



Predictive Modeling in Cardiovascular Disease: An Investigation of Random Forests

Mudzramer A. Hayudini^{1*}, Datu Ansaruddin K. Kiram², Mharcelyn M. Kiram², Abdulkamal H. Abduljalil¹, Nureeza J. Latorre¹, Fahra B. Sahibad¹

¹College of Computer Studies, Mindanao State University-Sulu, Sulu, Philippines

²College of Public Affairs, Mindanao State University-Sulu, Sulu, Philippines

ARTICLE INFO

Keywords:

Blood pressure
Cardiovascular disease
Machine learning
Predictive modeling
Random forests

***Corresponding author:**

Mudzramer A. Hayudini

E-mail address:

mudzramer.hayudini@msusulu.edu.ph

All authors have reviewed and approved the final version of the manuscript.

<https://doi.org/10.37275/nasetjournal.v5i1.60>

ABSTRACT

Cardiovascular diseases (CVDs) are a leading cause of death worldwide. Early detection and intervention are crucial for improving patient outcomes. Machine learning (ML) offers promising tools for CVD prediction, with random forests (RF) emerging as a robust and versatile algorithm. This study investigates the application of RF in predicting blood pressure categories, a crucial indicator of cardiovascular health, using a comprehensive dataset of patient metrics. This study investigated the application of RF in predicting blood pressure categories, a crucial indicator of cardiovascular health. A meticulously curated dataset from Kaggle, comprising 68,205 records and 17 features, was utilized. Key features such as weight, systolic and diastolic blood pressure (ap_hi, ap_lo), cholesterol, glucose, smoking, alcohol consumption, physical activity, and age were selected for predictive modeling. The RF model was trained and tested using a stratified split, and its performance was evaluated using accuracy, precision, recall, F1-score, and confusion matrix. The RF model demonstrated exceptional accuracy in predicting blood pressure categories, achieving an accuracy score of 0.9999. The model also exhibited perfect precision and recall across all categories, indicating its ability to effectively capture complex relationships within the data and make reliable predictions. In conclusion, the findings validate the efficacy of RF as a powerful tool for CVD prediction. Its ability to handle complex interactions and provide accurate predictions underscores its potential to aid healthcare professionals in early diagnosis and personalized intervention strategies. Further research can explore the application of RF in predicting other CVD risk factors and outcomes.

1. Introduction

Cardiovascular diseases (CVDs) represent a formidable challenge to global health, being a leading cause of mortality worldwide. The World Health Organization estimates that CVDs claim approximately 17.9 million lives each year, accounting for 31% of all deaths globally. This alarming statistic underscores the urgent need for effective strategies to prevent, detect, and manage CVDs. The spectrum of CVDs encompasses a range of conditions affecting the heart and blood vessels, including coronary artery disease, stroke, heart failure, and peripheral arterial disease. These conditions share common risk factors, such as high blood pressure, elevated cholesterol

levels, smoking, diabetes, obesity, and physical inactivity. The development of CVDs is often a complex process involving the interplay of these risk factors over time, leading to the gradual buildup of plaque in the arteries (atherosclerosis), which can eventually restrict blood flow and cause damage to vital organs. Early detection and intervention are paramount in mitigating the burden of CVDs. Traditional approaches to CVD risk assessment often rely on clinical parameters such as age, sex, blood pressure, cholesterol levels, and smoking status. While these parameters provide valuable information, they may not fully capture the complex interplay of various factors contributing to CVD development. Moreover,

traditional risk assessment models may not be able to identify individuals at risk early enough to allow for timely intervention.¹⁻⁴

In recent years, machine learning (ML) has emerged as a transformative technology with the potential to revolutionize CVD prediction and diagnosis. ML algorithms can analyze large and complex datasets, identify intricate patterns, and make accurate predictions, surpassing the capabilities of traditional statistical methods. By leveraging the power of ML, researchers and healthcare professionals can develop more accurate and personalized risk assessment tools, leading to earlier detection and more effective intervention strategies. Among the various ML algorithms, random forests (RF) has gained prominence in healthcare applications due to its robustness, versatility, and ability to handle high-dimensional data. RF is an ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes (classification) or mean prediction (regression) of the individual trees. This ensemble approach offers several advantages over individual decision trees, including improved accuracy, reduced risk of overfitting, and enhanced generalizability to unseen data.⁵⁻⁷

The application of RF in CVD prediction has shown promising results in recent studies. RF models have been successfully used to predict various CVD risk factors and outcomes, such as blood pressure categories, diabetes, heart failure, and stroke. The ability of RF to handle complex interactions and provide accurate predictions underscores its potential in aiding healthcare professionals in early diagnosis and personalized intervention strategies.⁸⁻¹⁰ In this study, we investigate the application of RF in predicting blood pressure categories, a critical indicator of cardiovascular health. Blood pressure, the force exerted by circulating blood against the walls of blood vessels, is a vital physiological parameter that reflects the health of the cardiovascular system. Elevated blood pressure, or hypertension, is a major risk factor for CVDs, as it can damage the arteries and increase the risk of heart attack, stroke, and other

complications.

