

# **Random Forest Report**

## **Introduction**

A Random Forest is an ensemble learning method that builds multiple decision trees and combines their outputs to improve accuracy and robustness. It operates by training each tree on a random subset of the data and aggregating their predictions. This report evaluates the performance of a Random Forest model applied to a given dataset.

## **Methodology**

A Random Forest follows these principles:

- It constructs multiple decision trees using bootstrapped samples of the dataset.
- Each tree selects a random subset of features to consider for splitting at each node.
- The final prediction is made by aggregating outputs from all individual trees (majority vote for classification, average for regression).

## **Steps Taken:**

### **1. Data Preprocessing:**

- Checked for missing values.
- Conducted exploratory data analysis (EDA) to examine feature distributions and correlations.
- Encoded categorical variables using one-hot encoding.
- Scaled numerical features if necessary.
- Split data into training (70%) and testing (30%) sets.

### **2. Model Training:**

- Implemented Random Forest using `sklearn.ensemble.RandomForestClassifier`.
- Experimented with different numbers of estimators (trees) and criteria (Gini and Entropy).
- Tuned hyperparameters such as `max_depth`, `min_samples_split`, and `max_features`.

### 3. Performance Evaluation:

- Measured using Accuracy, Precision, Recall, and F1-Score.
- Evaluated feature importance to identify the key contributors to predictions.

#### Random Forest Model Analysis

No	Model Name	criterion	max_depth	min_samples_leaf	max_features	min_samples_split	Results
01	Default Configuration	-	-	-	-	-	0.9403
02	Increased Depth	squared_error	20	2	sqrt	2	0.8030
03	More Trees	squared_error	None	2	log2	2	0.7991
04	Regularization	squared_error	10	2	auto	5	0.7653
05	Bootstrapping Disabled	squared_error	None	1	auto	2	0.8075
06	Custom Feature Selection	squared_error	None	1	0.5	2	0.8207
07	Increased Max Features	squared_error	None	4	auto	10	0.9161
08	Warm Start	squared_error	None	1	auto	2	0.8207

#### Conclusion

1. The analysis of various Random Forest model configurations reveals that the **Default Configuration** achieves the highest performance with a result of **0.9403**. This suggests that the default settings of the Random Forest Regressor are well-suited for the dataset used in this analysis. Other configurations, such as **Increased Depth**, **More Trees**, and **Regularization**, show lower performance, indicating that these modifications may lead to overfitting or underfitting in this specific context. The **Increased Max Features** configuration also performs relatively well, achieving a result of **0.9161**, which is close to the default configuration.

#### Recommendations

1. **Default Configuration:** Continue using the default configuration for the Random Forest Regressor, as it provides the best performance for the dataset.
2. **Increased Max Features:** Consider experimenting further with the max\_features parameter, as the configuration with increased max features showed promising results.

3. **Regularization:** If overfitting is a concern, the **Regularization** configuration (with max\_depth=10 and min\_samples\_split=5) could be explored further to balance model complexity and performance.
4. **Avoid Overcomplicating:** Avoid overly complex configurations (e.g., **Increased Depth** or **More Trees**) unless there is a specific need, as they may not improve performance and could increase computational cost.

### **Final Recommendation**

The **Default Configuration** of the Random Forest Regressor is the most effective for this dataset, achieving the highest result of **0.9403**. It is recommended to use this configuration for deployment or further analysis. If additional tuning is required, focus on the max\_features parameter, as it has shown potential for improvement without significantly compromising performance. Avoid unnecessary complexity in the model to maintain efficiency and generalizability.