

# Random Forests: Theory and Practice

September 10, 2025

## 1 Introduction

Random forests are ensemble models that aggregate many decision trees trained on bootstrap samples and with randomized feature selection. By averaging (regression) or majority voting (classification), they reduce variance, improve generalization, and retain robustness without extensive preprocessing.

## 2 Theory and Formulas

Random forests combine *bagging* and *feature subsampling*. Let  $\{T_b\}_{b=1}^B$  be trees trained on bootstrap samples  $\mathcal{D}_b$  and using a random subset of features at each split. For classification, the prediction is the majority vote

$$\hat{y} = \text{mode}(T_1(\mathbf{x}), \dots, T_B(\mathbf{x})), \quad (1)$$

and for regression, the average

$$\hat{y} = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x}). \quad (2)$$

Variance reduces as  $B$  increases when individual trees are not perfectly correlated (promoted by feature subsampling). A common choice is `max_features` =  $\sqrt{p}$  for classification with  $p$  features.

**Out-of-bag (OOB) estimation:** each bootstrap leaves out roughly  $\approx 36\%$  of samples. The model's OOB score uses these held-out samples as an internal validation, providing a nearly unbiased generalization estimate without extra cross-validation.

Key hyperparameters: number of trees  $B$  (`n_estimators`), split feature count (`max_features`), tree depth and leaf size (regularization), and bootstrap/OOB settings.

## 3 Applications and Tips

- **Pros:** strong out-of-the-box performance, handles mixed feature types, robust to outliers and monotone transforms, little scaling required.
- **Cons:** larger memory and inference cost than a single tree; less interpretable than a shallow tree.
- **Regularization:** tune `max_depth`, `min_samples_leaf`, `max_features`; increase `n_estimators` until OOB/test metrics stabilize.
- **Diagnostics:** use OOB score, learning curves vs trees, and check feature importances (and, when needed, permutation importance).

- **Preprocessing:** one-hot encode categorical features; no need to standardize numeric features.

## 4 Python Practice

Run the script in this chapter directory to generate figures into `figures/`.

Listing 1: Generate Random Forest figures

```
1 python gen_random_forest_figures.py
```

Listing 2: `gen_random_forest_figures.py`

```
1 """
2 Figure generator for the Random Forest chapter.
3
4 Generates illustrative figures and saves them into the chapter's 'figures/'
5 folder next to this script, regardless of current working directory.
6
7 Requirements:
8 - Python 3.8+
9 - numpy, matplotlib, scikit-learn
10
11 Install (if needed):
12     pip install numpy matplotlib scikit-learn
13
14 This script avoids newer or experimental APIs for broader compatibility.
15 """
16 from __future__ import annotations
17
18 import os
19 import numpy as np
20 import matplotlib.pyplot as plt
21 from matplotlib.colors import ListedColormap
22
23 try:
24     from sklearn.datasets import make_moons, make_classification
25     from sklearn.ensemble import RandomForestClassifier
26 except Exception:
27     raise SystemExit(
28         "Missing scikit-learn. Please install with: pip install scikit-learn"
29     )
30
31
32 def _ensure_figures_dir(path: str | None = None) -> str:
33     """Create figures directory under this chapter regardless of CWD."""
34     if path is None:
35         base = os.path.dirname(os.path.abspath(__file__))
36         path = os.path.join(base, "figures")
37     os.makedirs(path, exist_ok=True)
38     return path
39
40
41 def _plot_decision_boundary(ax, clf, X, y, title: str):
```

```

42 x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
43 y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5
44 xx, yy = np.meshgrid(
45     np.linspace(x_min, x_max, 400), np.linspace(y_min, y_max, 400)
46 )
47 Z = clf.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
48 cmap_light = ListedColormap(["#FFEEEE", "#EEEEFF"])
49 cmap_bold = ListedColormap(["#E74C3C", "#3498DB"])
50 ax.contourf(xx, yy, Z, cmap=cmap_light, alpha=0.8, levels=np.unique(Z).
    size)
51 ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold, edgecolors="k", s=20)
52 ax.set_title(title)
53 ax.set_xlabel("Feature 1")
54 ax.set_ylabel("Feature 2")
55
56
57 def fig_rf_decision_boundary_2class(out_dir: str) -> str:
58     np.random.seed(0)
59     X, y = make_moons(n_samples=500, noise=0.3, random_state=0)
60     clf = RandomForestClassifier(
61         n_estimators=150, max_depth=None, random_state=0
62     )
63     clf.fit(X, y)
64
65     fig, ax = plt.subplots(figsize=(5.5, 4.5), dpi=150)
66     _plot_decision_boundary(ax, clf, X, y, "Random Forest boundary (
        n_estimators=150)")
67     out_path = os.path.join(out_dir, "rf_decision_boundary_2class.png")
68     fig.tight_layout()
69     fig.savefig(out_path)
70     plt.close(fig)
71     return out_path
72
73
74 def fig_rf_n_estimators_compare(out_dir: str) -> str:
75     np.random.seed(1)
76     X, y = make_moons(n_samples=600, noise=0.28, random_state=1)
77     models = [
78         (RandomForestClassifier(n_estimators=5, random_state=1), "n_estimators
            =5"),
79         (RandomForestClassifier(n_estimators=200, random_state=1), "
            n_estimators=200"),
80     ]
81     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=True,
        sharey=True)
82     for ax, (m, title) in zip(axes, models):
83         m.fit(X, y)
84         _plot_decision_boundary(ax, m, X, y, f"Random Forest: {title}")
85     fig.suptitle("Effect of number of trees")
86     out_path = os.path.join(out_dir, "rf_n_estimators_compare.png")
87     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
88     fig.savefig(out_path)
89     plt.close(fig)
90     return out_path

