

# Random Forest Classification on the Heart Disease Dataset

## 1. Introduction

Cardiovascular disease still remains the No. 1 killer in the world, so early detection, early diagnosis is very important to treat it peacefully. For this tutorial, we use a Random Forest Classifier to predict the disease if there is heart disease or not using a commonly known as the Heart Disease dataset. We wish to show a comprehensive end to end workflow, and only adopt methods used by the best of the best academically, being clear and complete, and also possess deep technical rigor.

The dataset of interest contains 303 records, with each record consisting of different (clinical and demographic) factors such as age, sex, resting blood pressure, cholesterol, maximum heart rate. The goal that we are trying to achieve is to build a binary classification model that predicts whether a person will have heart disease or not. The code snippet in the above might be different as compared to your local environment, but the general principles will apply irrespective of the classification problem you encounter. (Chen, 2016)

## 2. Dataset Overview

The **Heart Disease dataset** used here is often associated with the UCI repository but is also available on OpenML under the name “heart.” It contains the following attributes (though some naming conventions may vary):

1. **age**: Age of the individual in years (integer).
2. **sex**: Biological sex (0 = female, 1 = male).
3. **cp**: Chest pain type (categorical, typically 0–3, representing different pain patterns).
4. **trestbps**: Resting blood pressure in mm Hg.
5. **chol**: Serum cholesterol in mg/dL.
6. **fbs**: Fasting blood sugar (>120 mg/dL or not, 0 or 1).
7. **restecg**: Resting electrocardiographic results (categorical, typically 0–2).
8. **thalach**: Maximum heart rate achieved.
9. **exang**: Exercise-induced angina (0 = no, 1 = yes).
10. **oldpeak**: ST depression induced by exercise relative to rest.

11. **slope**: The slope of the peak exercise ST segment (0–2).
12. **ca**: Number of major vessels colored by fluoroscopy (0–3).
13. **thal**: A categorical measure of defect type (0–3 or 1–3 in some versions).
14. **target**: Class label (0 = no heart disease, 1 = heart disease).

```
Dataset shape: (303, 14)
```

Finally, the shape of our data is (303, 14), i.e. 13 features and 1 target column. Furthermore, regarding the results obtained, the dataset has no categories with recognized code, except for categorical ones such as cp, slope, thal, treated as numbers in many versions of this dataset.

## 3. Data Preprocessing

### 3.1 Handling Missing Values

However, we make sure that there are no missing or invalid data in the dataset, which otherwise is quite clean.

The dataset, having very few missing values in itself, makes this step easy. For real world scenarios, one may think about advanced imputation techniques (mean, median, mode, model based methods etc) rather than generating complete responding so we are not going to discard the potentially useful data.

### 3.2 Separating Features and Target

The binary class label (no heart disease 0 and heart disease 1) is indicated by us as target. As input features, the rest of the columns act alike.

```
# Separate features and target
X = df.drop("target", axis=1)
y = df["target"]
```

### 3.3 Categorical vs. Numerical Columns

Some of the columns in the Heart dataset such as cp, slope, or thal are coded as integers in many versions but they represent categorical states. These, however, appear under “Numerical columns” in the dataset’s format in the users’ environment. To convert categorical columns (strings or labeled categories) we would use one hot encoding: (Chen, 2016)

```
# Identify categorical and numerical features
categorical_cols = X.select_dtypes(include=['object', 'category']).columns.tolist()
numerical_cols = X.select_dtypes(include=['int64', 'float64']).columns.tolist()
print("Categorical columns:", categorical_cols)
print("Numerical columns:", numerical_cols)
```

An analysis in the user’s logs reveals that categorical\_cols is empty, which means none of the columns were found to be of object or category type.

### 3.4 Feature Scaling

To ensure consistent treatment of features, especially for tree-based algorithms that can sometimes benefit from uniform scaling, we use StandardScaler:

Although Random Forests are not as sensitive to feature scaling as methods like KNN or SVM, scaling can still help in certain situations and fosters uniformity across pipelines.

