platform for data analysis and machine learning tasks, we were able to extract a performance summary, Table 3.1 that offers a clear snapshot of SVM's predictive prowess on our heart failure dataset.

Metric/Class	Survived	Died	Weighted Average
Precision	0.91	0.77	0.88
Recall	0.93	0.71	0.88
F1-Score	0.92	0.74	0.88
Support	46	14	60

Table 3.1: Performance Metrics of the Support Vector Machine (SVM) Classifier

#### Interpretation and Insights:

- **Precision:** This metric gives us an understanding of the model's accuracy in terms of false positives. For class 0 (Survived), the SVM model achieved a high precision of 0.91,

suggesting that when the model predicts a patient will survive, it's correct about 91% of the time. For class 1 (Died), the precision is 0.77, indicating a slightly lower confidence when predicting fatalities.

- **Recall:** Recall represents the model's capability to correctly identify true positives. A recall of 0.93 for class 0 implies that the model successfully identified 93% of the actual survivors. On the other hand, for class 1, the recall of 0.71 suggests that it identified 71% of the actual fatalities.
- **F1-Score:** This is the harmonic mean of precision and recall, giving a balanced measure of the model's overall performance. An F1-Score close to 1 is ideal. For our model, class 0 has an F1-Score of 0.92 and class 1 has 0.74, indicating a robust performance, especially for predicting survivors.
- **Support:** This indicates the actual number of occurrences of the class in the specified dataset. Here, 46 instances belong to class 0 and 14 to class 1.

The overall accuracy of the SVM model stands at 0.88, which means that in 88% of the cases, the model made the correct prediction. These insights set the foundation for understanding the model's capabilities and areas of improvement. As we delve deeper into other models, such comparisons will serve as a compass guiding our model selection and tuning endeavors.

## 3.2 Decision Trees:

Decision Trees (DT) are intuitive and interpretable machine learning models used for both classification and regression tasks. At their core, decision trees split the data into subsets based on feature values. These splits are made in a hierarchical manner, where each internal node represents a feature, each branch denotes a decision rule, and each leaf node holds an outcome. For our dataset, given the complex relationships between variables and the outcome, a decision tree can provide a straightforward visualization of the decisions made at each step. The inherent interpretability of decision trees is especially beneficial when trying to understand the key features affecting heart failure, the study on DT is (Vrutti Tanna, 2020) based.

### 3.2.1 Justification for using Decision Trees:

1. **Interpretable Model:** Given the complexity of medical data, having an easily interpretable model like a decision trees can be crucial for clinical decision-making.
2. **Handles Non-linearity:** As observed in variables like ejection fraction, the relationships can be non-linear. Decision trees inherently handle such complexities.
3. **Feature Importance:** Decision trees provide insights into which features are critical for predictions, aligning well with the correlational insights from our EDA.

Decision trees come with a set of flexible assumptions (Vrutti Tanna, 2020) that make them quite versatile. Firstly, they are adept at handling both categorical and continuous data. For continuous data, decision trees typically set thresholds to categorize them. Importantly, they don't make any assumptions about how the data is distributed, which classifies them as non-parametric. This means they don't get restricted by a fixed type of distribution, like some other models do. Additionally, they aren't limited to just modeling linear relationships. Even if data points curve and twist in all directions, decision trees can capture those non-linear patterns.

### 3.2.2 Mathematical Overview:

The decision trees algorithm works by recursively partitioning the dataset based on the feature that provides the best separation, as measured by an impurity metric. Common metrics include: Gini impurity:

$$Gini(p) = 1 - \sum_{i=1}^k p_i^2$$

Entropy:

$$Entropy(p) = - \sum_{i=1}^k p_i \log(p_i)$$

Classification error:

$$Error(p) = 1 - \max(p_i)$$

Where  $p_i$  is the proportion of samples belonging to class  $i$ .

After diving deep into the theoretical aspects of the Decision Trees algorithm, it was time to test its strength against our dataset. With the aid of Jupyter Notebook – an industry-recognized tool for data processing and machine learning – we derived a performance snapshot, Table 3.2 that highlights the efficiency of the Decision Tree in predicting heart failure outcomes.

Metric/Class	Survived	Died	Weighted Average
Precision	0.90	0.75	0.86
Recall	0.93	0.64	0.87
F1-Score	0.91	0.69	0.86
Support	46	14	60

Table 3.2: Performance Metrics of the Decision Trees

#### Interpretation and Insights:

- **Precision:** Precision reflects the model's accuracy in minimizing false positives. For class 0 (Survived), the Decision Trees model boasts a commendable precision of 0.90. This means that when predicting survival, it's on point about 90% of the time. When it predicts class 1 (Died), it has a precision of 0.75, suggesting a little more caution in predictions regarding fatalities.

- **Recall:** This metric portrays the model's proficiency in pinpointing true positives. With a recall of 0.93 for class 0, it confidently captured 93% of the actual survivors. For class 1, a recall of 0.64 indicates it managed to identify roughly 64% of the genuine fatalities.
- **F1-Score:** Serving as the harmonic mean of precision and recall, the F1-Score provides a consolidated view of the model's efficacy. An optimal F1-Score is close to 1. In our case, the model scored 0.91 for class 0 and 0.69 for class 1, showcasing a strong performance, particularly in predicting survival outcomes.
- **Support:** This represents the true count of occurrences for each class in our dataset. Here, we have 46 instances for class 0 and 14 for class 1.

To conclude, the Decision Trees model's overall accuracy stands impressively at 0.87. This signifies that in 87% of the scenarios, the model's predictions were spot on.

### 3.3 Random Forest:

Random Forest is an ensemble learning method, combining multiple decision trees to make more accurate and stable predictions. The model's strength lies in its ability to mitigate the overfitting commonly associated with individual decision trees, thanks to its ensemble nature, see (Stacey Ronaghan, 2018).

#### 3.3.1 Justification for Using Random Forest:

1. **Handling Overfitting:** Individual decision trees tend to overfit the data. Random Forest, with its ensemble approach, can average out biases and reduce variance, making predictions more generalizable.
2. **Feature Importance:** Random Forest provide a ranked list of feature importance, which can be invaluable in understanding which variables play a critical role in prediction.
3. **Model Complexity:** It can capture complex interactions between variables by combining the predictions of individual trees.

Random Forest operate on a couple of foundational assumptions (Stacey Ronaghan, 2018). Firstly, there's an underlying belief that a more significant number of trees in the ensemble will lead to a more robust prediction. This is termed as the 'Large Ensemble' assumption, emphasizing the strength derived from the collective rather than individual trees. Secondly, there's the 'Decorrelation' assumption, which underscores the necessity for individual decision trees to maintain a degree of independence from one another. To achieve this, randomness is introduced at two crucial stages: during the bootstrap sampling of data and the selection of features. Both these steps ensure that each tree in the ensemble captures different aspects or distinctions of the data, thereby enhancing the overall prediction's reliability and robustness.