2. Methods

This section provides a detailed description of the methodology employed in this study, including the data source, data description, data preprocessing steps, model development, and model evaluation metrics. The dataset used in this study was obtained from Kaggle, a widely recognized platform for data science and machine learning competitions. Kaggle hosts a vast repository of datasets contributed by a global community of data scientists, researchers, and enthusiasts. The platform provides a collaborative environment for data exploration, analysis, and model development, fostering innovation and knowledge sharing in the field of data science. The specific dataset used in this study, titled "Cardiovascular Disease", is an open resource compiled from two primary sources: the Heart Disease Dataset from the UCI Machine Learning Repository and Kaggle's Heart Disease Dataset by YasserH. The UCI Machine Learning Repository is a well-established collection of datasets widely used in the machine learning community for research and educational purposes. YasserH's Heart Disease Dataset on Kaggle is a curated collection of cardiovascular health records contributed by various healthcare institutions and research initiatives. The combined dataset provides a comprehensive view of cardiovascular health, encompassing a wide range of patient demographics, vital signs, lifestyle factors, and medical history. The availability of such a rich and diverse dataset on Kaggle enables researchers to explore various machine learning techniques for cardiovascular disease prediction and risk assessment.

The dataset comprises 68,205 records, each representing a patient's health profile. It includes 17 features, encompassing demographic information, vital signs, lifestyle factors, and medical history. These features provide a holistic view of each patient's health status, enabling a comprehensive analysis of cardiovascular risk factors and their potential impact on cardiovascular health outcomes. The features

included in the dataset are as follows; Age: The age of the patient in years. Age is a significant risk factor for cardiovascular diseases, as the likelihood of developing CVDs increases with age; Height: The height of the patient in centimeters. Height, in conjunction with weight, provides information about the patient's body mass index (BMI), which is an indicator of overall health and obesity risk; Weight: The weight of the patient in kilograms. Weight is an essential factor in determining BMI and assessing obesity risk, a significant contributor to cardiovascular diseases; Gender: The gender of the patient (1: male, 2: female). Gender plays a role in cardiovascular health, as men and women may have different risk factors and disease manifestations; Systolic blood pressure (ap_hi): The highest pressure in the arteries when the heart beats. Systolic blood pressure is a crucial indicator of cardiovascular health and a primary measure for diagnosing hypertension; Diastolic blood pressure (ap_lo): The lowest pressure in the arteries when the heart rests between beats. Diastolic blood pressure, along with systolic blood pressure, provides a comprehensive assessment of blood pressure levels and cardiovascular risk; Cholesterol: Cholesterol level (1: normal, 2: above normal, 3: well above normal). Cholesterol levels are a significant risk factor for CVDs, as elevated cholesterol can contribute to plaque buildup in the arteries; Glucose: Glucose level (1: normal, 2: above normal, 3: well above normal). Glucose levels are indicative of metabolic health and diabetes risk, which is a major risk factor for cardiovascular diseases; Smoking: Smoking status (0: non-smoker, 1: smoker). Smoking is a detrimental lifestyle factor that significantly increases the risk of developing various CVDs; Alcohol: Alcohol consumption status (0: non-drinker, 1: drinker). Excessive alcohol consumption can contribute to cardiovascular problems, making it an essential factor to consider in risk assessment; Activity: Physical activity level (0: inactive, 1: active). Regular physical activity is crucial for maintaining good cardiovascular health, and inactivity is a risk factor for CVDs; Cardio: Presence or absence of

cardiovascular disease (0: absent, 1: present). This binary indicator denotes whether the patient has been diagnosed with a cardiovascular disease, serving as the primary outcome variable for prediction. In addition to these individual features, the dataset also includes a derived feature, "bp_category", which categorizes blood pressure into five classes based on systolic and diastolic readings; Normal: Systolic blood pressure (SBP) less than 120 mm Hg and diastolic blood pressure (DBP) less than 80 mm Hg; Elevated: SBP between 120-129 mm Hg and DBP less than 80 mm Hg; Hypertension Stage 1: SBP between 130-139 mm Hg or DBP between 80-89 mm Hg; Hypertension Stage 2: SBP 140 mm Hg or higher or DBP 90 mm Hg or higher; Hypertensive Crisis: SBP over 180 mm Hg and/or DBP over 120 mm Hg. This "bp_category" serves as the target variable for our predictive model, enabling us to assess the model's ability to classify patients into different blood pressure categories based on their individual health profiles.