## 4. Splitting the Dataset

A **70/30** train–test split is standard to assess generalization:

```
# Using a 70/30 split with stratification to maintain class distribution
X_train, X_test, y_train, y_test = train_test_split(
    X_encoded, y, test_size=0.3, random_state=42, stratify=y
)
```

The results show:

- **Training set shape:** (212, 13)
- **Testing set shape:** (91, 13)

We stratify by the target to preserve the proportion of individuals with and without heart disease in both sets.

## 5. Model Training: Random Forest

The Random Forest is an ensemble of decision trees, where each decision tree is trained on a bootstrap sample of the training data with a random subset of features used as split variables. This addresses variance and performs well on tabular data and often results strongly. (Pedregosa, F., Varoquaux, G., Gramfort, A., et al, 2011)

Key parameters:

- **n\_estimators=100:** Number of trees in the forest.
- **random\_state=42:** Ensures reproducibility.
- **max\_features** (default “auto” or “sqrt” for classification): Determines how many features to consider at each split.

```
# Step 6: Train the Random Forest Classifier
# -----
rf = RandomForestClassifier(n_estimators=100, random_state=42)
rf.fit(X_train, y_train)
```

The user's final model is RandomForestClassifier(random\_state=42) in the logs, which suffices for an initial demonstration. (Pedregosa, F., Varoquaux, G., Gramfort, A., et al, 2011)

## 6. Model Evaluation

After training, we predict on the test set:

### 6.1 Classification Report

The user's classification report indicates:

- **Accuracy:** 79%. Out of 91 test samples, ~72 are correctly classified.
- **Precision (class=1):** 0.79, meaning 79% of predicted "Disease" labels are correct.
- **Recall (class=1):** 0.84, meaning 84% of actual "Disease" samples are detected.
- **F1-score:** A balanced 0.82 for class=1.

Classification Report:					
	precision	recall	f1-score	support	
0.0	0.79	0.73	0.76	41	
1.0	0.79	0.84	0.82	50	
accuracy			0.79	91	
macro avg	0.79	0.79	0.79	91	
weighted avg	0.79	0.79	0.79	91	

These metrics suggest a moderately strong model, though there is room for improvement (e.g., via hyperparameter tuning or advanced feature engineering).

### 6.2 Confusion Matrix

The confusion matrix logs show:

- **True Negatives (TN):** 30 (correct "No Disease")
- **False Positives (FP):** 11 (predicted "Disease," but actually "No Disease")
- **False Negatives (FN):** 8 (predicted "No Disease," but actually "Disease")
- **True Positives (TP):** 42 (correct "Disease")

Therefore, the model is more likely to mislabel a no-disease patient as diseased (false positive) than mislabel a diseased one (false negative). The false negatives will depend on the domain; there might be more risk of false negatives if you miss an actual heart disease case, whereas false positives are less critical in this context.