### 3.3.2 Mathematical Overview:

Random Forest build multiple decision trees during training and produce the class that is the mode of the classes output by individual trees for classification tasks. The algorithm involves:

- **Bootstrap Sampling:** A subset of data is taken (with replacement) from the training dataset.
- **Feature Randomness:** Instead of searching for the most optimal split among all features, a random subset of those features is selected.
- **Tree Building:** Using the above subsets, a decision tree is constructed. No trimming takes place, ensuring the full growth of trees.
- **Final Prediction:** For classification, the mode of all the predictions from individual trees is returned.

Building on the foundational concept of Decision Trees, Random Forest introduces an ensemble approach, which aggregates multiple decision trees to yield a more accurate and stable prediction. Harnessing the power of Jupyter Notebook, we have unveiled a performance summary, Table 3.3 that emphasizes the proficiency of Random Forest in discerning heart failure outcomes.

Metric/Class	Survived	Died	Weighted Average
Precision	0.93	0.73	0.89
Recall	0.91	0.79	0.88
F1-Score	0.92	0.76	0.88
Support	46	14	60

Table 3.3: Performance Metrics of the Random Forest

#### Interpretation and Insights:

- **Precision:** Precision quantifies the model's assurance in avoiding false positives. For class 0 (Survived), the Random Forest model shows an impressive precision of 0.93. This translates to the model being accurate about 93% of the time when it predicts survival. For fatalities (class 1), the model has a precision of 0.73, suggesting careful attention in predicting deaths.
- **Recall:** This metric showcases the model's capability in correctly identifying true positives. Achieving a recall of 0.91 for class 0, the model wisely recognized 91% of the actual survivors. Meanwhile, for class 1, the recall stands at 0.79, indicating it identified nearly 79% of the real fatalities.
- **F1-Score:** A balanced measure of precision and recall, the F1-Score offers a comprehensive view of the model's overall robustness. An F1-Score influencing 1 is seen as

exemplary. Here, the Random Forest model posted scores of 0.92 for class 0 and 0.76 for class 1, revealing a commendable overall performance, especially in identifying those who survived.

- **Support:** This metric provides the actual number of instances for each class in the test dataset. In this case, class 0 accounts for 46 instances, while class 1 represents 14.

In summary, Random Forest, with their ensemble mechanism building upon the Decision Trees model, accomplished an overall accuracy of 0.88. This means that the model delivered accurate predictions in 88% of the test cases.

### 3.4 XGBoost classifier:

XGBoost, short for "Extreme Gradient Boosting", is an advanced implementation of gradient boosted decision trees. The technique has risen in popularity due to its remarkable predictive performance across a variety of problems and datasets. At its core, XGBoost builds an ensemble of decision trees in a sequential manner, with each tree trying to correct the errors of its predecessor, see (Chen, T et.al, 2016).

Our dataset, with its diverse feature interactions and potential non-linear relationships, serves as a suitable playground for XGBoost. The algorithm's capability to weigh features, manage missing data, and its built-in regularization makes it a top contender for predictive modeling tasks in this context.

#### 3.4.1 Justification for using XGBoost:

1. **Feature Interactions:** Given the complexity and potential interactions among features, XGBoost can adaptively learn from the data, focusing on more influential predictors and interactions.
2. **Regularization:** The in-built regularization avoids overfitting, making it more generalizable to new, unseen data.
3. **Flexibility:** XGBoost can handle missing data, and its ability to be parallelized makes it faster, especially for larger datasets.

XGBoost makes certain key assumptions (Chen, T et.al, 2016) when modeling data. First, it goes a step beyond traditional gradient boosting by adding L1 and L2 regularization terms to its objective function. This means it has a built-in check against creating overly complex models that might overfit the data. Essentially, it tries to strike a balance between fitting the data well and keeping the model simple. Second, XGBoost operates on the idea of 'boosting'. It believes that by adding new trees, it can correct the mistakes made by the earlier ones, gradually refining its predictions as more trees are introduced.

### 3.4.2 Mathematical Overview:

XGBoost employs the gradient boosting framework where new trees are fit to the negative gradient (or error) of the loss function. This iterative addition aims to minimize the overall prediction error. The objective function in XGBoost comprises the sum of the loss function and a regularization term, which helps in penalizing complex models, thus preventing overfitting.

XGBoost optimizes the following objective function:

$$Obj(\Theta) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{i=1}^T \Omega(f_i)$$

Where:

$l$  : Loss function that measures the difference between the prediction and the true data

$\hat{y}_i$  : Prediction for sample  $i$

$\Omega$  : Regularization term to penalize the complexity, defined as:

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2$$

Here:

$T$  : Number of leaves in the tree

$w$  : Represents the scores on the leaves

$\gamma$  : Regularization parameter

$\lambda$  : Regularization parameter

The results crystallized in a performance summary, Table 3.4, reflecting the efficacy of XGBoost in predicting heart failure outcomes.

Metric/Class	Survived	Died	Weighted Average
Precision	0.89	0.64	0.83
Recall	0.89	0.64	0.83
F1-Score	0.89	0.64	0.83
Support	46	14	60

Table 3.4: Performance Metrics of the XGBoost Classifier

#### Interpretation and Insights:

- **Precision:** For class 0 (Survived), the precision stands at 0.89, indicating when XGBoost predicts a survival, it's accurate around 89% of the time. For class 1 (Died), the precision drops to 0.64, denoting a lowered assurance when predicting fatalities.
- **Recall:** A recall of 0.89 for class 0 reveals that 89% of actual survivors were correctly

identified by the model. For class 1, the recall mirrors the precision at 0.64.

- **F1-Score:** A balanced measure of precision and recall, the F1-Score for class 0 is 0.89, dropping to 0.64 for class 1. This again reiterates a stronger predictive power for survivors compared to fatalities.
- **Support:** This metric displays the actual occurrences of each class in the dataset, with 46 instances of class 0 and 14 of class 1.

Conclusively, the XGBoost model delivers an accuracy of 0.83, correctly predicting outcomes in 83% of the cases.

As we wrap up this chapter, it's evident that our journey doesn't conclude here. In the chapters to come, we will delve into the critical aspects of Performance Evaluation and Model Optimization. Our focus will shift towards comprehensively assessing the models we've developed and fine-tuning them for enhanced predictive accuracy. So, stay tuned as we continue our quest for precision and insight in the realm of heart failure prediction.

## Chapter 4

# Performance Evaluation and Model Optimization

Building on our earlier investigation, in which we identified the core components of each selected prediction model, we now delve further into evaluating their performances. Raw accuracy numbers and performance indicators give a quick overview of a model's capabilities, but a more detailed examination provides information that is essential for real-world applications. Understanding the intricate details of each prediction, the differences between sensitivity and specificity, and a model's tolerance to different data distributions are crucial.

The multidimensional evaluation metrics Confusion Matrices, Receiver Operating Characteristic (ROC) Curves, and Precision-Recall Curves will be the focus of this chapter. These measures will provide a more thorough picture of each model's capabilities by highlighting its advantages and potential weaknesses. Aware that models may frequently be enhanced by altering their parameters, we'll also explore optimization strategy like Grid Search to fully utilize the power of our predictive tools. The goal of this chapter is to carefully evaluate improve, and then identify the best model for our dataset, laying the groundwork for possible deployment in real-world scenarios.