```

```

91
92
93 def fig_rf_max_features_compare(out_dir: str) -> str:
94     np.random.seed(2)
95     X, y = make_moons(n_samples=600, noise=0.32, random_state=2)
96     models = [
97         (RandomForestClassifier(max_features=1, n_estimators=150, random_state
98             =2), "max_features=1"),
99         (RandomForestClassifier(max_features="sqrt", n_estimators=150,
100             random_state=2), "max_features=sqrt"),
101     ]
102     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=True,
103         sharey=True)
104     for ax, (m, title) in zip(axes, models):
105         m.fit(X, y)
106         _plot_decision_boundary(ax, m, X, y, f"Random Forest: {title}")
107     fig.suptitle("Effect of feature subsampling (max_features)")
108     out_path = os.path.join(out_dir, "rf_max_features_compare.png")
109     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
110     fig.savefig(out_path)
111     plt.close(fig)
112     return out_path
113
114 def fig_rf_feature_importances(out_dir: str) -> str:
115     X, y = make_classification(
116         n_samples=800,
117         n_features=10,
118         n_informative=4,
119         n_redundant=3,
120         n_repeated=0,
121         random_state=7,
122         shuffle=True,
123     )
124     clf = RandomForestClassifier(n_estimators=200, random_state=7)
125     clf.fit(X, y)
126     importances = clf.feature_importances_
127
128     fig, ax = plt.subplots(figsize=(7.0, 4.0), dpi=160)
129     idx = np.arange(importances.size)
130     ax.bar(idx, importances, color="#2ECC71")
131     ax.set_xticks(idx)
132     ax.set_xticklabels([f"f{i}" for i in idx])
133     ax.set_ylabel("importance")
134     ax.set_title("Random Forest feature importances")
135     ax.set_ylim(0, max(0.25, importances.max() + 0.05))
136     for i, v in enumerate(importances):
137         ax.text(i, v + 0.01, f"{v:.2f}", ha="center", va="bottom", fontsize=8)
138     out_path = os.path.join(out_dir, "rf_feature_importances.png")
139     fig.tight_layout()
140     fig.savefig(out_path)
141     plt.close(fig)
142     return out_path

