

Methods and Algorithms

Spring 2026

By the end of this lecture, you will be able to:

- ➊ Explain how decision trees partition feature space
- ➋ Implement Random Forests using bagging and feature randomization
- ➌ Interpret feature importance and out-of-bag error
- ➍ Apply ensemble methods to fraud detection problems

Finance Application: Fraud detection with interpretable feature importance

From single models to ensemble methods that combine many weak learners

Fraud Detection Challenge

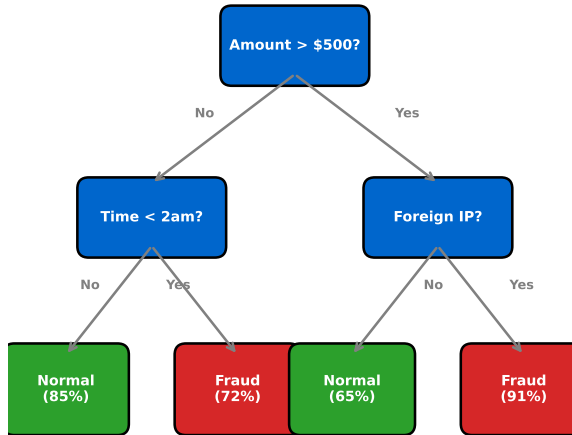
- Need high accuracy: fraudulent transactions cost millions
- Need interpretability: explain why transaction flagged
- Complex patterns: fraud evolves and adapts

Why Random Forests?

- Combines many trees for robust predictions
- Built-in feature importance ranking
- Handles non-linear relationships naturally

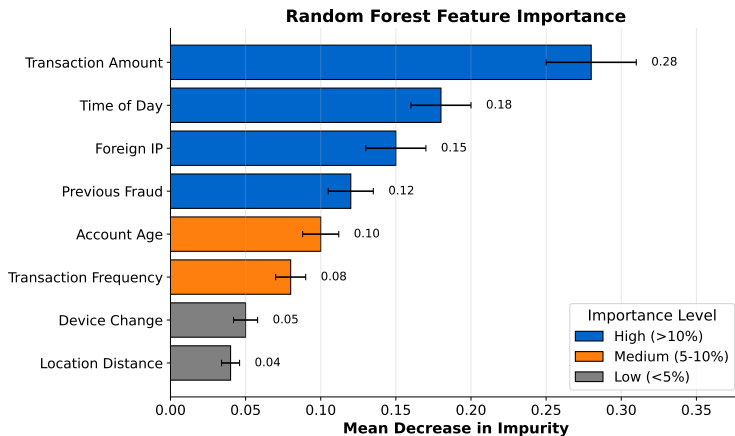
Ensemble methods: “wisdom of crowds” for machine learning

Decision Tree for Fraud Detection



https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L04_Random_Forests/01_decision_tree

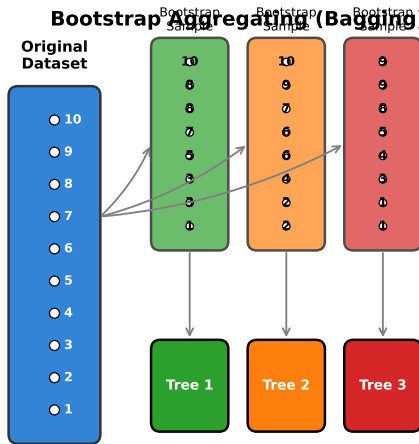
Trees split data using simple rules at each node until reaching a prediction



https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L04_Random_Forests/02_feature_importance

Random Forests automatically rank which features matter most for prediction

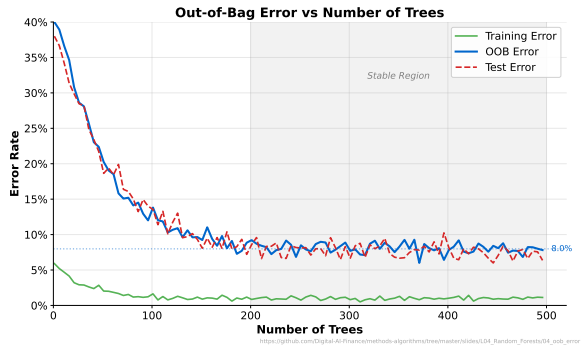
Bootstrap Aggregating (Bagging)



Each tree trained on ~63% unique samples (with replacement)

