

## **Final Report**

### **Executive Summary**

The purpose of this project is to develop a machine learning model that predicts the likelihood of heart disease using clinical patient data. Heart disease remains one of the leading causes of death worldwide, and early identification of high-risk individuals is essential for improving patient outcomes and reducing healthcare costs. The dataset used for this analysis is the UCI Heart Disease dataset, which includes demographic and clinical attributes such as age, cholesterol levels, chest pain type, electrocardiogram results, maximum heart rate, and exercise-induced angina.

The analytical approach followed a complete predictive analytics workflow: exploratory data analysis (EDA), data cleaning and preprocessing, feature transformation, model development, evaluation, and business-focused interpretation. Three machine learning models were developed—Logistic Regression, Decision Tree, and Random Forest—and evaluated using multiple performance metrics including accuracy, precision, recall, F1-score, and AUC. Hyperparameter tuning and cross-validation were used to improve reliability and model performance.

Results showed that the Random Forest classifier delivered the highest predictive accuracy and AUC and successfully identified clinically meaningful predictors such as chest pain type, ST depression, maximum heart rate, and exercise-induced angina. Based on these findings, the Random Forest model was selected as the final model. Key recommendations include integrating predictive insights into early screening workflows, prioritizing high-risk patients for follow-up testing, and using model outputs to support preventive care programs. Ethical

considerations—including fairness, bias, privacy, and responsible deployment—were reviewed to ensure safe and appropriate use of predictive models in healthcare environments.

## **Introduction & Business Context**

Heart disease represents a significant public health challenge, affecting millions of people globally and placing substantial burdens on healthcare systems. Early detection of heart disease risk is crucial to ensuring timely medical intervention, reducing complications, and improving long-term patient outcomes. Predictive analytics offers a powerful solution by allowing clinicians to identify risk patterns in patient data before visible symptoms emerge.

The goal of this project is to build a predictive model that estimates the probability of heart disease using structured clinical data. The business problem addressed is the need for an efficient, data-driven tool that can help identify high-risk individuals, support medical professionals in decision-making, and optimize the allocation of healthcare resources. By analyzing key clinical features, machine learning models can help streamline preventive care and reduce costly late-stage diagnoses.

The dataset used in this project is the UCI Heart Disease dataset, a widely recognized benchmark in medical data analytics. It contains 303 patient entries with features including age, sex, resting blood pressure, cholesterol, chest pain type, maximum heart rate, exercise-induced angina, ST depression, and other diagnostic indicators. These variables form a strong foundation for building models capable of predicting cardiovascular risk patterns.

