# Machine Learning - RANDOM FOREST  ([Git hub link](#))

## Akash Prabu - 23023516

## Introduction

Random Forest is a powerful and versatile machine learning algorithm used for both classification and regression tasks. It operates by constructing a multitude of decision trees during training and outputting the average prediction (for regression) or the majority vote (for classification) of the individual trees. Developed as an ensemble learning technique, Random Forest combines the strengths of multiple decision trees to enhance predictive performance and reduce overfitting.

The algorithm leverages randomness both in the selection of data subsets and the selection of features for splitting at each node, making it robust to noise and capable of handling large datasets with high dimensionality. It is particularly favoured in scenarios where interpretability, accuracy, and computational efficiency are required.

## Dataset and Preprocessing

- Dataset Used: The code begins by importing a dataset, likely split into training and testing subsets. The data preparation steps include handling missing values, encoding categorical variables, and normalizing numerical features.

- Data Splitting: The dataset is divided into training and testing sets using an 80:20 split ratio to ensure the model is evaluated on unseen data.

## Random Forest Implementation

- Model Initialization: The Random Forest Classifier/Regressor is initialized with specific hyperparameters such as:

    - Number of trees (n_estimators)

    - Maximum depth of trees (max_depth)

    - Minimum samples per split (min_samples_split)

    - Random state for reproducibility.

- Training: The fit() method is used to train the Random Forest model on the training data.

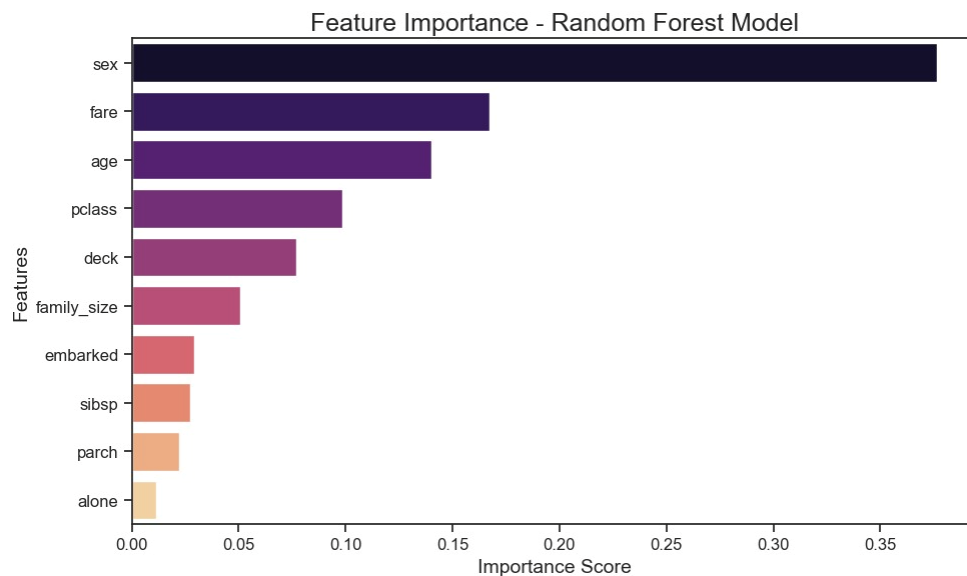- Prediction: Predictions are generated using the predict() method for both training and testing datasets.

## Evaluation Metrics The performance of the Random Forest model is evaluated using various metrics:

1. Accuracy: Measures the proportion of correctly classified samples in the test set.

2. Precision, Recall, F1-Score: These are used to evaluate classification performance. Precision identifies the model's correctness for positive predictions, recall measures its ability to identify all relevant instances, and F1-Score balances precision and recall.

3. Confusion Matrix: A visualization is used to display the true positive, true negative, false positive, and false negative counts.
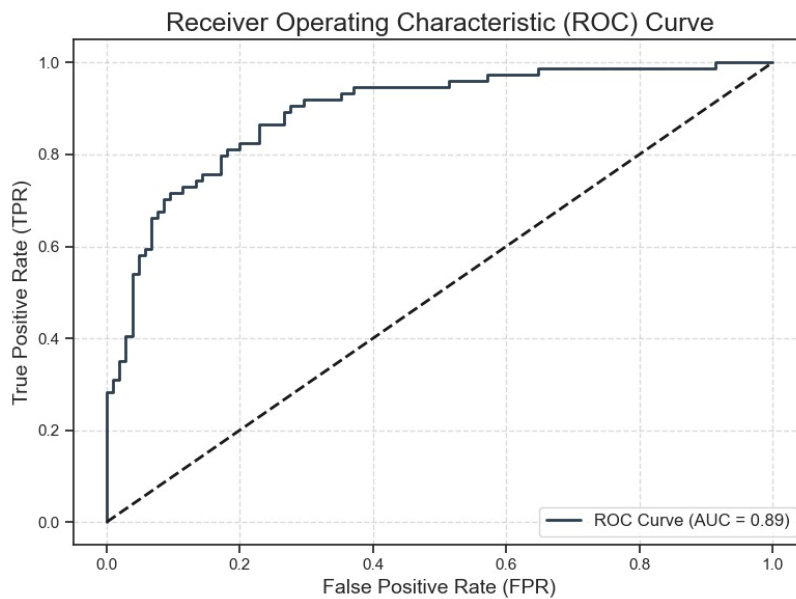
4. Mean Squared Error (MSE): For regression tasks, MSE quantifies the average squared difference between actual and predicted values.