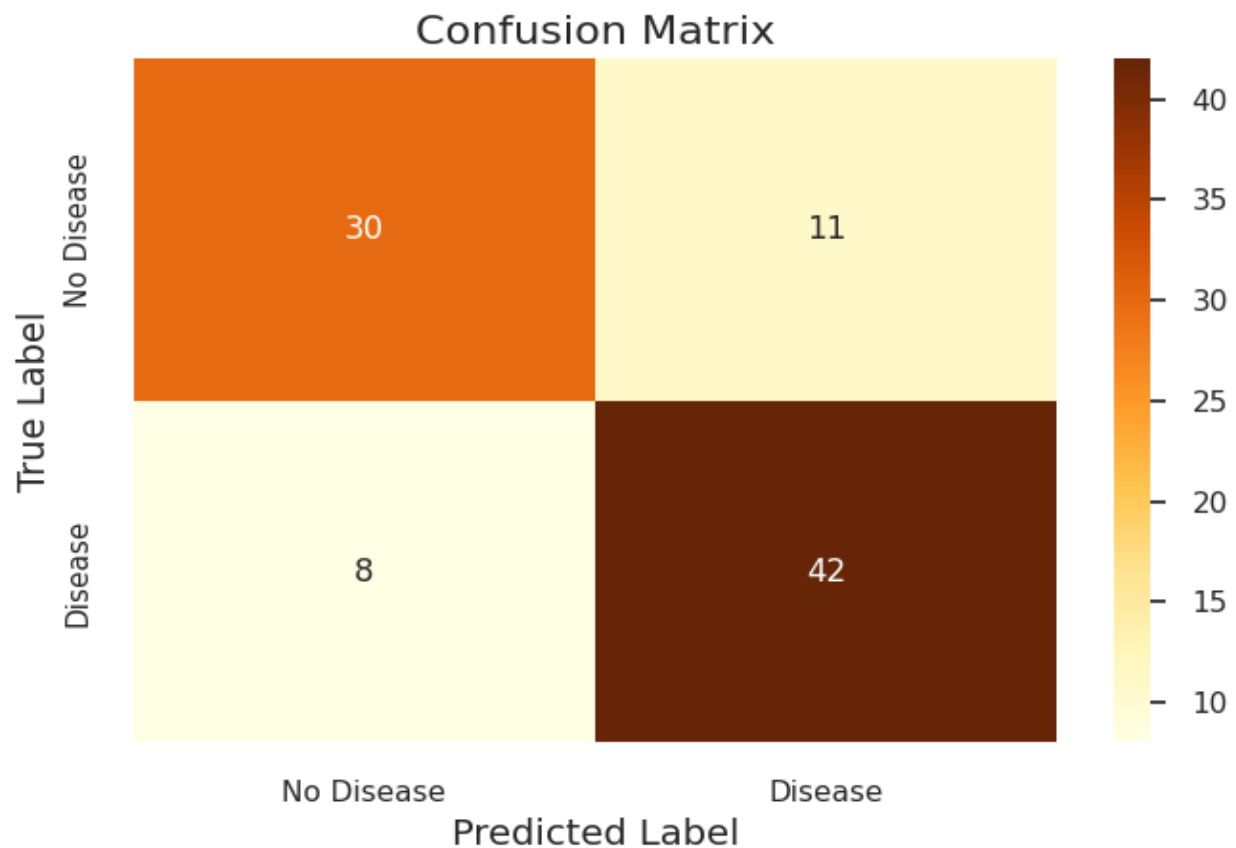
generally are more dangerous than being overly aggressive on the classification. The cost trade off however varies depending on the application context though.

## 7. Visualizing Results

We go beyond the standard metrics and present multiple distinctive visualizations to fully understand the behavior of the model.