```

```

142
143 def fig_rf_oob_curve(out_dir: str) -> str:
144     np.random.seed(3)
145     X, y = make_classification(
146         n_samples=1200,
147         n_features=15,
148         n_informative=5,
149         n_redundant=5,
150         random_state=3,
151     )
152
153     trees = np.unique(np.linspace(5, 300, 15).astype(int))
154     oob_scores = []
155     for n in trees:
156         # OOB requires bootstrap=True
157         rf = RandomForestClassifier(
158             n_estimators=n, oob_score=True, bootstrap=True, random_state=3
159         )
160         rf.fit(X, y)
161         oob_scores.append(rf.oob_score_)
162
163     fig, ax = plt.subplots(figsize=(6.5, 4.0), dpi=160)
164     ax.plot(trees, oob_scores, marker="o", color="#9B59B6")
165     ax.set_xlabel("n_estimators")
166     ax.set_ylabel("OOB score")
167     ax.set_title("Out-of-bag score vs number of trees")
168     ax.grid(True, linestyle=":", alpha=0.4)
169     out_path = os.path.join(out_dir, "rf_oob_curve.png")
170     fig.tight_layout()
171     fig.savefig(out_path)
172     plt.close(fig)
173     return out_path
174
175
176 def main():
177     out_dir = _ensure_figures_dir(None)
178     generators = [
179         fig_rf_decision_boundary_2class,
180         fig_rf_n_estimators_compare,
181         fig_rf_max_features_compare,
182         fig_rf_feature_importances,
183         fig_rf_oob_curve,
184     ]
185     print("Generating figures into:", os.path.abspath(out_dir))
186     for gen in generators:
187         try:
188             p = gen(out_dir)
189             print("Saved:", p)
190         except Exception as e:
191             print("Failed generating", gen.__name__, ":", e)
192
193
194 if __name__ == "__main__":
195     main()

```

## 5 Result

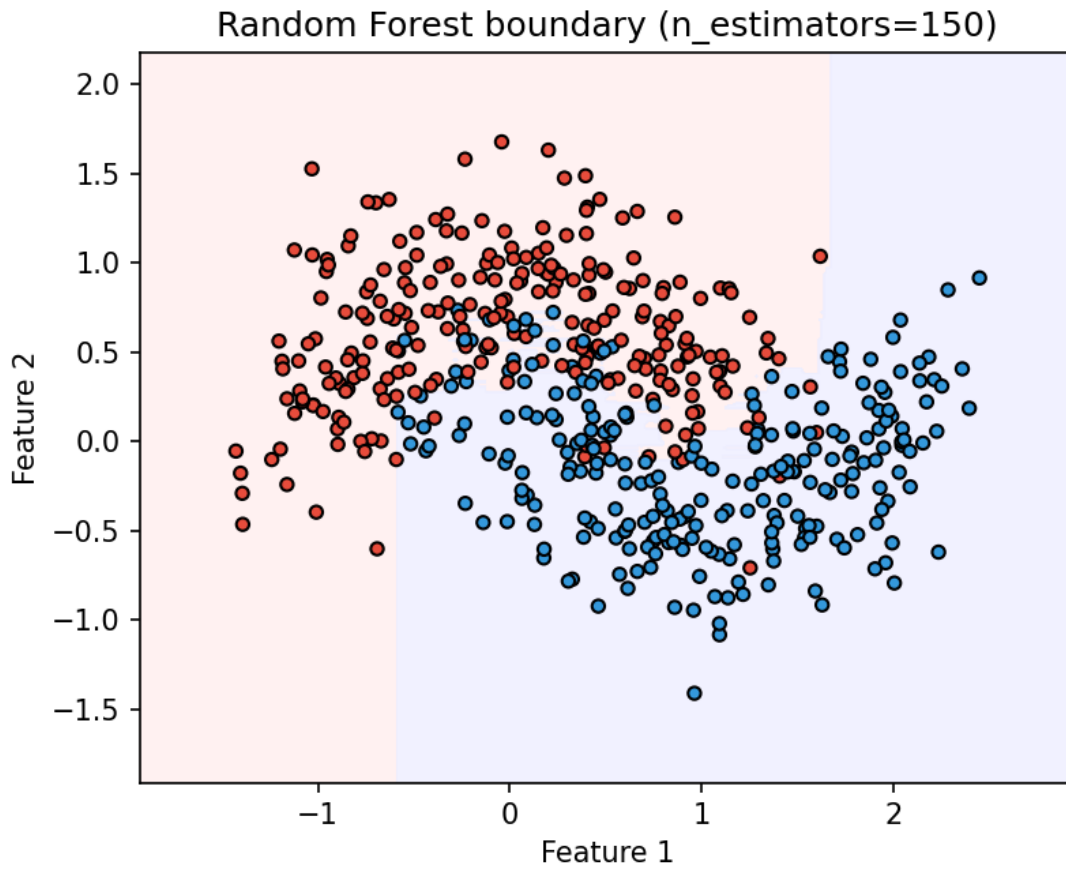


Figure 1: Random forest decision boundary on a 2-class dataset.