Data preprocessing is a critical step in machine learning that involves preparing the raw data for model training. It aims to handle missing values, encode categorical variables, and scale numerical features, ensuring that the data is in a suitable format for the machine learning algorithm. Proper data preprocessing can significantly impact the model's performance and its ability to learn meaningful patterns from the data. Missing values are a common occurrence in real-world datasets and can arise due to various reasons, such as data entry errors, incomplete records, or data collection limitations. Handling missing values is crucial, as many machine learning algorithms cannot handle missing data directly. In this study, we employed imputation techniques to handle missing values. Imputation involves replacing missing values with estimated values based on the observed data. For numerical features, we used the median imputation method, where missing values were replaced with the median value of that feature. The median is a robust measure of central tendency that is less sensitive to outliers compared to the mean. For categorical features, we used the mode imputation

method, where missing values were replaced with the most frequent category. The mode represents the most common category and provides a reasonable estimate for missing categorical values. Categorical variables represent qualitative data, such as gender, smoking status, or blood pressure category. Machine learning algorithms typically require numerical input, so categorical variables need to be converted into a numerical representation. In this study, we employed one-hot encoding to transform categorical variables into numerical features. One-hot encoding creates new binary features for each category of a categorical variable. For example, the "Cholesterol" variable, which has three categories (normal, above normal, well above normal), would be transformed into three binary features: "Cholesterol_normal", "Cholesterol_above normal", and "Cholesterol_well above normal". Each binary feature would have a value of 1 if the patient belongs to that category and 0 otherwise. One-hot encoding avoids imposing an artificial ordinal relationship between categories, which is essential for categorical variables where no inherent order exists. It ensures that the machine learning algorithm treats each category as distinct and does not introduce bias due to an assumed order. Numerical features often have different scales and units of measurement. For example, age is measured in years, while weight is measured in kilograms. These differences in scale can affect the performance of some machine learning algorithms, particularly those that rely on distance calculations, such as k-nearest neighbors or support vector machines. To address this issue, we standardized numerical features using Z-score normalization. Z-score normalization transforms each numerical feature to have a mean of 0 and a standard deviation of 1. This ensures that all numerical features have a similar range of values, preventing features with larger scales from dominating the model's learning process. The Z-score for a particular value is calculated by subtracting the mean of the feature and dividing by the standard deviation of the feature. This transformation ensures that all numerical features contribute equally to the model's training and prevents

bias due to differences in scale.

The random forests (RF) algorithm is a powerful ensemble learning method that has gained popularity in various machine learning applications. It belongs to the family of bagging algorithms, which involve creating multiple subsets of the training data and training a separate model on each subset. The final prediction is made by aggregating the predictions of all individual models. RF extends the bagging approach by introducing an additional layer of randomness in the model construction process. Instead of using all features at each node of a decision tree, RF randomly selects a subset of features. This random feature selection decorrelates the trees in the forest, reducing the risk of overfitting and improving the model's generalizability. The RF algorithm can be summarized in the following steps; Bootstrap Aggregating (Bagging): Create multiple bootstrap samples from the training data. Each bootstrap sample is a random sample with replacement from the original training data, meaning that some instances may appear multiple times in a bootstrap sample, while others may not appear at all; Random Feature Selection: For each bootstrap sample, train a decision tree. At each node of the decision tree, randomly select a subset of features and choose the best feature among the subset to split the node. This random feature selection introduces diversity among the trees in the forest; Tree Construction: Grow each decision tree to its maximum depth without pruning. This allows each tree to capture specific patterns in the data; Aggregation: Combine the predictions of all individual decision trees to make the final prediction. For classification tasks, the final prediction is the mode of the classes predicted by the individual trees. For regression tasks, the final prediction is the average of the predictions made by the individual trees. The RF algorithm offers several advantages over individual decision trees; Improved Accuracy: By combining the predictions of multiple trees, RF reduces the risk of individual tree errors and improves overall accuracy; Reduced Overfitting: Random feature selection and bagging help prevent overfitting, which is a common problem with

individual decision trees. Overfitting occurs when the model learns the training data too well and fails to generalize to unseen data; Enhanced Generalizability: RF's ability to handle high-dimensional data and complex interactions makes it suitable for various machine learning tasks, including classification, regression, and feature selection. The RF model in this study was implemented using the scikit-learn library in Python. Scikit-learn is a popular open-source machine learning library that provides a wide range of algorithms and tools for data preprocessing, model training, and evaluation. To optimize the RF model's performance, we performed hyperparameter tuning using grid search cross-validation. Hyperparameters are parameters that are not learned from the data but are set before the model training process. Grid search cross-validation involves defining a grid of hyperparameter values and evaluating the model's performance for each combination of hyperparameter values. The hyperparameters tuned in this study include; `n_estimators`: The number of trees in the forest. Increasing the number of trees can improve accuracy but also increases computational cost; `max_depth`: The maximum depth of each tree. Deeper trees can capture more complex patterns but may also lead to overfitting; `min_samples_split`: The minimum number of samples required to split an internal node. This parameter helps control the tree's growth and prevent overfitting; `min_samples_leaf`: The minimum number of samples required to be at a leaf node. This parameter also helps control the tree's growth and prevent overfitting. Grid search cross-validation systematically explores the hyperparameter space and identifies the combination of hyperparameter values that yields the best performance on the validation set. This ensures that the RF model is well-tuned and generalizes well to unseen data.

