

Introduction to Random Forest

→ a bunch of trees

1. Introduction

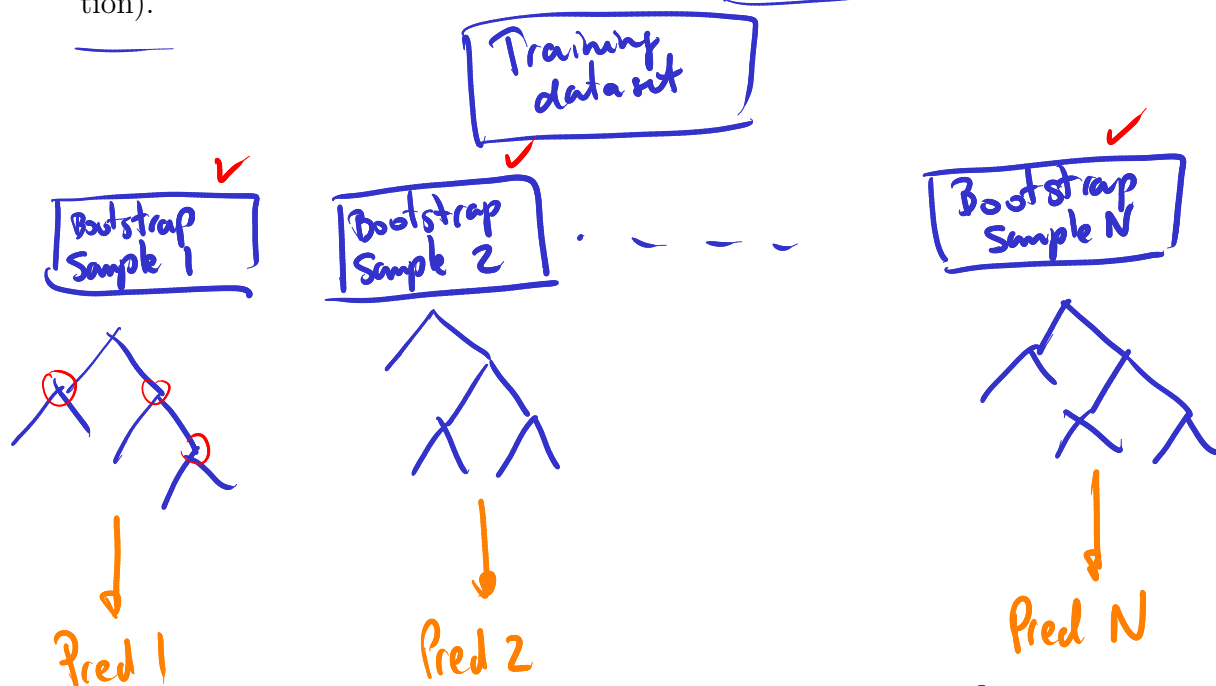
→ machine learning technique that combines model predictions

Random Forests are a type of ensemble learning method used for classification and regression. They operate by constructing multiple decision trees and outputting the mode of the classes (for classification) or mean prediction (for regression) of the individual trees.

2. How It Works

Random Forests combine bagging (bootstrap aggregation) with decision trees:

- Multiple samples are drawn with replacement from the training set. ✓
- A decision tree is trained on each sample. ✓
- At each node, a random subset of features is considered for splitting. ✓
- The final prediction is obtained by averaging (regression) or majority vote (classification).



$$\text{Random Forest Prediction} = \frac{\text{Pred 1} + \text{Pred 2} + \dots + \text{Pred N}}{N}$$

3. Applications in Data Science

Random Forests are widely used in:

- ✓ • **Credit Scoring:** Classifying customers as low/high credit risk.
- ✓ • **Medical Diagnosis:** Predicting disease outcomes.
- ✓ • **Customer Churn:** Identifying customers likely to leave.
- ✓ • **Image Classification:** Feature-rich image datasets.
- ✓ • **Feature Importance:** Identifying the most predictive variables.

4. Advantages

- Handles both classification and regression tasks. ✓
- Reduces overfitting compared to single decision trees. ✓
- Provides estimates of feature importance. ✓
- Robust to outliers and noise. ✓

5. Disadvantages

- Can be computationally expensive with large datasets. ✓
- Less interpretable than single decision trees. ✓
- May not perform well with sparse data or high-dimensionality without tuning.

6. Conclusion

Random Forests are an essential tool in the data scientist's toolkit due to their performance, flexibility, and ability to reduce overfitting. While they may lack interpretability, their predictive power often outweighs this drawback in many practical applications.