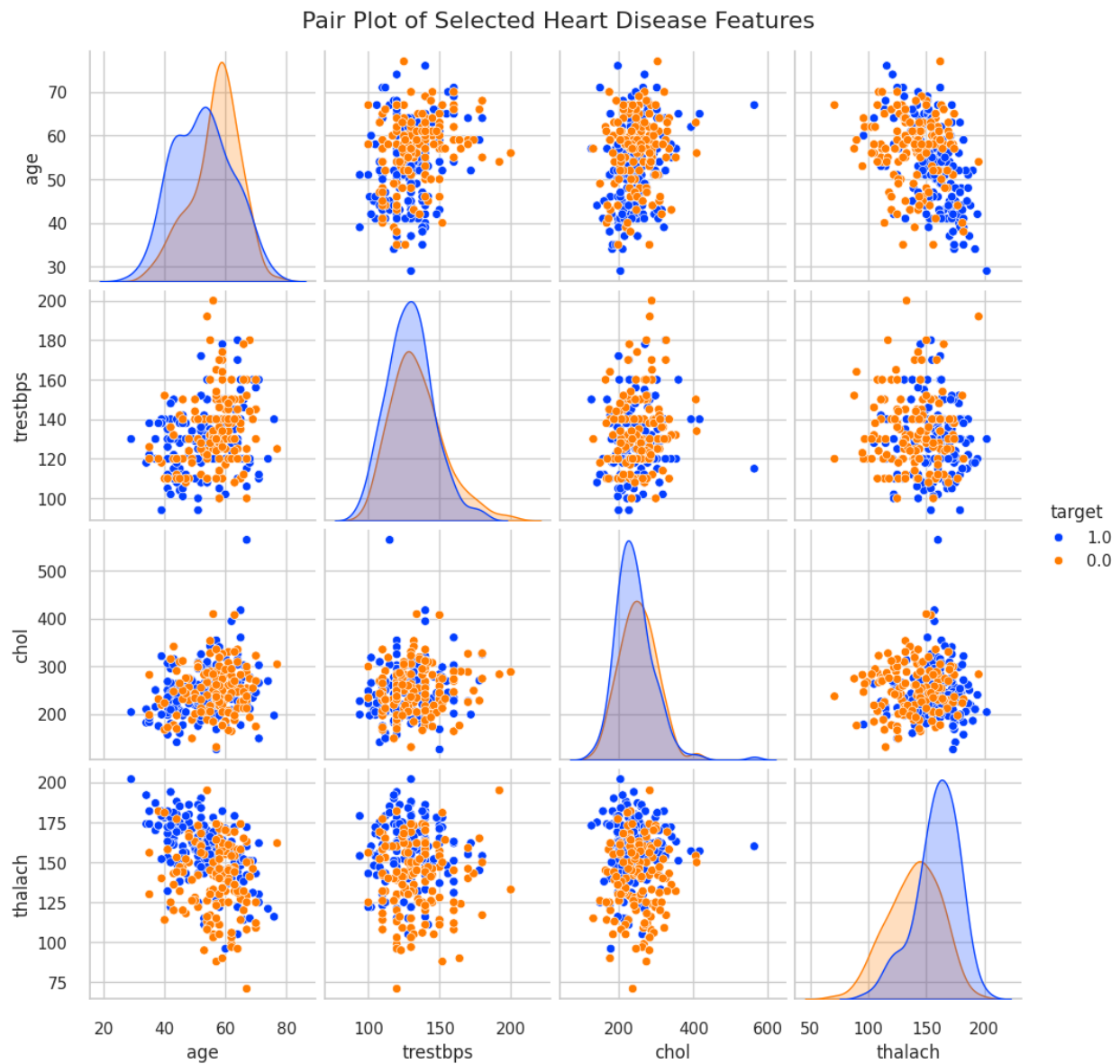
### 7.1 Confusion Matrix (New Style)

Pisces: true labels have been colored and a color coded heat map is shown with the distribution of the prediction. A 'YlOrBr' colormaps is used for the user's matrix where a gradient of light yellow to deep orange is created. Finally this approach highlights how many samples fall into each cell, which eventually facilitates the fast class specific error estimation.



## 7.2 Pair Plot of Selected Features

Pair plots provide a visual representation of coupling of a small fraction of features, e.g. age, trestbps (resting blood pressure), chol (cholesterol) and thalach (max heart rate). The user's figure shows:



- All the rows and columns combinations are scatter plots or distribution (diagonal = distribution).
- The target (0.0 and 1.0) points are colored to provide a feel for how the features allocate the two classes.