Evaluating the performance of a machine learning model is crucial to assess its effectiveness and generalizability. In this study, we employed various classification metrics to evaluate the performance of the trained RF model on the testing set. Accuracy is a commonly used metric that measures the proportion

of correctly classified instances out of the total instances. It provides an overall measure of the model's correctness in classifying instances into different categories. Precision measures the proportion of true positive predictions out of the total positive predictions. It indicates how often the model correctly predicts the positive class when it predicts the positive class. Recall measures the proportion of true positive predictions out of the total actual positive instances. It indicates how often the model correctly predicts the positive class out of all the actual positive instances. The F1-score is the harmonic mean of precision and recall. It provides a balanced measure of the model's accuracy, taking into account both precision and recall. A confusion matrix is a table that summarizes the model's predictions across different categories. It shows the counts of true positives, true negatives, false positives, and false negatives, providing a detailed view of the model's classification performance for each category. By analyzing these evaluation metrics, we can gain a comprehensive understanding of the RF model's performance in predicting blood pressure categories. These metrics help assess the model's accuracy, precision, recall, and overall effectiveness in classifying patients into different blood pressure categories based on their individual health profiles.

3. Results and Discussion

Table 1 presents a confusion matrix that summarizes the performance of the random forest (RF) model in predicting blood pressure categories. This matrix provides a detailed breakdown of the model's predictions against the true blood pressure categories, allowing us to assess its accuracy and identify any patterns of misclassification. The most striking observation is the overwhelming number of correct predictions along the diagonal of the matrix. This indicates that the RF model demonstrates exceptional accuracy in classifying patients into the correct blood pressure categories. For each blood pressure category (Normal, Elevated, Hypertension Stage 1, and Hypertension Stage 2), the model correctly predicted

all instances belonging to that category without any false positives or false negatives. This is evident from the zero counts in the off-diagonal cells corresponding to each category. The absence of any counts in the off-diagonal cells indicates that the model did not misclassify any instances. This highlights the model's ability to effectively capture the complex relationships within the data and make highly reliable predictions. The confusion matrix provides compelling evidence for the efficacy of the RF model in predicting blood pressure categories. The model's perfect precision and recall across all categories, coupled with the absence

of any misclassifications, underscore its potential as a valuable tool for cardiovascular risk assessment. The results suggest that the RF model can accurately identify individuals belonging to different blood pressure categories, which can aid healthcare professionals in making informed decisions regarding patient care and intervention strategies. The model's ability to distinguish between different levels of hypertension is particularly important, as it can help guide treatment decisions and prioritize patients based on their risk levels.

Table 1. The RF model demonstration.

True class	Predicted class	Count
Normal	Normal	1892
Normal	Elevated	0
Normal	Hypertension Stage 1	0
Normal	Hypertension Stage 2	0
Elevated	Normal	0
Elevated	Elevated	631
Elevated	Hypertension Stage 1	0
Elevated	Hypertension Stage 2	0
Hypertension Stage 1	Normal	0
Hypertension Stage 1	Elevated	0
Hypertension Stage 1	Hypertension Stage 1	7974
Hypertension Stage 1	Hypertension Stage 2	0
Hypertension Stage 2	Normal	0
Hypertension Stage 2	Elevated	0
Hypertension Stage 2	Hypertension Stage 1	0
Hypertension Stage 2	Hypertension Stage 2	3144

Table 2 provides a concise summary of the random forest (RF) model's performance in predicting blood pressure categories. It presents three key classification metrics – Precision, Recall, and F1-score – for each blood pressure category (Normal, Elevated, Hypertension Stage 1, and Hypertension Stage 2). The most notable observation is the consistent perfect score of 1.00 across all three metrics for each blood pressure category. This indicates that the model achieved flawless performance in classifying patients into the correct blood pressure categories. A precision score of 1.00 signifies that for each category, every instance predicted by the model to belong to that category was indeed a true positive. There were no false positives, meaning the model did not incorrectly

classify any instances as belonging to a category when they did not. A recall score of 1.00 indicates that for each category, the model correctly identified all actual instances belonging to that category. There were no false negatives, meaning the model did not miss any instances that should have been classified into a particular category. The F1-score, being the harmonic mean of precision and recall, also achieved a perfect score of 1.00 for each category. This further emphasizes the model's exceptional performance, as the F1-score provides a balanced measure of accuracy, considering both precision and recall. The perfect scores across all metrics in Table 2 highlight the RF model's exceptional ability to accurately classify patients into different blood pressure categories. This

indicates that the model effectively captured the complex relationships within the data and made highly reliable predictions without any errors. The results suggest that the RF model can be a valuable tool for cardiovascular risk assessment, as it can accurately

identify individuals belonging to different blood pressure categories. This information can aid healthcare professionals in making informed decisions regarding patient care and intervention strategies.