### 4.1 Confusion Matrix:

A confusion matrix, often used with classification models, is a tabular representation that describes the performance of a predictive model on a set of data for which the true values are known. It provides insights into the number of correct and incorrect predictions made by the model, allowing for a more granulated assessment of its capabilities. The reference is taken from (Sarang Narkhede, 2018).

#### 4.1.1 Components of a Confusion Matrix:

A standard confusion matrix for binary classification comprises four primary components:

- **True Positives (TP):** Instances that were positive and were correctly predicted as positive by the model.
- **True Negatives (TN):** Instances that were negative and were correctly predicted as negative by the model.
- **False Positives (FP):** Instances that were negative but were incorrectly predicted as positive by the model.
- **False Negatives (FN):** Instances that were positive but were incorrectly predicted as negative by the model.

In the following Table 4.1, we will present the confusion matrix in tabular format, providing a comprehensive view of the model's true positives, true negatives, false positives, and false negatives.

Actual / Predicted	Predicted Positive	Predicted Negative
Actual Positive	TP (True Positive)	FN (False Negative)
Actual Negative	FP (False Positive)	TN (True Negative)

*Table 4.1: Confusion Matrix*

In the realm of predictive modeling, the accuracy of our models holds immense importance, particularly when dealing with critical healthcare decisions such as predicting heart failure. To assess the effectiveness of our models, we rely on a pair of metrics that can make or break the utility of our predictions: Type I and Type II errors.

#### **The Significance of Type I and Type II Errors in Model Evaluation:**

- **Type I Error (False Positive):** These errors occur when our model erroneously identifies a negative case as positive. While a Type I error is indeed an error, its consequences are usually less severe. In the context of heart failure prediction, it might lead to additional tests or medical examinations for patients flagged by the model. While an inconvenience, this can serve as a safety net, ensuring that potential heart failure cases do not go unnoticed.
- **Type II Error (False Negative):** On the other hand, Type II errors are of grave concern, especially in healthcare predictions. A Type II error transpires when our model fails to recognize a positive case as such, essentially missing the signs of a potentially critical condition like heart failure. In these instances, early detection and timely intervention are vital. A Type II error might mean that individuals with heart failure risk factors are not receiving the necessary medical attention they urgently require. While the numerical count of Type II errors might appear modest, the real-world implications can be profound, potentially affecting patient outcomes and well-being.

As we recognize the critical role of Type I and Type II errors in model evaluation, we now stand ready to examine the performance of each model in our set. In the forthcoming sections, we will employ the powerful tool of the confusion matrix to dissect our models' predictions, gaining valuable insights into their strengths and limitations.

#### 4.1.2 Confusion Matrix for SVM:

The Figure 4.1 shows the confusion matrix for SVM. The SVM model demonstrates a promising overall performance. A high True Positive rate of 43 indicates that the model correctly identified 43 instances as positive out of the actual positive cases. Additionally, a robust True Negative rate of 10 means that out of the actual negative cases, 10 of them were correctly predicted as negative. These numbers suggest a strong predictive capability of the SVM model for both classes.

##### Type I and Type II Errors:

- **Type I Error (False Positive):** The SVM model produced 4 false positives. This means that in 4 instances, the model incorrectly predicted a negative case as positive. While a false positive is an error, its implications might be less severe as it might just lead to additional tests or examinations.
- **Type II Error (False Negative):** The model had 3 false negatives. This is particularly concerning, especially in the context of predicting heart failure. A false negative means that the model failed to identify 3 potential heart failure risks, suggesting that these patients might be overlooked in real-life scenarios. This can be detrimental as early detection and intervention are critical in heart failure cases. Thus, while the number might seem small, the implications of these misclassifications can be significant.

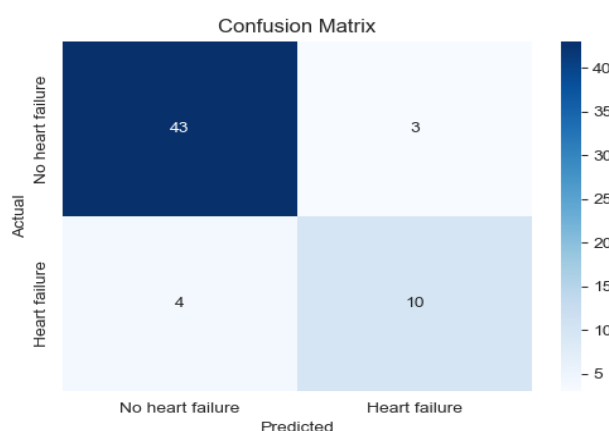


Figure 4.1: Confusion Matrix for SVM

### 4.1.3 Confusion Matrix for Decision Trees:

The Decision Trees model showcases a commendable overall predictive performance Figure 4.2. With a True Positive rate of 43, it correctly identified 43 cases as positive out of the actual positive instances. Furthermore, a True Negative count of 9 denotes that the model accurately predicted 9 out of the actual negative cases as negative. These figures suggest that the Decision Trees model has a robust predictive capability across both classes, similar to the SVM model.

#### Type I and Type II Errors:

- **Type I Error (False Positive):** The Decision Trees model resulted in 5 false positives. This implies that there were 5 occasions where the model erroneously flagged a negative instance as positive. In a clinical context, while these cases would lead to unnecessary additional evaluations, they aren't as critical as false negatives. Nevertheless, any misclassification can strain medical resources and patient trust.
- **Type II Error (False Negative):** With 3 false negatives, the model missed identifying 3 patients who were at potential risk of heart failure. This is especially worrisome in a medical setting, given the importance of early detection in heart failure cases. Missing these cases could potentially mean delaying essential interventions, leading to adverse patient outcomes.

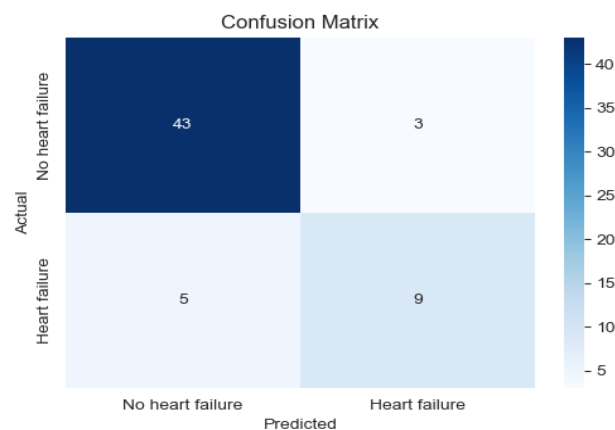


Figure 4.2: Confusion Matrix for Decision Trees

### 4.1.4 Confusion Matrix for Random Forest:

The Random Forest model delivers an impressive predictive ability, Figure 4.3. With a True Positive count of 42, the model correctly identified 42 actual positive cases. Moreover, a True Negative tally of 11 signifies the accurate classification of 11 out of the actual negative instances. These numbers reflect the model's robust capability to predict across both classes and show the



strength of ensemble techniques like Random Forest in medical diagnostics.