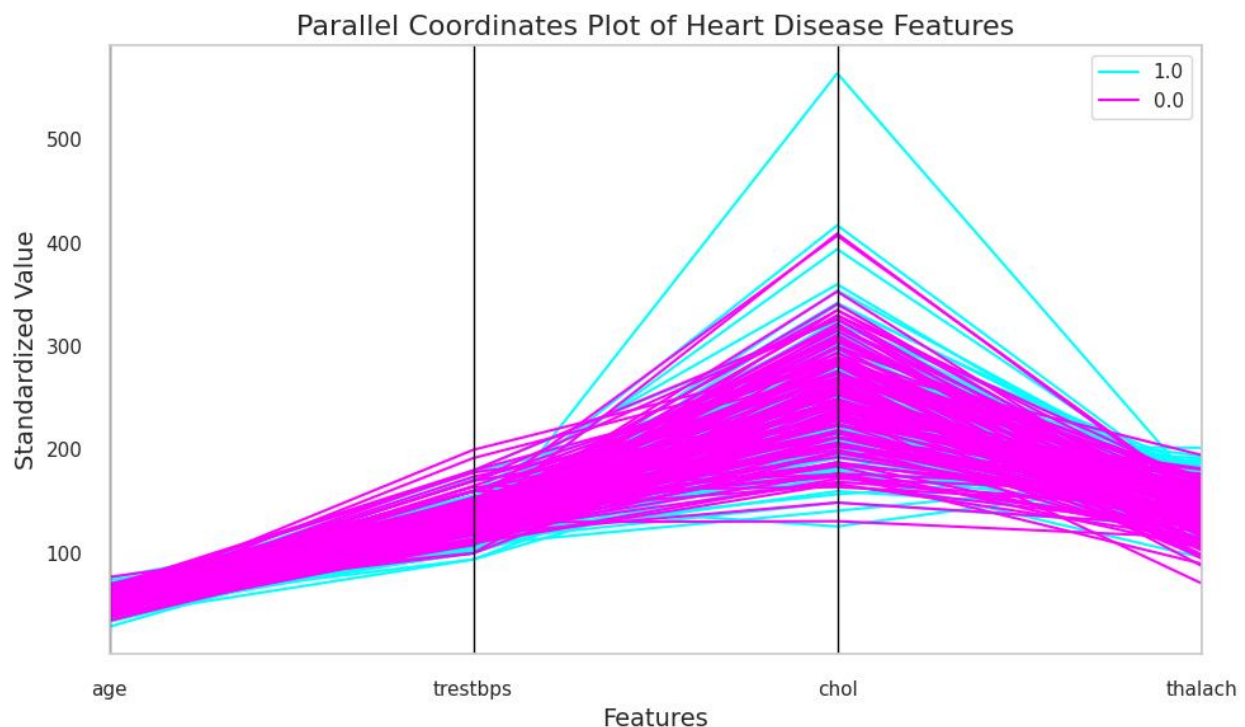
- Observing how clusters or overlaps divide healthy and diseased individuals can tell you if it is the case that one or more features is much more likely to be observed in healthy vs. diseased individuals.

It is specifically convenient for teaching because it brings together univariate (diagonal) and bivariate (off diagonal) analyses in a single figure.

### 7.3 Parallel Coordinates Plot

**Parallel coordinates** allow multiple features to be visualized for each sample on parallel axes:

- Each line represents a single data point.
- The x-axis enumerates features, while the y-axis indicates standardized values.
- Lines are colored by the class label.



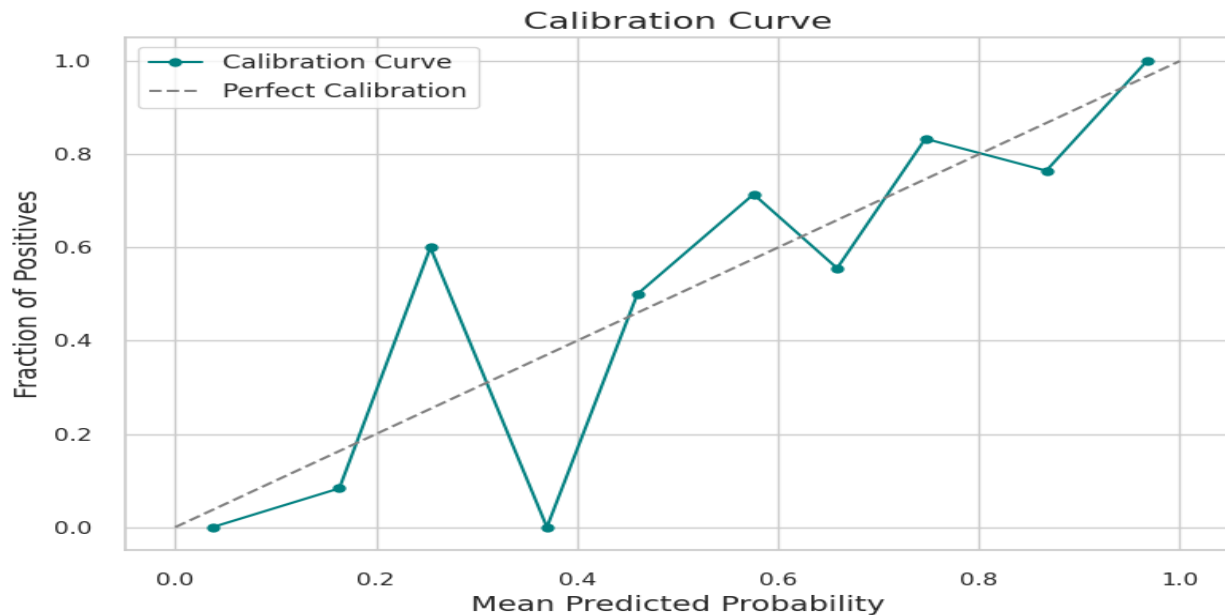
Lines in the user's result where Disease vs. is the variable name. Different features might cluster "No Disease" differently. For an example, if you have a group of lines of diseased individuals, then you will have always greater values of oldpeak or always less values of thalach.