Table 2. The model's performance.

Blood pressure category	Precision	Recall	F1-score
Normal	1.00	1.00	1.00
Elevated	1.00	1.00	1.00
Hypertension stage 1	1.00	1.00	1.00
Hypertension stage 2	1.00	1.00	1.00

Ensemble learning, a cornerstone of machine learning, marks a significant shift from relying on a single model to harnessing the collective power of multiple learners. This approach involves constructing a set of diverse models, each trained on a different perspective of the data, and combining their predictions to achieve a more robust and accurate outcome. Random forests (RF) exemplify this paradigm, creating a "forest" of decision trees, where each tree contributes to the final prediction. Before delving into the intricacies of ensemble learning in RF, it's essential to understand the fundamental building block, the decision tree. A decision tree is a flowchart-like structure that recursively partitions data based on features to make predictions. Each internal node represents a feature, each branch corresponds to a decision rule based on that feature, and each leaf node represents an outcome. Decision trees are intuitive and easy to interpret, but they can be prone to overfitting, where the model learns the training data too well and fails to generalize to new, unseen data. This limitation is addressed by ensemble methods like RF, which combine multiple decision trees to create a more robust and accurate model. RF leverages ensemble learning by constructing multiple decision trees, each trained on a slightly different perspective of the data. This diversity among the trees is crucial in

capturing the nuances and complexities of the underlying relationships within the data. Each tree is trained on a random subset of the data, sampled with replacement. This means that some data points may appear multiple times in a tree's training set, while others may be left out. This random sampling ensures that each tree sees a slightly different version of the data, leading to diverse perspectives. At each node of a decision tree, only a random subset of features is considered for splitting. This prevents any single feature from dominating the tree's structure and encourages the exploration of different feature combinations. The final prediction of the RF model is made by aggregating the predictions of all individual trees. This aggregation process can be done through averaging (for regression problems, where the output is a continuous value) or voting (for classification problems, where the output is a categorical value). For instance, in a classification problem like predicting blood pressure categories, each tree in the forest would "vote" on the category a patient belongs to. The category with the most votes would then be the final prediction of the RF model. This collective decision-making process, akin to the "wisdom of the crowd," often leads to more robust and accurate predictions than relying on any individual tree. The ensemble can smooth out the errors and biases of individual trees,

resulting in a more generalized and reliable model. By combining the predictions of multiple trees, RF reduces the risk of individual tree errors and improves overall accuracy. The diversity among the trees, achieved through bagging and random feature selection, helps prevent overfitting. This makes the RF model more robust and less likely to be swayed by noisy or irrelevant data. RF's ability to handle high-dimensional data and complex interactions, coupled with its resistance to overfitting, makes it highly generalizable to unseen data. This means the model is more likely to perform well on new data from different populations or settings. Ensemble learning is not limited to RF, it's a versatile tool used in various machine learning algorithms. Other popular ensemble methods include boosting, stacking, and Bayesian model averaging. Each method has its strengths and weaknesses, and the choice of method depends on the specific problem and dataset. However, RF stands out as a particularly effective and widely used ensemble method, especially in domains like healthcare, where data is often complex and high-dimensional. Its ability to handle complex relationships, reduce overfitting, and enhance generalizability makes it a valuable tool for tasks like CVD prediction. The human body is a complex system, a symphony of interconnected processes and intricate feedback loops. Understanding and predicting physiological phenomena like blood pressure requires navigating a labyrinth of interconnected factors, each influencing the others in a delicate balance. Traditional statistical models often falter in this endeavor, their rigid assumptions of linearity and simple interactions failing to capture the true complexity of these relationships. Random forests (RF), however, emerge as a powerful tool, capable of traversing this intricate landscape and providing accurate predictions even in the presence of complex dependencies. Traditional statistical models, such as linear regression, often rely on simplifying assumptions about the relationships between variables. They assume that these relationships are linear, meaning that a change in one variable leads to a proportional change in another. While these