#### Type I and Type II Errors:

- **Type I Error (False Positive):** The Random Forest model produced 3 false positives. This indicates that on 3 occasions, the model mistakenly predicted negative instances as positive. While this could lead to additional, possibly unjustified medical interventions, it's not as grave a concern as false negatives in this context. Still, it's crucial to minimize such errors to optimize healthcare resources.
- **Type II Error (False Negative):** With 4 false negatives, the model overlooked 4 patients who were indeed at risk of heart failure. This misclassification is critical, especially in heart failure scenarios where timely intervention can be life-saving. Missing out on such cases can result in delayed treatments, potentially compromising patient care.

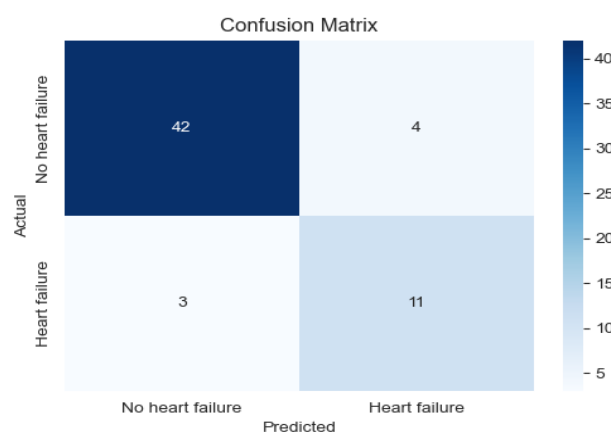


Figure 4.3: Confusion Matrix for Random Forest

#### 4.1.5 Confusion Matrix for Xgboost Classifier:

The XGBoost Classifier demonstrates commendable predictive capabilities, Figure 4.4. With a True Positive count of 41, the model rightly flagged 41 cases as positive. Additionally, the model accurately pinpointed 9 negative instances, signifying its balanced capability to discern across both classes. Given that XGBoost is renowned for its advanced gradient boosting mechanism, these numbers underline its effectiveness, especially in predictive medical settings.

#### Type I and Type II Errors:

- **Type I Error (False Positive):** The XGBoost model recorded 5 false positives, indicating that on five occasions, the model erroneously identified negative instances as positive.

Such misclassifications, while less detrimental than false negatives in heart failure prediction scenarios, might lead to excessive medical examinations, contributing to unnecessary healthcare expenditure.

- **Type II Error (False Negative):** The presence of 5 false negatives in the XGBoost model's predictions is a matter of concern. These instances denote that the model missed five patients who were genuinely at risk. This type of error is especially worrisome in the realm of heart failure predictions, where failing to identify at-risk individuals can result in missed or delayed treatments, compromising patient health.

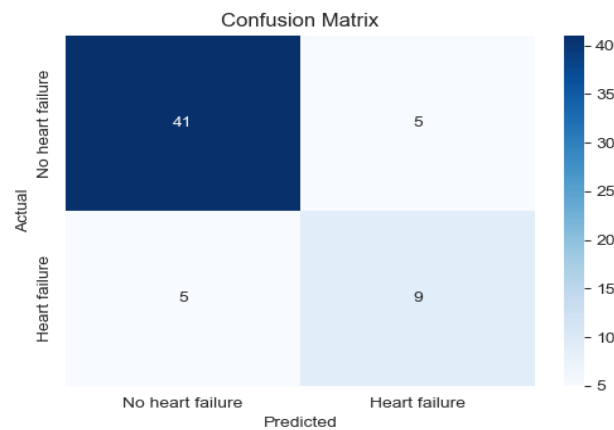


Figure 4.4: Confusion Matrix for Xgboost Classifier

## 4.2 ROC Curve Analysis:

The Receiver Operating Characteristic (ROC) curve is a graphical representation that illustrates the performance of a binary classification system. The true positive rate (Sensitivity) is plotted on the Y-axis against the false positive rate (1-Specificity) on the X-axis. The area under the curve (AUC) provides a single metric that evaluates the model's overall classification ability. The reference for this topic is taken from (SPSS Statistics, 2021).

### Interpret AUC values:

- **AUC = 1.0:** Perfect classifier; the model separates positive and negative cases perfectly.
- **AUC > 0.5:** Better than random guessing; the model has some discriminative power.
- **AUC = 0.5:** Equivalent to random guessing; the model has no discriminative power.
- **AUC < 0.5:** Worse than random guessing; the model performs inversely.

In the upcoming section, we will delve into ROC analysis for each model. Here, we will examine the outcomes, offer interpretations, and present graphical representations of the results.

### 4.2.1 ROC Curve Analysis for SVM and XGBoost:

The Receiver Operating Characteristic (ROC) curve, a critical evaluation tool for binary classification models, showcases the compromise between the true positive rate and the false positive rate. Analyzing the curves for SVM and XGBoost provides us with revealing insights. The SVM, with its AUC score of 0.91, demonstrates a steep and commendable rise in its true positive rate for only a minor surge in the false positive rate, refer Figure 4.5. This implies that the SVM, in the context of heart failure predictions, effectively maximizes correct identifications while minimizing erroneous ones.

On the other hand, XGBoost, with an AUC of 0.93, stands slightly superior, refer Figure 4.6. This gradient boosted algorithm's sequential learning technique ensures that its curve approaches closer to the ideal top-left corner, demonstrating its heightened ability to predict heart failure. In summary, both classifiers exhibit strong predictive performances; however, XGBoost holds a marginal advantage over SVM in terms of its ROC curve positioning and AUC score.

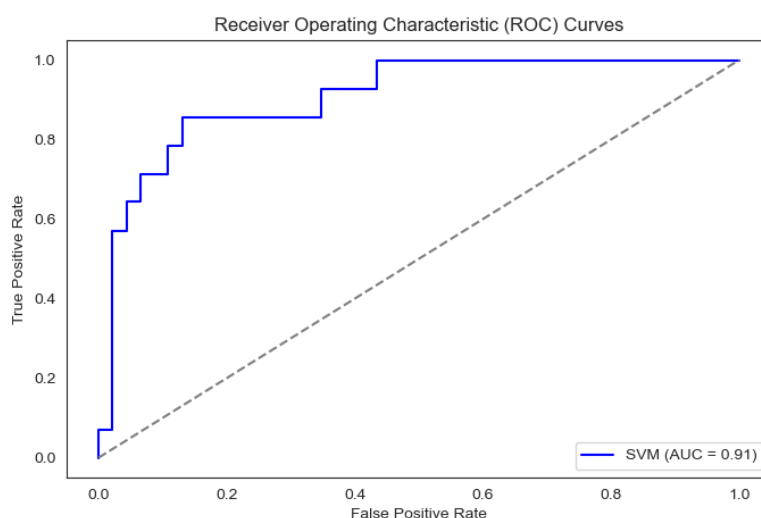


Figure 4.5: Receiver operating Characteristic curves for SVM

### 4.2.2 ROC Curve Analysis for Decision Trees and Random Forest:

The ROC curve clearly shows a model's classification abilities by looking at the balance between the true positive rate (sensitivity) and the false positive rate (1-specificity). When comparing Decision Trees and Random Forest, we can see clear differences, refer Figure 4.7.