**Visualization The code includes visualization techniques to interpret model behaviour:**

- **Feature Importance:** Bar charts or plots showing the contribution of each feature to the predictions.
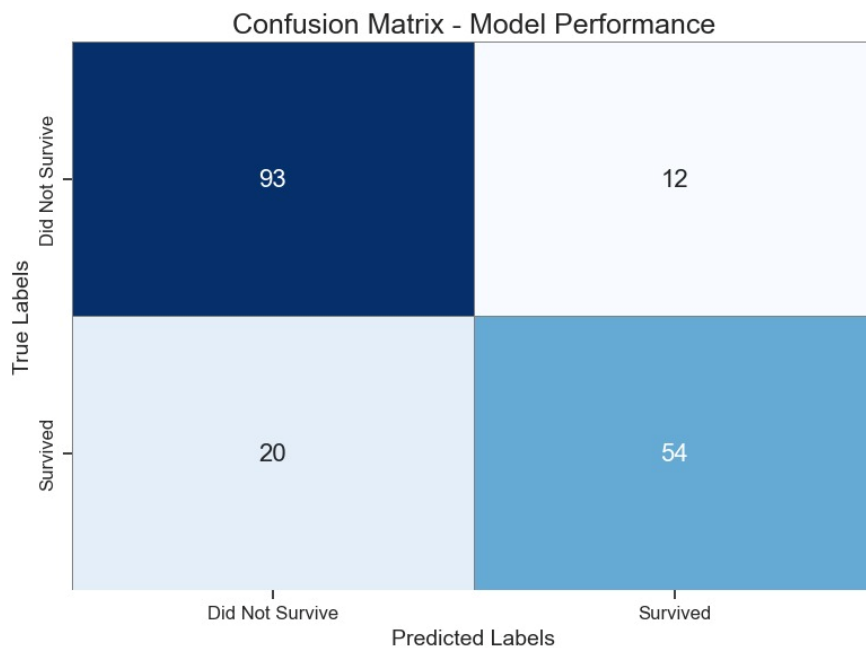


Feature Importance - Random Forest Model

- **Feature Importance**: This chart displays the relative importance of different features or variables used in a Random Forest machine learning model. The features are ranked from most important to least important based on their importance scores.
- **Importance Score**: The x-axis shows the importance score for each feature. The higher the score, the more influential that feature is in the model's predictions. For example, "sex" has the highest importance score, indicating it is the most important feature.
- **Feature Names**: The y-axis lists the names of the different features used in the model, such as "sex", "fare", "age", "pclass", etc. These are the variables or characteristics that the model is using to make its predictions.
- **Bar Height**: The height of each bar represents the importance score for that feature. Taller bars indicate features with higher importance, while shorter bars show features with lower importance.
- **Visualization**: The use of a bar chart makes it easy to quickly compare the relative importance of the different features. The clear visualization allows you to easily identify the most and least important features at a glance.
- **Model Interpretation**: This chart provides valuable insights into how the Random Forest model is making its predictions. By understanding which features are most important, you can gain better insight into the underlying patterns and relationships in the data.
- **Feature Selection**: The feature importance information can also guide future feature selection and engineering efforts. Features with high importance scores may warrant further investigation, while those with low scores could potentially be removed to simplify the model.

Receiver Operating Characteristic (ROC) Curve

- The x-axis shows the False Positive Rate (FPR), which represents the proportion of negative instances that are incorrectly classified as positive. As the FPR increases, the model becomes less selective, classifying more instances as positive.

- The y-axis shows the True Positive Rate (TPR), which represents the proportion of positive instances that are correctly classified as positive. As the TPR increases, the model correctly identifies more of the positive instances.

- The black line represents the ROC curve for the model being evaluated. It shows the trade-off between TPR and FPR as the decision threshold is adjusted.