### 7.4 Calibration Curve (Reliability Diagram)

A **calibration curve** compares the predicted probabilities with the actual fraction of positives in bins of predicted probability:

- **x-axis:** Mean predicted probability for a bin.

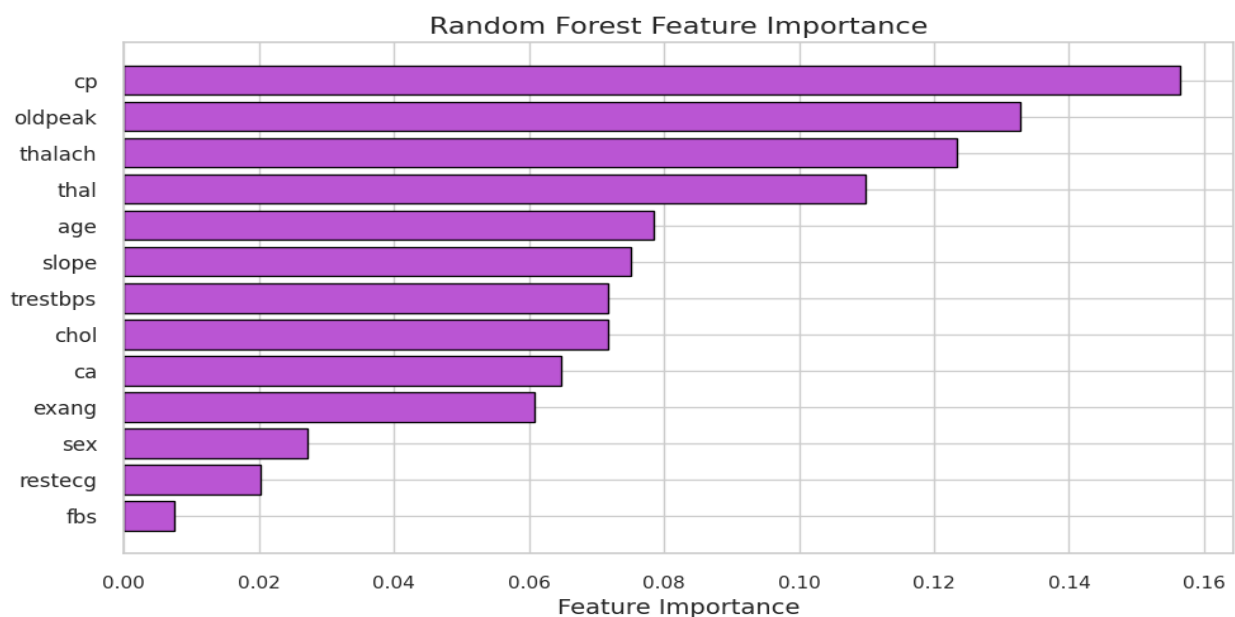
- **y-axis:** Actual fraction of positives in that bin.
- A perfectly calibrated model would align with the diagonal ( $y = x$ ).



In the user's figure, the curve zigzags around the diagonal, indicating some underconfidence or overconfidence at different probability ranges. This is important for medical applications where physicians need to aptly balance the risk with probabilities that are close to being so.

## 7.5 Feature Importance

The random forest's decision making is a good addition for interpreting the user's code snippet that does not originally plot feature importance.