For Decision Trees, an AUC score of 0.87 indicates a decent model, but it's not perfect. This model can identify true positives effectively, but it might also pick up some false positives. The curve probably goes up at a steady pace, reflecting the model's simple method of making predictions based on individual feature thresholds.

On the other hand, Random Forest, with an AUC of 0.95, show a much better performance.

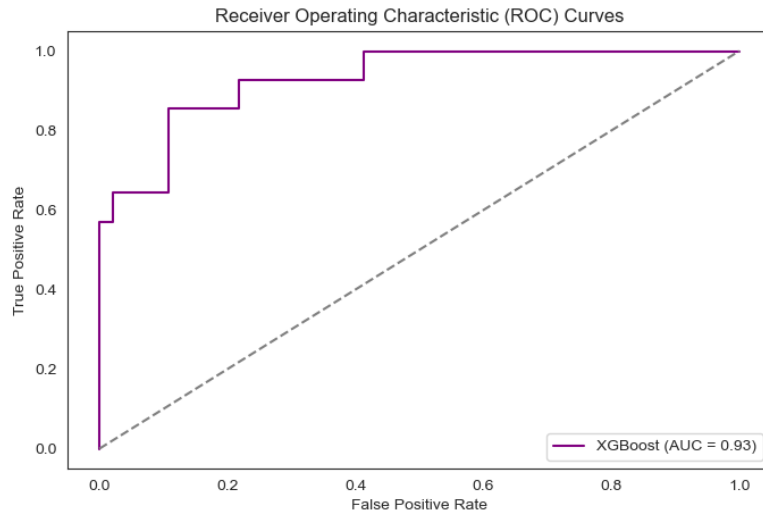


Figure 4.6: Receiver operating Characteristic curves for XGBoost

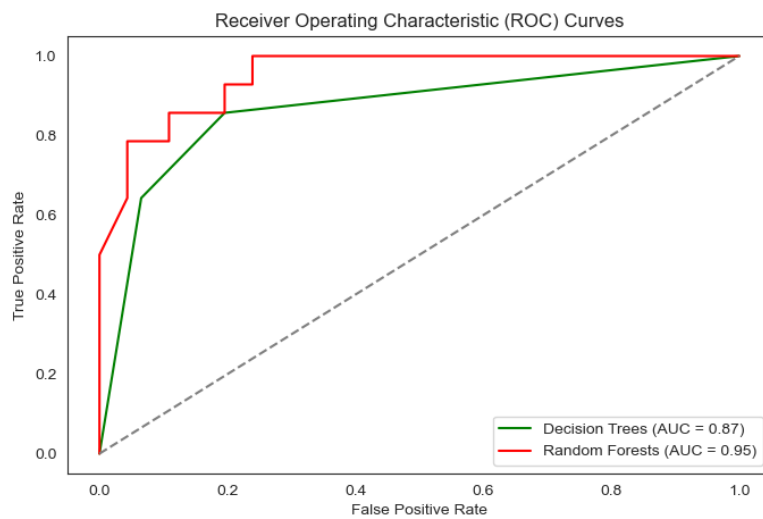


Figure 4.7: Receiver operating Characteristic curves for Decison Trees and Random Forest

Random Forest use the combined strength of many decision trees, which helps in reducing errors and the chances of overfitting. The high AUC suggests that the ROC curve for this model goes up quickly close to the y-axis, which means fewer false positives. So, when it comes to predicting heart failure, the Random Forest model seems to be a more trustworthy tool, capable of finding real cases with fewer mistakes.

In short, while Decision Trees provide a good classification performance, the combined approach of Random Forest makes them perform even better, as seen in their ROC curve and AUC score.

### 4.3 Precision-Recall Curve Analysis:

Precision-Recall curve analysis is a crucial assessment tool, especially when we're dealing with imbalanced datasets. Precision and recall, both vital in evaluating classification models, offer distinct insights into their performance. Precision measures the model's accuracy in predicting positive instances by calculating the ratio of true positives to the total predicted positives. High precision signifies a low occurrence of false positives, indicating precise identification of positive cases. Conversely, recall, also known as sensitivity or the true positive rate, assesses the model's ability to capture all actual positive instances, determined by the ratio of true positives to the total actual positives. High recall indicates a minimal rate of false negatives, showcasing the model's proficiency in identifying most positive cases. Together, these metrics provide a comprehensive assessment of a model's classification performance, especially valuable in scenarios where the impact of false positives and false negatives varies significantly, see (Teemu Kanstrén, 2020).

**Creating the Precision-Recall Curve:** To create a Precision-Recall curve, we begin by setting a classification threshold, usually 0.5 in binary classification. We calculate precision and recall at this threshold and then vary it across a range from 0 to 1. This variation yields different precision and recall values, which are plotted as points on the Precision-Recall curve. This curve visually illustrates how precision and recall change with different thresholds, aiding in the evaluation and selection of an appropriate operating point for the model.

**Interpreting the Precision-Recall Curve:** The Precision-Recall curve typically appears as a uneven line, commencing at the point (0,0) and concluding at (1,1). Its performance is quantified by the area under the curve, denoted as AUC-PR. A larger AUC-PR value signifies superior model performance.

As discussed in the ROC Curve Analysis section, evaluating a model's performance on different metrics ensures we select the most robust and reliable model for our application. For our heart failure prediction task, the models performed as follows, refer Figure 4.8:

SVM showcased an AUC score of 0.78 on the Precision-Recall curve, demonstrating a good

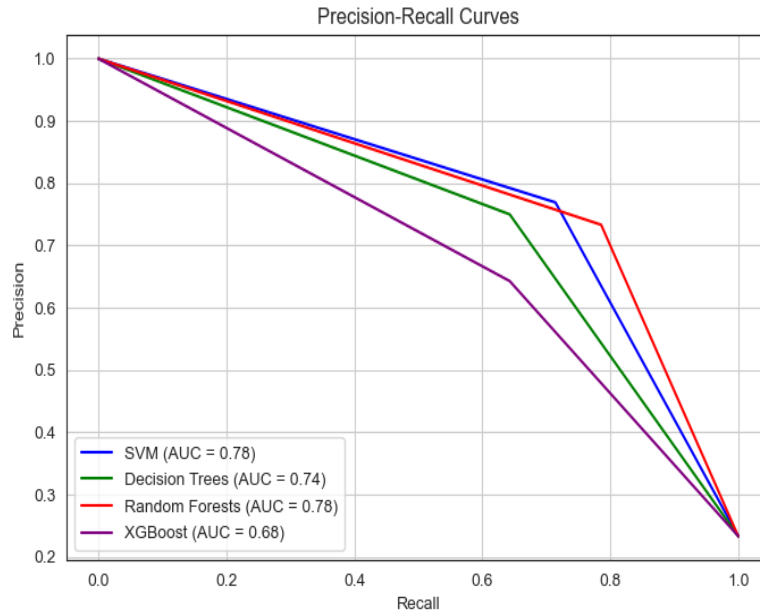


Figure 4.8: Model Performance Evaluation via Precision-Recall Curves

balance between precision and recall. This suggests that it has a reliable ability to confidently predict true positives and at the same time keeping the false positives in check. XGBoost, on the other hand, recorded an AUC of 0.68, indicating a somewhat lesser harmony between precision and recall compared to SVM. Decision Trees attained an AUC of 0.74, representing a reasonable performance, while Random Forest paralleled SVM's performance with an AUC of 0.78. This accentuates the strength of ensemble methods, particularly Random Forest, in achieving a fine balance in the precision and recall metrics.