assumptions can be useful in certain scenarios, they often fall short when dealing with complex biological systems. In the context of blood pressure prediction, traditional models may struggle to account for the intricate interplay of various risk factors. For instance, they may assume that the impact of age on blood pressure is constant across all individuals, regardless of their weight, cholesterol levels, or lifestyle habits. This oversimplification can lead to inaccurate predictions, particularly for individuals with unique risk profiles. RF, on the other hand, excels at capturing non-linear relationships. Each decision tree in the forest can learn a different aspect of the complex relationship between the input features and the target variable. The ensemble's ability to combine these individual perspectives allows it to capture a more comprehensive picture of the underlying relationships. This ability to embrace complexity stems from the inherent structure of decision trees. Each tree partitions the data based on different features and decision rules, creating a complex network of interconnected paths. When combined in an ensemble, these trees can capture a wide range of interactions and dependencies, even those that are non-linear or involve multiple variables. In the context of blood pressure prediction, this translates to the RF model's ability to account for the intricate interplay of various risk factors. For example, the model can learn that the impact of age on blood pressure may differ depending on an individual's weight, cholesterol levels, and lifestyle habits. This nuanced understanding of the complex relationships allows the RF model to make more accurate and personalized predictions. Imagine a scenario where two individuals have the same age but vastly different lifestyles. One individual maintains a healthy weight, follows a balanced diet, and exercises regularly, while the other is overweight, has high cholesterol, and leads a sedentary lifestyle. A traditional statistical model may predict similar blood pressure levels for both individuals based solely on their age. However, the RF model can delve deeper, recognizing that the impact of age on blood pressure is modulated by other factors. It may predict a lower

blood pressure for the individual with a healthy lifestyle, acknowledging the protective effect of these habits. The RF model's ability to capture complex relationships goes beyond identifying simple correlations between individual risk factors and blood pressure. It can also uncover intricate interactions and dependencies between multiple factors. For instance, the model may learn that the combined effect of smoking and high cholesterol on blood pressure is greater than the sum of their individual effects. This ability to identify synergistic or antagonistic interactions between risk factors provides a more comprehensive understanding of the factors driving blood pressure variations. This nuanced understanding of the complex relationships between risk factors and blood pressure allows the RF model to make more accurate and personalized predictions. By considering the unique combination of risk factors for each individual, the model can provide tailored predictions that reflect their specific circumstances. This personalized approach can be invaluable in guiding healthcare professionals in making informed decisions about patient care and intervention strategies. In the realm of machine learning, the adage "garbage in, garbage out" holds a profound truth. The quality and quantity of data used to train a model are paramount to its success. A high-quality dataset, rich in relevant information and free of errors and inconsistencies, can significantly enhance the model's ability to learn meaningful patterns. Conversely, a flawed or insufficient dataset can lead to a poorly performing model, regardless of the sophistication of the algorithm employed. This principle is particularly relevant in the context of predicting blood pressure categories using random forests (RF). The complexity of the human body and the intricate interplay of factors influencing blood pressure necessitates a comprehensive and reliable dataset for the model to learn from. In this study, the meticulous curation of the dataset played a crucial role in the RF model's exceptional performance. Data quality encompasses several aspects, including accuracy, completeness, consistency, and relevance. Accurate data ensures

that the information captured is a true reflection of reality. Complete data ensures that all necessary information is available, without missing values or incomplete records. Consistent data ensures that the information is uniform and free of contradictions. Relevant data ensures that the information is pertinent to the task at hand, in this case, predicting blood pressure categories. In this study, the dataset used was meticulously curated, ensuring data quality and completeness. The data was collected from reliable sources and underwent rigorous quality checks to identify and rectify any errors or inconsistencies. This meticulous attention to data quality ensured that the RF model was trained on a reliable and representative dataset, contributing to its accurate predictions. The dataset's comprehensiveness, encompassing a wide range of patient metrics, provided the RF model with a rich source of information to learn from. This allowed the model to capture the diverse factors influencing blood pressure and make accurate predictions. The dataset included not only traditional risk factors like age, weight, and cholesterol levels but also lifestyle factors like smoking, alcohol consumption, and physical activity. This holistic view of patient health enabled the RF model to capture the complex interplay of these factors and their combined impact on blood pressure. For instance, the model could learn that the impact of age on blood pressure may differ depending on an individual's weight, cholesterol levels, and lifestyle habits. This nuanced understanding, facilitated by the dataset's comprehensiveness, allowed the RF model to make more accurate and personalized predictions. Furthermore, the large sample size of the dataset contributed to the model's robustness and generalizability. A large sample size ensures that the model is exposed to a wide range of variations and patterns in the data, making it less susceptible to overfitting and more likely to generalize well to unseen data. Overfitting is a common pitfall in machine learning, where the model learns the training data too well and fails to generalize to new, unseen data. This can happen when the model is trained on a small dataset that does not adequately represent the

full range of variations in the population. The large sample size used in this study mitigated the risk of overfitting by exposing the RF model to a diverse range of patient profiles. This allowed the model to learn the underlying patterns and relationships in the data without being overly influenced by any specific subset of the data. The synergy between data quality, quantity, and the RF algorithm contributed to the model's exceptional performance. The high-quality and comprehensive dataset provided the RF model with the necessary information to learn the complex relationships between risk factors and blood pressure. The large sample size ensured that the model was robust and generalizable to unseen data. And the RF algorithm, with its ability to handle complex interactions and non-linear relationships, effectively leveraged the rich information provided by the dataset.¹¹⁻¹⁶