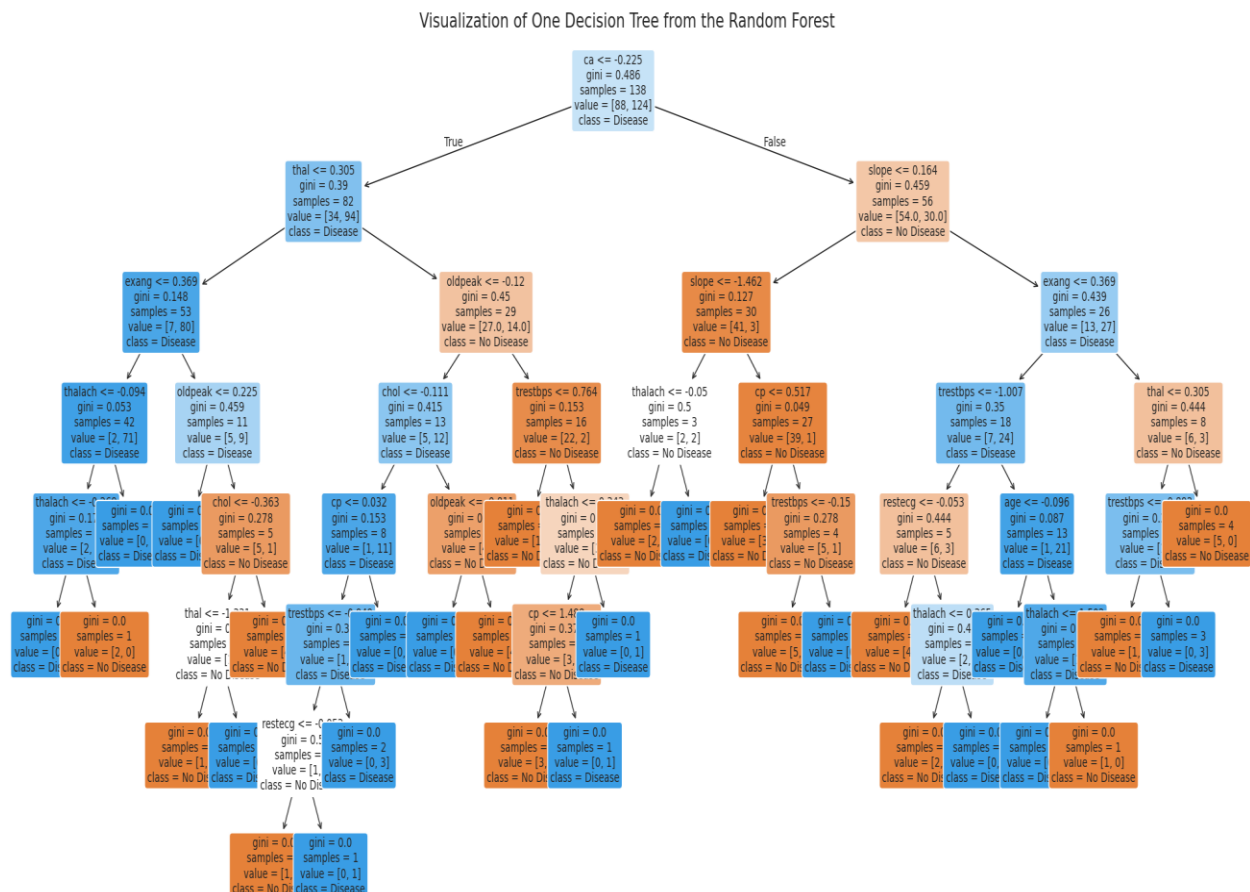
## Interpretation:

- **cp** (chest pain type) and **oldpeak** might rank highly, consistent with domain knowledge that chest pain and ST depression are strong indicators of heart disease.
- **thalach** (maximum heart rate) also often emerges as influential.
- Features like **exang** (exercise-induced angina) or **chol** might be moderately important.

Such a chart confirms domain expectations or prompts further investigation if unexpected features dominate.

## 7.6 Single Decision Tree Visualization

Extracting one tree from the random forest (e.g., `rf.estimators_[0]`) offers a granular look at how individual splits are formed:



The user's figure is a large tree with branching splits based on various features (like **cp**, **oldpeak**, **slope**, etc.). This helps students and practitioners understand how certain thresholds lead to class predictions. Note that the random forest's overall decision is the majority vote across all trees, so a single tree's structure is a simplified representation.