Given the critical nature of our task - predicting heart failure - it's essential to ensure the selected model minimizes both false positives (incorrectly identifying a healthy person as at risk) and false negatives (missing a potential heart failure risk). Between the models, both SVM and Random Forest appear to offer a marginally better balance, making them strong candidates for further tuning and deployment in real-world applications.

#### 4.4 Model Selection: Deciding the Best Fit

After meticulously analyzing the results of our models from various perspectives, including the confusion matrices, ROC curves, and Precision-Recall curves, it's pivotal to collate our insights and draw a conclusive perspective on the best-performing model. This section aims to holistically compare the performances and discern the most suitable model for predicting heart failure.

From the confusion matrices, we discerned that each model presented varying degrees of

true positives and true negatives, with Random Forests showing a particularly strong balance. The ROC curves, which evaluate the equilibrium between sensitivity (true positive rate) and specificity (false positive rate), revealed Random Forest as the top performer with an AUC of 0.95, closely followed by XGBoost and SVM. Precision-Recall curve analysis, crucial for imbalanced datasets like ours, also showcased strong performances from SVM and Random Forest, both registering an AUC of 0.78.

Taking all metrics into account, while SVM, XGBoost, and Decision Trees demonstrated commendable performances, the Random Forests model consistently emerged as the most robust and balanced model across the board. Its ensemble nature, which aggregates results from multiple decision trees, undoubtedly contributes to its strong, consistent performance. Given its superior predictive accuracy and balance in both the ROC and Precision-Recall evaluations, we recommend the **Random Forest** model as the prime candidate for further optimization and real-world deployment in the realm of heart failure prediction.

This conclusion, however, does not downplay the effectiveness of the other models. Each algorithm has its own merits and could be the model of choice under different circumstances or datasets. But, for our current dataset and objectives, Random Forest take the lead.

**Path Forward:** Even though Random Forests stand out as the preferred choice, it's essential to remember that every model can be further improved. Model optimization, especially through technique like grid search, can fine-tune hyperparameters, potentially enhancing predictive accuracy. In our upcoming sections, we'll dive deep into the world of model optimization, aiming to squeeze out even more performance from our chosen model. Stay tuned!

## 4.5 Model Optimization:

Model optimization is the next logical step after establishing a baseline model. While our models have showcased promising results, there's always room for improvement. This process ensures we maximize the efficiency, accuracy, and generalization capability of our chosen algorithm by fine-tuning its hyperparameters.

As we strive to further enhance our models' performances, delving into hyperparameter tuning emerges as the pivotal step. In this section, we will be applying grid search method across all our models: SVM, XGBoost, Decision Trees, and Random Forest. Through this comprehensive analysis, we aim to extract the utmost predictive power from each model and then make an informed decision about which optimized model stands out as the best. Performance of a model is greatly influenced by hyperparameters. The training process is governed by these pre-established setups. Hyperparameters are predetermined before the learning phase, in contrast to model parameters, which are learned from data. The accuracy of a model's predictions can be considerably improved by their proper configuration, see (Jeremy Jordan, 2017).

### 4.5.1 Hyperparameter Tuning Technique:

The fine-tuning of hyperparameters is an essential step in machine learning model optimization. The right choice of hyperparameters can significantly boost the model's performance.

**Grid Search:** Grid Search stands out as a widely adopted technique for hyperparameter tuning, particularly for its systematic and comprehensive approach. By exhaustively evaluating every possible combination of hyperparameters specified in a predefined grid, it ensures a thorough exploration of the parameter space.

In the grid search methodology, one specifies a set of possible values for each hyperparameter of interest. The method then constructs a grid of all combinations of these hyperparameter values and evaluates the model's performance for each combination. For instance, if you're optimizing two hyperparameters, one with 5 possible values and the other with 4 possible values, grid search will evaluate the model with all 20 possible combinations, see (Jason Brownlee, 2020).

The primary advantage of grid search is its precision. By considering every potential set of hyperparameters, it assures finding the combination that offers the best model performance, as measured by a predefined metric, such as accuracy or F1-score.

However, this precision comes at a cost. The computational demands can be significant, especially if the grid is large or if the model itself is computationally intensive. Hence, while grid search is ideal for scenarios where optimal performance is a priority, and computational resources are available, one must consider its compromises in terms of time and computational expense. The primary aim of this section is to optimize the hyperparameters to enhance its predictive performance, let us see for each model:

### 4.5.2 SVM:

- **Parameters Explored:**

1. **C (Regularization Parameter):** The regularization parameter, C, determines the trade-off between maximizing the margin of the hyperplane and minimizing the classification error. It essentially controls the balance between underfitting and overfitting.

Values explored: [0.1, 1, 10]

2. **gamma (Kernel Coefficient):** In kernel methods such as the Radial Basis Function (RBF) (Jason Brownlee, 2020) in SVM, the parameter gamma plays a pivotal role in determining the extent to which individual training samples influence the model. A low gamma value can restrict the model's capacity, possibly resulting in underfitting, whereas a high gamma can render the model excessively flexible, potentially leading to overfitting.

Values explored: ['scale', 'auto', 0.1, 1]



- **Optimal Parameters Found:** After conducting an extensive grid search, the model determined the optimal hyperparameters for the SVM as:

**C: 1**

**gamma: 'scale'**

The 'scale' setting for gamma is calculated as  $\frac{1}{n \cdot \text{features} \cdot X.\text{var}()}$  for the input data  $X$ .

- **Resulting Performance:** With these optimal hyperparameters in place, the SVM model achieved an accuracy of 0.90 on the test dataset. This high accuracy 90% underscores the SVM's robustness and enhanced predictive capacity post optimization.

### 4.5.3 Decision Trees:

- **Parameters Explored:**

1. **max\_depth:** This parameter dictates the maximum depth of the tree. A shallow tree might result in underfitting, while a tree that's too deep could overfit the data.

Values Explored: [20, 25, 27]

2. **min\_samples\_split:** Specifies the minimum number of samples required to split an internal node. A smaller value might lead to a deeper tree, potentially increasing the risk of overfitting.

Values Explored: [2, 3, 4, 5, 6, 7]

3. **min\_samples\_leaf:** Determines the minimum number of samples required to be at a leaf node. Setting this parameter prevents the tree from growing too deep by ensuring a minimum number of samples in the leaves.