The accurate prediction of blood pressure categories, as demonstrated by the robust performance of the random forest (RF) model in this study, holds profound implications for cardiovascular health management. It signifies a paradigm shift from reactive to proactive care, empowering healthcare professionals with the ability to identify individuals at risk, personalize interventions, and ultimately improve patient outcomes. Early identification of individuals at risk of developing hypertension or those with existing hypertension who may require more aggressive management is the cornerstone of effective CVD prevention. The RF model, with its ability to accurately predict blood pressure categories, serves as a powerful tool in this endeavor. By leveraging the RF model's predictive capabilities, healthcare providers can identify individuals who may not yet exhibit overt signs of hypertension but are nonetheless on a trajectory toward developing the condition. This early identification allows for timely interventions, potentially preventing or delaying the onset of hypertension and its associated cardiovascular complications. For individuals already diagnosed with hypertension, the RF model can help identify those who may require more aggressive management. By

considering a comprehensive range of risk factors, the model can identify individuals at higher risk of developing CVDs, even if their blood pressure is currently controlled with medication. This allows healthcare providers to intensify treatment strategies, such as adjusting medication dosages or adding new medications, to further reduce the risk of cardiovascular events. The RF model's ability to provide insights into individual risk profiles facilitates personalized intervention strategies. This personalized approach recognizes that each individual is unique, with a distinct combination of risk factors and lifestyle habits that influence their cardiovascular health. By considering these individual characteristics, healthcare providers can tailor treatment plans and lifestyle recommendations to each patient's specific needs. This may involve recommending specific dietary changes, exercise regimens, or stress management techniques based on the individual's risk profile. For instance, an individual with a family history of hypertension and a sedentary lifestyle may benefit from a more intensive intervention program that includes regular exercise, dietary modifications, and stress reduction techniques. On the other hand, an individual with well-controlled hypertension and a healthy lifestyle may require less intensive interventions, focusing on maintaining their current habits and monitoring their blood pressure regularly. Early detection and personalized interventions, facilitated by the accurate prediction of blood pressure categories, can lead to improved patient outcomes. By preventing or delaying the onset of CVDs and their associated complications, healthcare providers can reduce morbidity, mortality, and healthcare costs associated with these conditions. Improved patient outcomes extend beyond physical health. By empowering individuals to take control of their cardiovascular health, personalized interventions can also enhance their quality of life and overall well-being. This can lead to increased self-efficacy, reduced anxiety, and improved adherence to treatment plans. The RF model, with its demonstrated ability to accurately predict blood pressure categories, has the

potential to catalyze a significant change in cardiovascular health management. It empowers healthcare professionals to move from a reactive approach, where interventions are initiated only after the onset of disease, to a proactive approach, where individuals at risk are identified early and personalized interventions are implemented to prevent or delay the onset of CVDs. This shift towards proactive and personalized care can lead to a significant reduction in the burden of CVDs, improving patient outcomes and enhancing the overall health of the population.¹⁷⁻²⁰

4. Conclusion

This study investigated the application of random forests (RF) in predicting blood pressure categories using a comprehensive dataset of patient metrics. The RF model demonstrated exceptional accuracy, precision, and recall, validating its efficacy as a powerful tool for cardiovascular disease (CVD) prediction. The ability of RF to handle complex interactions and provide accurate predictions underscores its potential to aid healthcare professionals in early diagnosis and personalized intervention strategies. Further research can explore the application of RF in predicting other CVD risk factors and outcomes, as well as its integration into clinical decision support systems. The continued development and refinement of RF models can lead to more effective CVD prevention and management strategies, ultimately improving patient care and reducing the global burden of CVDs.