Effect of number of trees

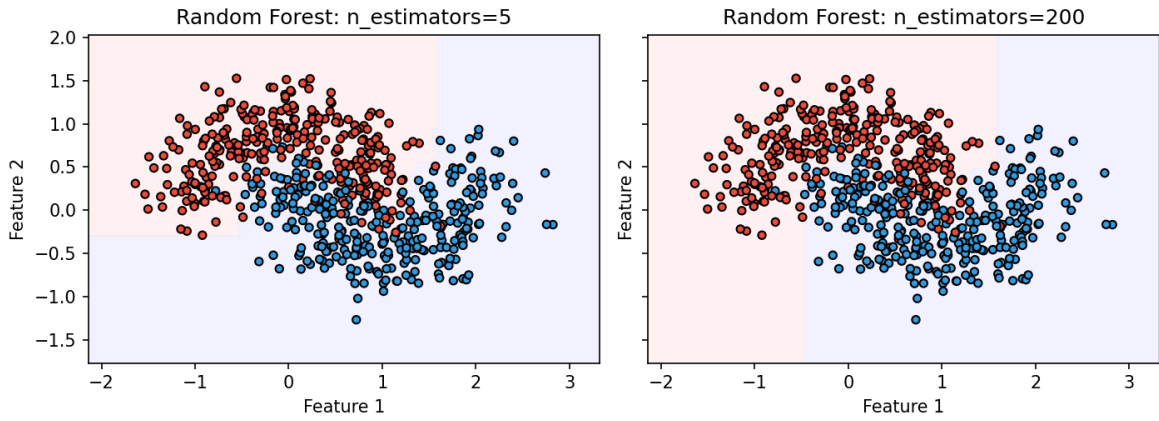


Figure 2: Effect of number of trees: small vs large ensemble.

Effect of feature subsampling (max\_features)

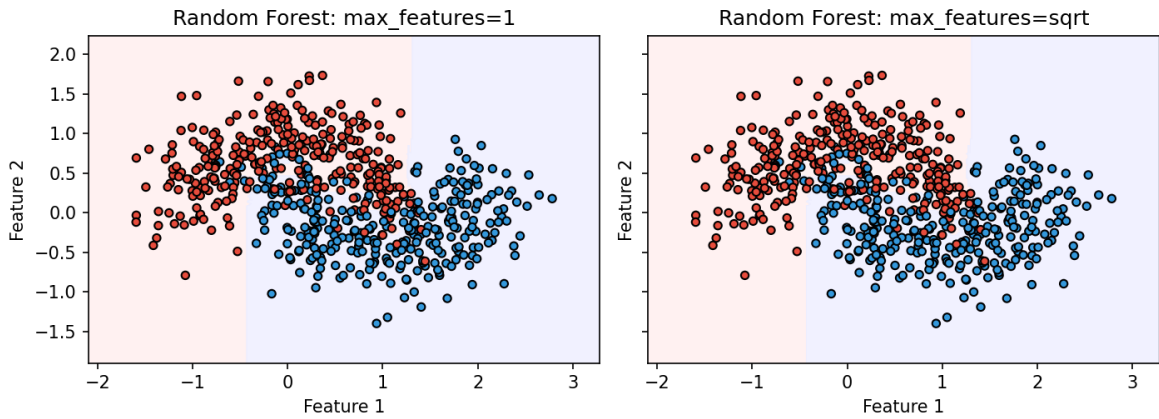


Figure 3: Decision boundaries with different max\_features.