Values Explored: [1, 2, 10, 20, 30, 40, 50, 60, 100]

- **Optimal Parameters Found:** After an exhaustive grid search, the optimal hyperparameters for the Decision Trees classifier were discerned as:

**max\_depth: 20**

**min\_samples\_split: 2**

**min\_samples\_leaf: 30**

- **Resulting Performance:** Upon utilizing these optimal hyperparameters, the Decision Trees classifier exhibited an accuracy of 0.87 on the test dataset. This accuracy (87%) indicates the effectiveness of the Decision Trees classifier, particularly after the hyperparameter optimization process.

### 4.5.4 Random Forest:

- **Parameters Explored:**

1. **n\_estimators:** Represents the number of trees in the forest. Increasing the number of trees can lead to a more robust model but may also increase computational cost.  
Values Explored: [100, 200, 300]
  2. **max\_depth:** This parameter sets the maximum depth of the tree. Adjusting this parameter helps to prevent overfitting, with a balanced value offering optimal generalization.  
Values Explored: [10, 20, 30]
  3. **min\_samples\_split:** Specifies the minimum number of samples required to split an internal node. Modifying this parameter can help in controlling the tree depth.  
Values Explored: [2, 5, 6, 7]
  4. **min\_samples\_leaf:** Determines the minimum number of samples required for a leaf node. It ensures that the tree remains robust and prevents it from being overly deep.  
Values Explored: [2, 3, 4, 5, 6]
- **Optimal Parameters Found:** Post rigorous grid search, the best hyperparameters for the Random Forest classifier were identified as:  
**n\_estimators: 300**  
**max\_depth: 10**  
**min\_samples\_split: 6**  
**min\_samples\_leaf: 5**
  - **Resulting Performance:** With the above-identified optimal hyperparameters, the Random Forest classifier achieved an accuracy of 0.92 on the test data. This elevated accuracy (92%) underscores the proficiency of the Random Forest classifier, especially after the hyperparameter optimization process.

#### 4.5.5 XGBoost classifier:

- **Parameters Explored:**
  1. **max\_depth:** This parameter represents the maximum depth of a tree, influencing over-fitting as deeper trees can capture more specific patterns.  
Values Explored: [2, 3, 4, 6, 9]
  2. **learning\_rate:** Step size at each iteration while moving towards a minimum of the loss function. A lower value makes the optimization more robust.  
Values Explored: [0.01, 0.05, 0.1, 0.3]
  3. **n\_estimators:** Number of boosting rounds or trees to be run. It's important to tune it properly to avoid underfitting or overfitting.  
Values Explored: [100, 150, 200, 250, 270]

- 4. **min\_child\_weight:** Used to control over-fitting. Higher values make the algorithm more conservative.

Values Explored: [1, 3, 4, 5, 6, 7]

- **Optimal Parameters Found:** After an exhaustive grid search, the optimal hyperparameters for the XGBoost classifier were found to be:

**max\_depth: 2**

**learning\_rate: 0.05**

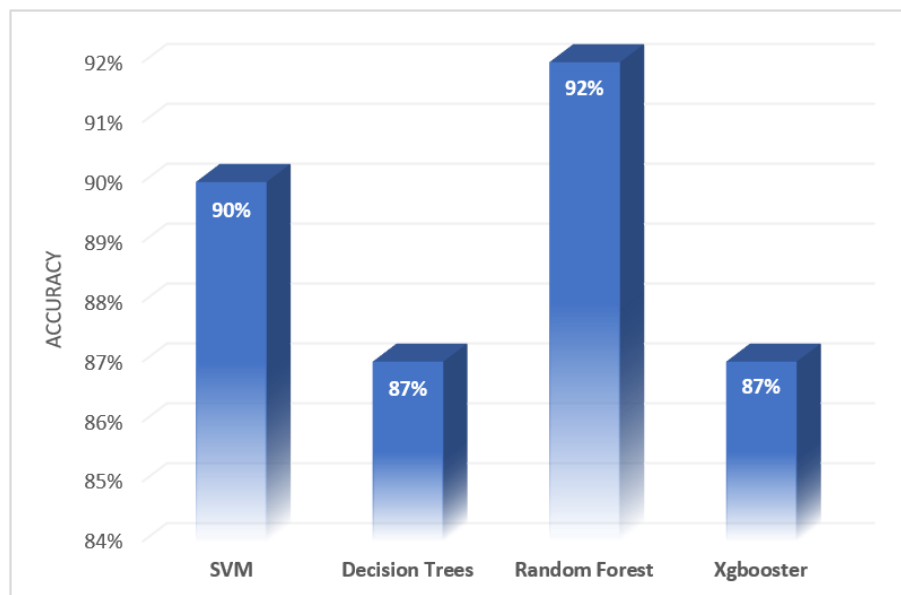
**n\_estimators: 100**

**min\_child\_weight: 7**

- **Resulting Performance:** Employing these optimal hyperparameters, the XGBoost classifier achieved an accuracy of 0.87 on the test set. This accuracy (87%) showcases the strength of the XGBoost model when appropriately tuned.

The comprehensive process of model development, evaluation, and hyperparameter tuning for the four machine learning algorithms has rendered a vivid picture of their predictive capabilities. When the accuracies are plotted on a bar graph, refer Figure 4.9, it's evident that while each model offers commendable performances, some minor differences become highlighted.

The **Random Forest** model takes the lead with an impressive accuracy of **92%**, followed closely by the SVM model at 90%. Both the Decision Trees and XGBoost models align with an accuracy of 87%. These results underscore the significance of employing diverse algorithms to uncover the most suitable model for a given dataset and the task at hand.



*Figure 4.9: Comparative Accuracy of Optimized Machine Learning Models*

In essence, the optimization phase has been pivotal in refining each model's performance, emphasizing the importance of hyperparameter tuning in machine learning workflows. As we

transition to the concluding chapter of this dissertation, we'll encapsulate the entirety of our journey, shedding light on the key findings and hinting at potential future directions in this realm of study.

## Chapter 5

# Conclusion

### 5.1 Recapitulation

At the beginning of this research, the primary objective was to employ machine learning models to make accurate predictions on the given dataset, and, as the journey progressed, it delved deeper into the details of model evaluation, optimization, and interpretation. In the earlier chapters, a comprehensive exploration of the dataset was undertaken, setting the stage for the subsequent modeling process.

Different machine learning models, namely Support Vector Machine (SVM), Decision Trees, Random Forests, and XGBoost, were introduced, and their assumptions were explained. The importance of understanding the underlying mechanics and assumptions of each model was emphasized, as it constructed the way for a clearer interpretation of the results.

The performance of each model was then assessed through a series of metrics, such as confusion matrices, ROC curves, and Precision-Recall curves. These evaluations provided critical insights into the strengths and weaknesses of each model in predicting the desired outcomes. To further enhance the models' performance, hyperparameter tuning was executed, focusing mainly on Grid Search methodology. This optimization step aimed to fine-tune the models and extract their maximum potential in terms of accuracy. By comparing the results from the models before and after optimization, a clearer understanding of the improvements achieved was obtained. This rigorous process was instrumental in guiding the research towards identifying the best-performing model for the task at hand. The journey through this research has been a meticulous endeavor, balancing between model understanding, evaluation, and optimization, all while keeping the end goal of accurate prediction in sight.