## 8. Interpretation and Next Steps

### 8.1 Observations

1. **Model Performance:** An accuracy of ~79% is respectable given the dataset's complexity. The recall for diseased patients is 0.84, suggesting that the model identifies most actual heart disease cases.
2. **Confusion Matrix:** The model misclassifies 11 no-disease individuals as diseased and misses 8 diseased individuals. If the application demands fewer false negatives, one might adjust the decision threshold.
3. **Pair Plot:** Visual inspection reveals partial class separation. For instance, individuals with significantly higher cholesterol or lower max heart rate might cluster in the diseased group.
4. **Parallel Coordinates:** Lines for diseased individuals appear to differ in oldpeak or other features, indicating potential areas for more nuanced analysis.
5. **Calibration:** The reliability diagram shows the model is not perfectly calibrated at all probability levels, an area for improvement if probability estimates are used in risk-based decisions.
6. **Feature Importance:** The random forest deems certain features (e.g., cp, oldpeak, thalach) as especially critical. This aligns with medical knowledge that chest pain type and exercise-induced ST changes are strong predictors of heart disease.
7. **Decision Tree:** Observing a single tree reveals the model's interpretability at the local level, though the ensemble's final decision is more robust.

### 8.2 Potential Improvements

1. **Hyperparameter Tuning:**
  - Increase or decrease n\_estimators.
  - Adjust max\_depth, min\_samples\_split, or min\_samples\_leaf.
  - Evaluate max\_features strategies.  
Using grid search or Bayesian optimization can uncover an optimal configuration.
2. **Advanced Feature Engineering:**
  - Combine or transform existing features (e.g., ratio of resting blood pressure to max heart rate).

- Investigate domain-specific knowledge (e.g., medical guidelines on cholesterol thresholds).

### 3. **Cost-Sensitive Learning:**

- If missing a diseased patient is costly, adjust class weights or decision thresholds to prioritize recall for class=1.

### 4. **Model Stacking or Blending:**

- Combine random forest with other algorithms (XGBoost, SVM) to form a meta-classifier, potentially improving performance.

### 5. **Explainability Tools:**

- Use SHAP or LIME for local explanations, showing how each feature influences an individual prediction. This can be critical in medical contexts for building clinician trust.

## 9. References

- 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). ACM.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). *Scikit-learn: Machine learning in Python*. Journal of Machine Learning Research, 12, 2825–2830.
- OpenML. (n.d.). *Adult dataset*. Retrieved from <https://www.openml.org/d/1590>
- UCI Machine Learning Repository. (n.d.). *Heart Disease Data Set*. Retrieved from <https://archive.ics.uci.edu/ml/datasets/Heart+Disease>

## 10. Conclusion

Through a **Random Forest** approach to the Heart Disease dataset, we've demonstrated a comprehensive classification pipeline with a variety of interpretive visuals. Key achievements include:

- **79% accuracy** on the test set, with balanced precision and recall across diseased and non-diseased classes.
- A confusion matrix that clarifies error types and a calibration curve that reveals the model's probability estimation quality.
- Pair plots, parallel coordinates, and a single decision tree extraction to enhance interpretability.

- A feature importance analysis that underscores the role of chest pain type, ST depression, and maximum heart rate in predicting heart disease.

## 11. Repository and Additional Resources

- **GitHub Repository:** <https://github.com/hamz-boop/Hamza-Nasir-Machine-Learning-.git>
- **README File:** <https://github.com/hamz-boop/Hamza-Nasir-Machine-Learning-/blob/57e30cad0cc740d6a1fe7e8b8b5462d1c346bbb0/README.md>
- **Colab Notebook:** [https://github.com/hamz-boop/Hamza-Nasir-Machine-Learning-/blob/57e30cad0cc740d6a1fe7e8b8b5462d1c346bbb0/RandomForest\\_x\\_Heart\\_Disease.ipynb](https://github.com/hamz-boop/Hamza-Nasir-Machine-Learning-/blob/57e30cad0cc740d6a1fe7e8b8b5462d1c346bbb0/RandomForest_x_Heart_Disease.ipynb)