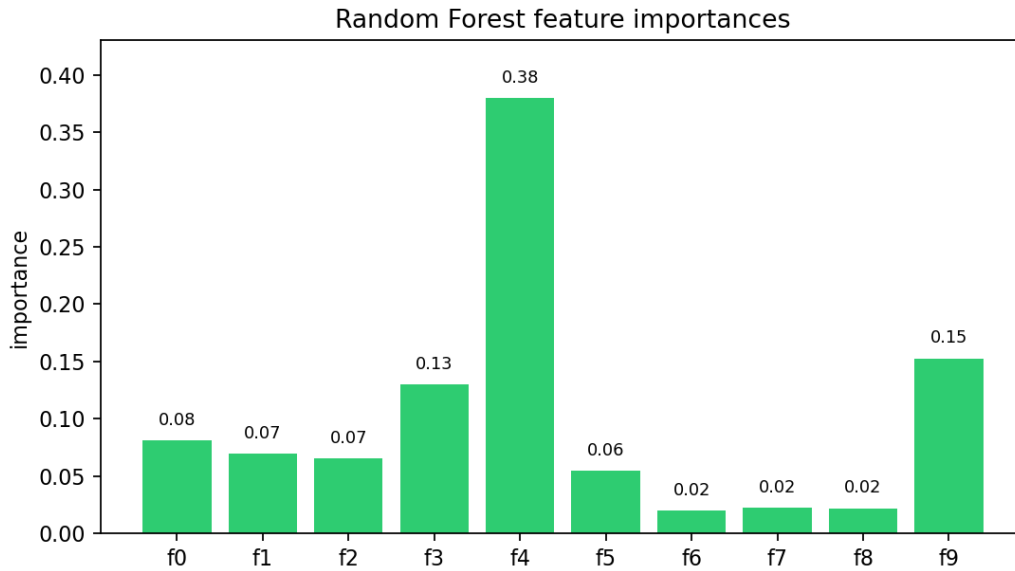


Figure 4: Feature importances from a random forest.

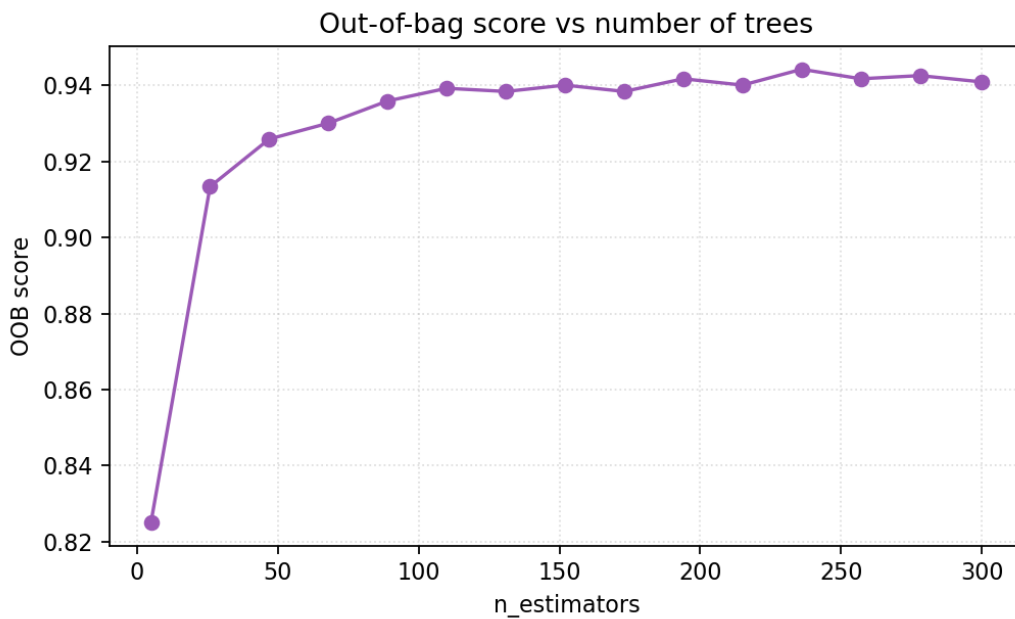


Figure 5: Out-of-bag score vs number of trees.

## 6 Summary

Random forests provide a strong, reliable baseline for many tasks. They reduce variance through bagging and feature subsampling, offer OOB validation, and yield useful importance measures. Tune tree count and regularization for the best trade-off between accuracy and efficiency.