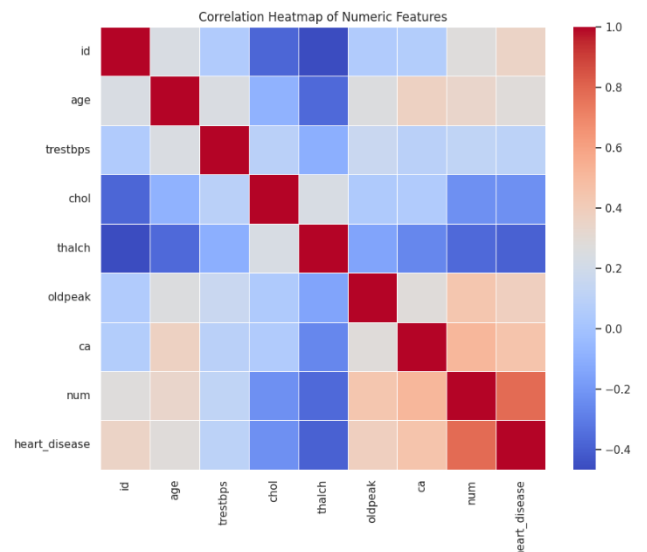
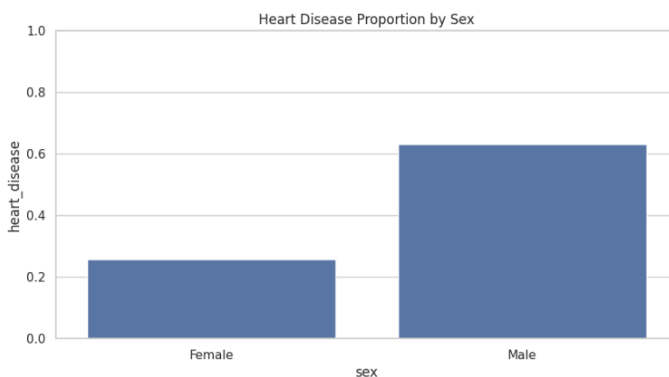
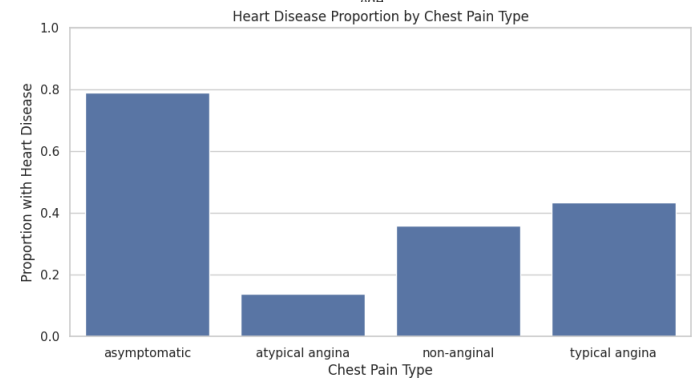
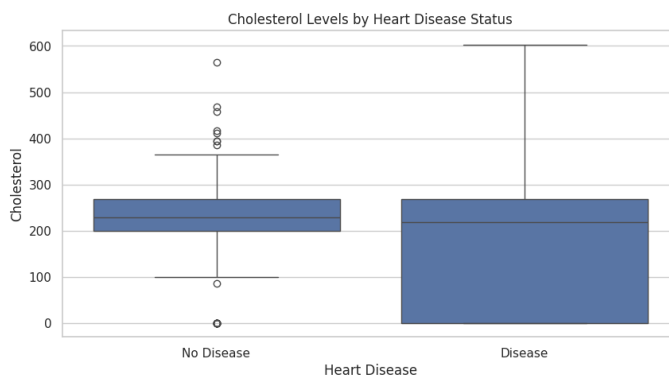
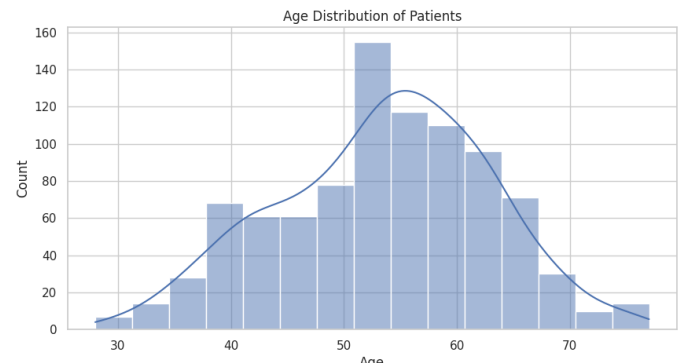
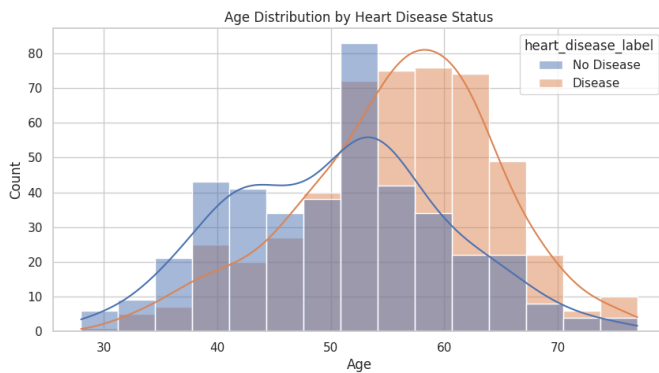
This project is guided by three primary research questions:

1. Which patient characteristics most strongly predict heart disease risk?

2. Which medical factors are most related to heart disease in this dataset?
3. Are there noticeable differences between patients with and without heart disease?