### 5.2 Major Takeaways

1. **Insights from Data Analysis:** As the narrative unfolds, life's journey takes a distinct turn beyond the age of 40. During this phase, the tendency to heart-related issues becomes increasingly apparent, emphasizing the pivotal role of age in this narrative. Then, the heart's

ejection fraction comes into focus. The story gets clearer when the ejection fraction goes below 40%, suggesting its key role in pointing towards heart-related complications.

As we move forward, time plays a central role. There's a specific window, from 25 to 100 days, where the risk becomes more noticeable, highlighting the need to closely watch over patients during this time. Lastly, the tale introduces two critical elements: serum creatinine and sodium levels. Their levels tell a clear story. Serum creatinine levels above 1.2 and serum sodium figures between 127-145 show increased chances of heart problems. This highlights the importance of regular health check-ups and monitoring, ensuring a better outcome for patients.

## 2. Model Performances:

- **Pre-optimization**, all models showcased noteworthy predictive capabilities, with the Random Forest and SVM models standing out due to their superior initial accuracies.
- **Post-optimization**, through hyperparameter tuning via Grid Search, each model experienced enhancements in performance. The Random Forest model, in particular, reached an impressive accuracy, solidifying its position as a reliable tool for this dataset.

## 3. Significance of Optimization:

- Hyperparameter tuning underscored the importance of understanding the parameters governing each model. By fine-tuning these, the study was able to push the boundaries of accuracy, revealing the latent potential within each model.
- The use of grid search emerged as a game-changer, systematically scanning through parameter combinations to identify the best set for each model. The effort invested in this step was rewarded with notable accuracy boosts.

## 5.3 Limitations and Challenges:

Throughout our research, we faced certain challenges and limitations, highlighting areas that could benefit from further attention in future studies.

- **Dataset limitations:** The dataset used for our models had certain limitations due to its finite size, which could have been more extensive to provide more detailed insights, particularly beneficial for ensemble methods.
- **Feature selection:** While we selected features based on domain knowledge, there's a possibility that we missed important ones, potentially affecting the accuracy of our models.

- **Hyperparameter tuning challenges:** The grid search method, while comprehensive, demanded significant computational resources. With numerous hyperparameters to consider, especially for complex models like Random Forest and XGBoost, we may not have explored all the optimal parameter combinations.
- **Overfitting concern:** Despite our efforts to prevent overfitting, it remains a concern in modeling exercises, especially with intricate models that can closely fit to training data, capturing random fluctuations rather than the actual pattern.

Acknowledging these challenges doesn't take away from our research's value. Instead, it provides a more rounded understanding of our findings and sets the stage for potential improvements in future studies.

## 5.4 Future Work

1. **Expanded Dataset:** To enhance the reliability and accuracy of our models, gathering a larger, more diverse dataset would be invaluable. Not only would this provide a richer data environment, but it could also potentially unveil patterns and trends that our current dataset might have missed.
2. **Feature Engineering and Selection:** While our current set of features has served us well, further work can be undertaken in the realm of feature engineering. New features could be derived or existing ones transformed to provide deeper insights. Utilizing techniques like Principal Component Analysis (PCA) or feature importance from ensemble methods could assist in more optimal feature selection.
3. **Alternative Algorithms:** There are several state-of-the-art machine learning algorithms and techniques emerging. Neural networks, especially deep learning models like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), could be explored in future studies, given their expertise in handling vast amounts of data and complex relationships.
4. **Emerging Technologies:** With the rapid advancements in the field of AI and machine learning, it's crucial to stay knowledgeable of new developments. Techniques like transfer learning, where pre-trained models are adapted for new tasks, or the use of quantum computing in machine learning, might soon become relevant for research in our domain.

## 5.5 Acknowledgments:

First and foremost, I would like to extend my heartfelt gratitude to Professor Tim Heaton. His consistent guidance, invaluable feedback, and unwavering support have been instrumental in the

completion of this research. His deep insights and seasoned expertise played a pivotal role in shaping this project, and for that, I am profoundly thankful.

I would also like to express my sincere appreciation to my friends, whose encouragement and companionship throughout this journey have been a source of strength. Their timely advice, genuine critiques, and countless brainstorming sessions were crucial in refining my ideas and methodologies.

Lastly, I would like to thank everyone who indirectly contributed to this research. Whether it was through providing datasets, assisting in software troubleshooting, or offering moral support during challenging phases, every gesture was significant. In conclusion, research is seldom a solitary endeavor, and the contributions of many have enriched this work. I am genuinely grateful to each one for their role in making this project a reality.

## **Appendix**

You can find the code for this project on GitHub at the following link:<https://github.com/Shashidhar-Rolex/Heart-Failure-Dataset-Analysis.git>



# Bibliography

1. World Health Organization. Cardiovascular diseases (CVDs). Retrieved from [https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds)), 2019.
2. Cardiovascular disease. Retrieved from <https://www.nhs.uk/conditions/cardiovascular-disease/>
3. The Economist, in an article titled "The world's most valuable resource" from 6 May 2017. Retrieved from <https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data>
4. Piwek, L., Ellis, D. A., Andrews, S., & Joinson, A. (2016). The rise of consumer health wearables: promises and barriers. *PLOS Medicine*, 13(2), e1001953.
5. Ahmad T, Munir A, Bhatti SH, Aftab M, Raza MA. Survival analysis of heart failure patients: a case study. *PLoS ONE*. 2017;12(7):0181001.
6. Jain, A., Nandakumar, K., & Ross, A. (2005). Score normalization in multimodal biometric systems. *Pattern recognition*, 38(12), 2270-2285.
7. Correlation and regression. Retrieved from <https://www.bmj.com/about-bmj/resources-readers/publications/statistics-square-one/11-correlation-and-regression>
8. National Institute of Diabetes and Digestive and Kidney. Retrieved from <https://www.niddk.nih.gov/health-information/diabetes/overview/preventing-problems/heart-disease-stroke>
9. Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297.
10. Correlation and regression. Retrieved from <https://shuzhanfan.github.io/2018/05/understanding-mathematics-behind-support-vector-machines/>
11. (Part – I). Retrieved from <https://www.datascienceprophet.com/understanding-the-mathematics-behind-the-decision-tree-algorithm-part-i/>

12. Towards Data Science. Retrieved from <https://t.ly/-6Xpb>
13. Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 785–794). San Francisco, CA, USA: ACM.
14. Towards Data Science. Retrieved from <https://towardsdatascience.com/understanding-confusion-matrix-a9ad42dcfd62>
15. IBM SPSS Statistics. Retrieved from <https://www.ibm.com/docs/en/spss-statistics/beta?topic=analysis-roc-statistics>
16. Towards Data Science. Retrieved from <https://towardsdatascience.com/a-look-at-precision-recall-and-f1-score-36b5fd0dd3ec>
17. Jeremy Jordan, Data Science. Retrieved from <https://www.jeremyjordan.me/hyperparameter-tuning/>
18. Machine Learning Mastery, Data Science. Retrieved from <https://t.ly/ERII6/>
19. Cardiovascular disease. Retrieved from <https://www.nhs.uk/conditions/cardiovascular-disease/>