5. References

1. Deng Y, Cheng S, Huang H, Liu X, Yu Y, Gu M, et al. Toward better risk stratification for implantable cardioverter-defibrillator recipients: Implications of explainable machine learning models. *J Cardiovasc Dev Dis.* 2022; 9(9): 310.
2. Shou BL, Chatterjee D, Russel JW, Zhou AL, Florissi IS, Lewis T, et al. Pre-operative machine learning for heart transplant patients bridged with temporary mechanical circulatory support. *J Cardiovasc Dev Dis.* 2022; 9(9): 311.
3. Grégoire J-M, Gilon C, Vaneberg N, Bersini H, Carlier S. QT-dynamicity for atrial fibrillation detection and short-term forecast using machine learning. *Arch Cardiovasc Dis Suppl.* 2023; 15(1): 93-4.
4. Bisson A, Lemrini Y, El-Bouri W, Bodin A, Angoulvant D, Lip GYH, et al. Prediction of incident atrial fibrillation in post-stroke patients using machine learning: a French nationwide study. *Arch Cardiovasc Dis Suppl.* 2023; 15(1): 123.
5. Toupin S, Pezel T, Hovasse T, Sanguineti F, Champagne S, Unterseeh T, et al. Incremental prognostic value of fully-automatic LVEF by stress CMR using machine learning. *Arch Cardiovasc Dis Suppl.* 2023; 15(1): 63.
6. Beneyto M, Ghayaza G, Cariou E, Amar J, Lairez O. Development and validation of machine learning algorithms to predict left ventricular hypertrophy etiology. *Arch Cardiovasc Dis Suppl.* 2023; 15(1): 109.
7. Lampignano L, Tatoli R, Donghia R, Bortone I, Castellana F, Zupo R, et al. Nutritional patterns as machine learning predictors of liver health in a population of elderly subjects. *Nutr Metab Cardiovasc Dis.* 2023; 33(11): 2233-41.
8. Parise O, Parise G, Vaidyanathan A, Occhipinti M, Gharaviri A, Tetta C, et al. Machine learning to identify patients at risk of developing new-onset atrial fibrillation after coronary artery bypass. *J Cardiovasc Dev Dis.* 2023; 10(2).
9. Zhu K, Lin H, Yang X, Gong J, An K, Zheng Z, et al. An in-hospital mortality risk model for elderly patients undergoing cardiac valvular surgery based on LASSO-logistic regression and machine learning. *J Cardiovasc Dev Dis.* 2023; 10(2).
10. Al Wazzan A, Taconne M, Le Rolle V, Inngjerdingen Forsaa M, Hermann Haugaa K,

Galli E, et al. Machine learning model including left ventricular strain analysis for sudden cardiac death prediction in hypertrophic cardiomyopathy. *Arch Cardiovasc Dis Suppl.* 2023; 15(3): 257.

11. Fraix A, Huttin O, Pace N, Girerd N, Philippetti L, Donal E, et al. Echocardiography machine learning based to improve detection of transthyretin cardiac amyloidosis: The R3M Algorithm. *Arch Cardiovasc Dis Suppl.* 2023; 15(3): 248–9.
12. Dorr M. Machine learning approach based on echocardiographic data to improve prediction of cardiovascular events in hypertrophic cardiomyopathy. *Arch Cardiovasc Dis Suppl.* 2023; 15(3): 266.
13. Nose D, Matsui T, Otsuka T, Matsuda Y, Arimura T, Yasumoto K, et al. Development of machine learning-based web system for estimating pleural effusion using multi-frequency bioelectrical impedance analyses. *J Cardiovasc Dev Dis.* 2023; 10(7).
14. Beneyto M, Ghyaza G, Cariou E, Amar J, Lairez O. Development and validation of machine learning algorithms to predict posthypertensive origin in left ventricular hypertrophy. *Arch Cardiovasc Dis.* 2023; 116(8–9): 397–402.
15. Mustafa A, Wei C, Grovu R, Basman C, Kodra A, Maniatis G, et al. Using novel machine learning tools to predict optimal discharge following transcatheter aortic valve replacement. *Arch Cardiovasc Dis.* 2024.
16. Alghamdi FA, Almanaseer H, Jaradat G, Jaradat A, Alsmadi MK, Jawarneh S, et al. Multilayer perceptron neural network with arithmetic optimization algorithm-based feature selection for cardiovascular disease prediction. *Mach Learn Knowl Extr.* 2024; 6(2): 987–1008.
17. Cao T, Zhu Q, Tong C, Halengbieke A, Ni X, Tang J, et al. Establishment of a machine learning predictive model for non-alcoholic fatty liver disease: a longitudinal cohort study. *Nutr Metab Cardiovasc Dis.* 2024; 34(6): 1456–66.
18. Yi J, Wang L, Song J, Liu Y, Liu J, Zhang H, et al. Development of a machine learning-based model for predicting individual responses to antihypertensive treatments. *Nutr Metab Cardiovasc Dis.* 2024; 34(7): 1660–9.
19. Vu T, Kokubo Y, Inoue M, Yamamoto M, Mohsen A, Martin-Morales A, et al. Machine learning approaches for stroke risk prediction: Findings from the Suita study. *J Cardiovasc Dev Dis.* 2024; 11(7): 207.
20. Delpino FM, Costa ÂK, César do Nascimento M, Dias Moura HS, Geremias Dos Santos H, Wichmann RM, et al. Does machine learning have a high performance to predict obesity among adults and older adults? A systematic review and meta-analysis. *Nutr Metab Cardiovasc Dis.* 2024; 34(9): 2034–45.