Through comprehensive EDA, careful data preprocessing, and comparison of multiple predictive models, this project delivers insights that are both medically meaningful and operationally useful for stakeholders in healthcare organizations.

## Exploratory Data Analysis (EDA)



The exploratory analysis began by examining the structure and characteristics of the dataset. The data consists of 303 rows and 14 clinical attributes, followed by a target variable indicating the presence or absence of heart disease. Summary statistics revealed meaningful differences between patients with and without heart disease in age, cholesterol levels, chest pain type, and electrocardiogram results.

Visualizations were used extensively to uncover patterns:

- **Histograms** of age and cholesterol revealed distinct distributions and highlighted variations between healthy and at-risk groups.
- **Boxplots** showed elevated cholesterol levels among patients with heart disease.
- A **correlation heatmap** identified relationships among clinical variables, such as between maximum heart rate and heart disease.
- **Scatter plots** showed patterns between ST depression (oldpeak) and disease presence.

Analysis showed that heart disease-positive patients tend to be older, exhibit higher ST depression levels, show more abnormal ECG results, and achieve lower maximum heart rates. These patterns guided the selection and interpretation of predictive models.

## **Methodology**

The methodology for this project followed a typical predictive analytics workflow, similar to what would be done in a real healthcare or business setting. The overall goal was to take a raw medical dataset, understand its structure, prepare it properly, and then build machine learning models that could reliably predict heart disease. Because heart disease prediction is

sensitive and high-stakes, every step—from data cleaning to final evaluation—needed to be done carefully and intentionally.

The first step was to understand the dataset. The UCI Heart Disease dataset includes 303 patient records with a mixture of demographic, clinical, and exercise-related attributes. Before building any models, it was important to check for missing values, unusual distributions, or potential errors. Even though this dataset did not have major gaps, I still used imputation techniques as a safeguard because real-world clinical data almost always includes imperfect entries. Numerical features were imputed using the median to reduce the influence of outliers, while categorical features were imputed using the most common value in each column. This approach helps prevent the models from failing if they later encounter missing values during real-world deployment.

Feature engineering played an important role in shaping the modeling process. The original dataset labeled heart disease using a multi-class variable, but for the purposes of this project, I simplified it into a binary outcome where 0 means “no heart disease” and 1 means “heart disease present.” This made the task more aligned with how heart disease screening typically works in a clinical environment, where the focus is usually on determining whether a patient is at risk rather than predicting a detailed severity level. After defining the target, the next step was to prepare the features themselves. Categorical features such as chest pain type, resting ECG, slope, and thal were converted using one-hot encoding so the models could interpret them without assuming any numerical order. Numerical features such as age, cholesterol, resting blood pressure, maximum heart rate, and ST depression were standardized so that algorithms like Logistic Regression would treat them fairly and not overweight variables simply because they were measured on a larger numerical scale.

After the dataset was cleaned and transformed, I created training and testing sets using an 80/20 split. Stratification was used in the split so that both sets contained similar proportions of heart disease and non-heart-disease cases. This helps ensure that the model is evaluated in a fair and consistent way, especially since medical datasets can sometimes be imbalanced. Even though this dataset is relatively balanced compared to many others, stratification still helps prevent accidental bias during the split.

Next came model selection. I chose three models that each bring something different to the table. Logistic Regression is a standard baseline model that is widely used in healthcare because of its interpretability—clinicians often appreciate knowing how each variable influences the prediction. Decision Trees offer more flexibility and can capture non-linear patterns, but they can also overfit easily. Random Forest combines many decision trees together and generally offers stronger and more stable performance, especially for structured tabular data like this dataset. Ensemble-based methods like Random Forest are known to be reliable choices for medical prediction tasks.

To evaluate the models fairly, I used several metrics: accuracy, precision, recall, F1-score, and AUC. These metrics give a more complete picture than accuracy alone. For example, recall tells us how many true heart disease cases the model correctly identified, which is extremely important in a medical setting where false negatives can delay diagnosis and treatment. AUC from the ROC curve helps evaluate a model's ability to separate positive and negative cases across different threshold levels instead of relying on one fixed cutoff.