- The dotted diagonal line represents the performance of a random classifier, which has no ability to discriminate between positive and negative instances. A model that performs better than random would have an ROC curve that bends towards the top-left corner of the plot.

- The area under the ROC curve (AUC) is a summary statistic that represents the overall performance of the model. The AUC value shown in the plot is 0.89, indicating a relatively strong model performance.

This ROC curve visualization allows you to assess the model's performance at different operating points, helping you choose the appropriate decision threshold based on the desired balance between true positives and false positives. It is a valuable tool for evaluating and comparing the performance of binary classification models.

The clear and intuitive presentation of the ROC curve makes it an effective way to communicate model performance to stakeholders, data analysts, and machine learning researchers. By understanding the trade-offs between TPR and FPR, you can make informed decisions about model deployment and optimization.

## Confusion Matrix - Model Performance

|  | Did Not Survive | Survived |
|---|---|---|
| **Did Not Survive** | 93 | 12 |
| **Survived** | 20 | 54 |

*True Labels (rows) / Predicted Labels (columns)*

The above confusion matrix is interpreted as below:

- The rows represent the true labels - "Did Not Survive" and "Survived". These are the actual outcomes in the dataset.

- The columns represent the predicted labels - also "Did Not Survive" and "Survived". These are the model's predictions.

- The values in each cell represent the number of instances that fall into that particular combination of true and predicted labels.

- The top-left cell (93) indicates the number of instances that were correctly predicted as "Did Not Survive".

- The top-right cell (12) indicates the number of instances that were incorrectly predicted as "Survived" when the true label was "Did Not Survive" (false positives).

- The bottom-left cell (20) indicates the number of instances that were incorrectly predicted as "Did Not Survive" when the true label was "Survived" (false negatives).

- The bottom-right cell (54) indicates the number of instances that were correctly predicted as "Survived".

This confusion matrix provides a comprehensive view of the model's performance. It allows you to assess not just the overall accuracy, but also the model's ability to correctly identify both positive and negative instances. This information is crucial for understanding the trade-offs and making informed decisions about model deployment.

The clear presentation of the true and predicted labels in a tabular format makes the confusion matrix an effective tool for communicating model performance to stakeholders, data analysts, and machine learning researchers. By understanding the specific types of errors the model is making, you can identify areas for improvement and optimize the model accordingly.

**Results and Analysis**

- Training vs. Testing Accuracy: A comparison of training and testing accuracy highlights whether the model is overfitting or generalizing well.

- Feature Importance: The most significant features influencing the prediction are highlighted, offering interpretability.

- Misclassification Analysis: Instances where the model fails to predict correctly are analysed to identify trends (e.g., ambiguous data points, class imbalance)

**Applications**

Random Forest is widely used in a variety of domains due to its versatility and robustness. Some common applications include:

1. **Medical Diagnosis**:
   - Predicting diseases based on patient history and medical tests.
   - Analysing risk factors for specific conditions.

2. **Financial Analysis**:
   - Credit scoring and fraud detection.
   - Forecasting stock market trends.

3. **Marketing**:
   - Customer segmentation and churn prediction.
   - Recommender systems based on user behaviour.

4. **Environmental Studies**:
   - Predicting climate changes and ecological impacts.
   - Classifying land cover using satellite imagery.

5. **Natural Language Processing (NLP)**:
   - Sentiment analysis and text classification tasks.

6. **Image and Video Processing**:
   - Object detection and recognition tasks.
   - Predicting image categories in computer vision applications.

**Advantages of Random Forest**

1. **Robustness**:
   - Handles missing and imbalanced data well.
   - Resistant to overfitting compared to single decision trees.

2. **Feature Importance**:

o   Provides insights into which features are most influential in predictions.

3. **Scalability**:

    o   Efficient for large datasets with high dimensionality.

4. **Flexibility**:

    o   Works for both classification and regression tasks.

5. **Handles Non-linearity**:

    o   Capable of modelling complex non-linear relationships.

**Challenges and Limitations**

1. **Computationally Intensive**:

    o   Training multiple trees can require significant computational resources.

2. **Interpretability**:

    o   While individual decision trees are interpretable, the ensemble nature of Random Forest makes it harder to understand the overall model.

3. **Risk of Overfitting on Noisy Data**:

    o   Despite being robust, overfitting can still occur if not tuned properly.

**Conclusion**:

The Random Forest model implemented in this code effectively leverages ensemble learning to achieve high accuracy and robust predictions. By combining model evaluation metrics with insightful visualizations, the assignment demonstrates a solid understanding of machine learning principles. With further optimization of hyperparameters and interpretability enhancements, this implementation could be extended to more complex datasets and applications.