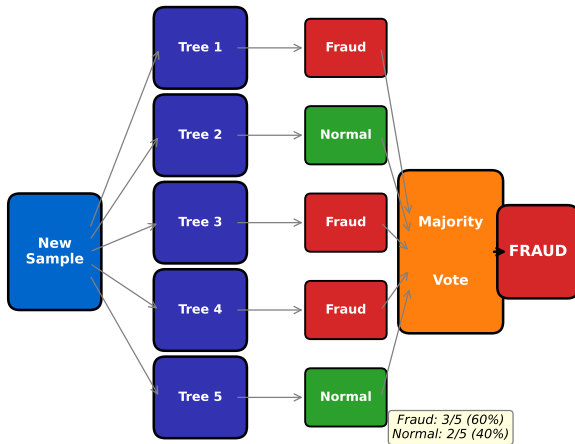
https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L04_Random_Forests/03_bootstrap

Each tree trains on a random sample, reducing overfitting



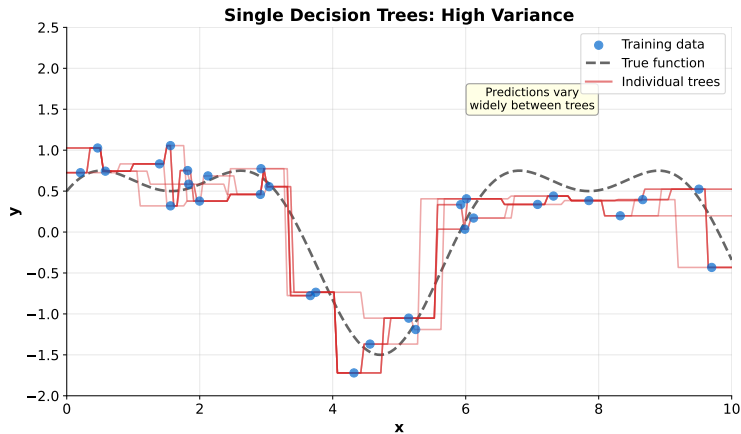
OOB error provides free cross-validation without held-out test set

Ensemble Voting (Classification)



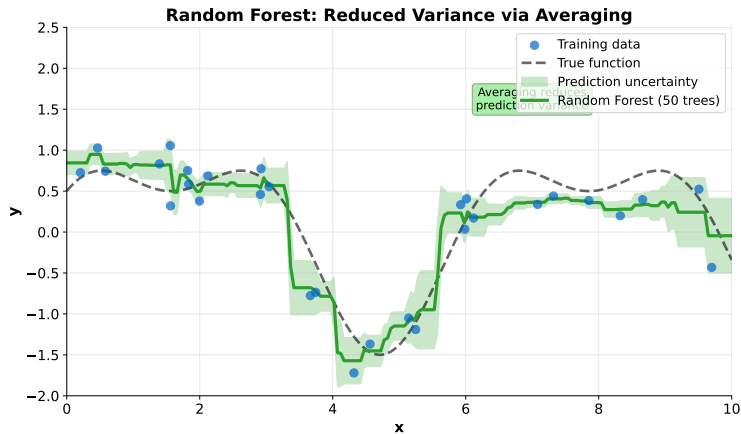
https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L04_Random_Forests/05_ensemble_voting

Final prediction combines votes from all trees (majority for classification)



https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L04_Random_Forests/06a_single_tree_variance

Each tree trained on different bootstrap sample produces different predictions



https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L04_Random_Forests/06b_random_forest_variance

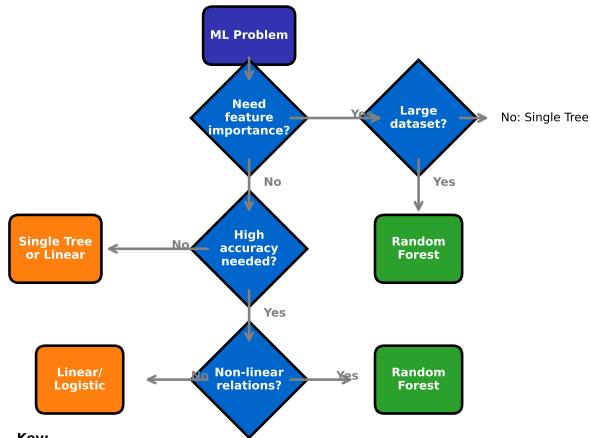
Averaging many high-variance trees produces low-variance ensemble

Open the Colab Notebook

- Exercise 1: Train a decision tree on credit data
- Exercise 2: Build a random forest and analyze feature importance
- Exercise 3: Tune hyperparameters with cross-validation

Link: <https://colab.research.google.com/> [TBD]

When to Use Random Forests



Key:

Random Forest: Best for accuracy + feature importance

Alternative: When interpretability is paramount

https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L04_Random_Forests/07_decision_flowchart

Random Forests excel when accuracy and feature importance both matter