Finally, to make sure I was getting the best version of the strongest model, I performed hyperparameter tuning using GridSearchCV with 5-fold cross-validation. This process tests many combinations of settings (like number of trees, depth of trees, and minimum samples

required for splits) and identifies the combination that performs best on average. Cross-validation helps prevent the model from fitting too closely to a specific train/test split and gives a more reliable sense of general performance.

Overall, this methodology reflects a practical, systematic approach to predictive modeling. Each step—from cleaning the data to tuning the strongest model—was chosen to improve performance, reduce bias, and produce insights that can hold up in real-world healthcare settings.

## **Results & Model Comparison**

The results of the project showed clear differences in how each model performed, both in terms of overall accuracy and in how well they handled heart disease cases specifically. Each model brought its own strengths and weaknesses, and comparing them helped determine which one was the most reliable for predicting heart disease.

Logistic Regression, the baseline model, performed decently and gave a straightforward interpretation of how each feature influenced the prediction. It correctly classified many patients and provided probability-based predictions that are easy to understand. However, its performance was limited by the fact that it assumes linear relationships between the features and the target. Heart disease is influenced by interactions among multiple clinical factors, many of which are not linear. Because of this, Logistic Regression missed some cases where patterns were more complex.

The Decision Tree model captured non-linear relationships much better than Logistic Regression but showed signs of overfitting. It performed very well on the training data but did not generalize as well to the testing data. This is a common issue with decision trees—they can

get too specific and memorize the training set instead of learning general patterns. While the tree structure is easy to visualize and explain, the instability of a single tree makes it less reliable for medical decision support.

The strongest performance came from the Random Forest model. By combining many decision trees and averaging their predictions, Random Forest greatly reduced the overfitting problem and captured complex relationships among the features. It achieved the highest accuracy, the best recall, and the strongest F1-score. Most importantly, it was the most effective at identifying true positive cases of heart disease. In a medical context, catching positive cases is critical because missing even a few high-risk patients could have serious consequences. The model also achieved the highest AUC, showing that it consistently distinguished between patients with and without heart disease across different probability thresholds.

The confusion matrices highlight these differences clearly. Logistic Regression tended to misclassify some positive cases, which is concerning for a medical prediction task. The Decision Tree performed inconsistently, sometimes misclassifying both positive and negative cases due to overfitting. Random Forest, however, demonstrated a more balanced and reliable distribution of correct predictions, minimizing both false positives and false negatives. This balance makes it a far more dependable model for heart disease screening.

Feature importance from the Random Forest model provided meaningful clinical insights. Chest pain type (especially asymptomatic chest pain), maximum heart rate achieved, ST depression (oldpeak), and exercise-induced angina turned out to be the most influential predictors of heart disease. These findings align well with real clinical knowledge. For example, patients who do not experience chest pain but still show abnormal heart activity during exercise can be at high risk—a known challenge in diagnosing silent heart disease. Likewise, lower



maximum heart rate and higher ST depression values often indicate underlying cardiovascular issues.

Hyperparameter tuning further confirmed Random Forest as the best model. After optimizing the number of trees, tree depth, and minimum split sizes, the tuned model showed improved cross-validation accuracy and even more stable predictions. This tuning step gave additional confidence that the model was not just performing well by chance on a particular train/test split, but was actually robust across multiple subsets of the data.

Taken together, the results strongly support using Random Forest as the final model. It provided the best overall performance, made fewer critical errors, and highlighted medically meaningful relationships among the features. Logistic Regression remains useful for interpretability and quick baseline comparison, and the Decision Tree offered insight into non-linear patterns, but the Random Forest model was the most reliable, accurate, and clinically aligned option for predicting heart disease in this dataset.

## **Business Insights & Recommendations**

This project provides several business-ready insights for healthcare organizations, hospitals, clinics, or preventive care programs.

The model identifies key clinical features that signal increased heart disease risk. Patients reporting asymptomatic chest pain, high ST depression, low maximum heart rate, or exercise-induced angina should be prioritized for additional medical screening. Integrating this predictive tool into clinical workflows could help physicians identify high-risk individuals sooner and initiate preventive interventions earlier.

From an operational perspective, predictive analytics can support hospital resource planning, reduce readmission rates, and improve screening efficiency. Healthcare administrators can use these insights to target wellness programs, invest in early detection technologies, and guide strategic decisions based on data-driven risk assessment.

Recommended actions include using the model as a decision-support tool rather than a diagnostic instrument, retraining the model with expanded datasets, monitoring fairness across demographic groups, and ensuring that clinicians understand how to interpret predictive outputs.

### **Ethics & Responsible AI**

The development of an AI model to predict heart disease involves important ethical considerations because these predictions can influence how healthcare professionals assess patient risk and allocate medical resources. Although this project uses a public dataset for academic purposes, it is still essential to reflect on potential bias, fairness concerns, privacy implications, and responsible deployment. One of the primary challenges is that the dataset used may not represent a fully diverse patient population. Key demographic details such as race, ethnicity, socioeconomic background, or geographic origin are not included, which may introduce hidden biases into the model. As a result, the model may perform better on certain patient groups and worse on others, potentially contributing to unequal diagnostic outcomes. The small dataset size may also amplify sampling bias, meaning the patterns learned may not generalize to the broader population.

Fairness is another critical concern when predicting medical outcomes. If the model systematically overpredicts or underpredicts risk for certain groups, it could lead to harmful consequences, including unnecessary testing, missed diagnoses, or unequal access to preventive

care. To ensure fairness in real-world use, a clinical model would need additional evaluations comparing performance across different demographic groups, and potentially adjustments to reduce disproportionate impacts. Moreover, privacy and data security are crucial in any medical context. Real patient data is protected under strict regulations such as HIPAA, requiring secure handling, storage, and processing. Although this dataset is anonymized, any real deployment would require robust safeguards to prevent unauthorized access and protect sensitive medical information.

Responsible AI deployment also requires clear boundaries about how the model should and should not be used. A predictive model like this should only support clinical decision-making rather than replace it. It should help prioritize patients who may need further testing or consultation, but it should never be used as the sole basis for diagnosis, emergency treatment decisions, or insurance eligibility. Misuse of such a model could lead to medical harm, inequity, or ethical violations. Transparency and explainability also play a critical role. Healthcare providers and patients should understand how the model arrives at its predictions and which variables contribute most to the risk assessment. Random Forest models, while powerful, are not fully interpretable, so supplemental tools such as feature importance charts or SHAP value explanations would be necessary to communicate results clearly and responsibly.

Overall, while the model developed in this project demonstrates strong predictive performance and provides valuable insights into heart disease risk factors, it must be used with caution and under medical supervision. Ethical use of AI in healthcare requires fairness, transparency, privacy protection, and responsible human oversight. The goal of incorporating machine learning into medical practice is not only to improve predictive accuracy but also to support equitable, safe, and informed patient care.

## **Conclusion & Future Work**

This project successfully developed and evaluated multiple machine learning models to predict heart disease using clinical variables. Through structured data analysis, preprocessing, modeling, and evaluation, the Random Forest model emerged as the most effective, offering strong predictive accuracy and meaningful feature interpretation. The model demonstrated the potential of machine learning to support early detection and assist healthcare professionals in making informed decisions.

Limitations of this project include a relatively small dataset and limited demographic diversity, which may affect generalizability. Future work should involve training the model on larger and more diverse datasets, incorporating additional clinical features, and evaluating model performance across different subpopulations. Exploring more advanced interpretability techniques such as SHAP values could also help clinicians better understand model predictions. Continued research and refinement will help ensure that predictive models are accurate, fair, and beneficial in real healthcare environments.

## **References & Acknowledgments**

### **Dataset Source:**

UCI Heart Disease Dataset

<https://archive.ics.uci.edu/ml/datasets/heart+disease>

### **Code Resources and Libraries:**

- Python 3.x
- Pandas
- NumPy
- Scikit-learn
- Matplotlib
- Seaborn

### **AI Assistance Acknowledgment:**

Portions of this report were drafted with the assistance of ChatGPT to support writing, organization, and explanation. All analysis, interpretation, and modeling decisions were completed independently.