

随机森林：理论与实践

2025 年 9 月 10 日

1 引言

随机森林（Random Forest）是基于“装袋”（bagging）与“特征子采样”的集成学习方法。它通过对多个随机化的决策树进行投票（分类）或求均值（回归），显著降低方差、提升泛化能力，同时对数据预处理的要求较低。

2 原理与公式

随机森林的两大关键：自助采样（bootstrap）与特征子采样（每次划分仅在随机抽取的特征子集上搜索最优划分）。设有 B 棵树 $\{T_b\}_{b=1}^B$ ，每棵树在各自的自助样本 \mathcal{D}_b 上训练。分类任务的集成预测为多数表决：

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

回归任务的集成预测为简单平均：

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

由于特征子采样降低了各基学习器的相关性，随着树数 B 的增加，整体方差可进一步下降。分类中常用的经验选择是 `max_features`= \sqrt{p} （ p 为特征数）。

袋外（OOB）估计：每次自助采样平均会留下约 $\approx 36\%$ 的“袋外”样本。利用袋外样本做内部验证可得到接近无偏的泛化评估，无需额外划分验证集或交叉验证。

常用超参数包括：树的数量 `B(n_estimators)`、每次划分考虑的特征数（`max_features`）、树深与叶子最小样本数（正则化）、以及 `bootstrap/oob_score` 设置等。

3 应用与技巧

- **优点：**默认表现强、鲁棒性好、对特征尺度不敏感、能处理混合类型特征。

- **缺点：**预测与存储开销较单棵树更大；可解释性低于浅层树。
- **正则化与调参：**调整 `max_depth`、`min_samples_leaf`、`max_features`；增大 `n_estimators` 直至 OOB/测试指标收敛。
- **诊断：**使用 OOB 分数、随树数变化的学习曲线；查看特征重要性（必要时用置换重要性）。
- **预处理：**类别特征需独热编码；数值特征一般无需标准化。

4 Python 实战

在本章节目录运行下述命令，图片将保存到 `figures/`：

Listing 1: 生成随机森林配图

```
1 python gen_random_forest_figures.py
```

Listing 2: `gen_random_forest_figures.py` 源码

```
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 """
17
18 from __future__ import annotations
19
20 import os
21 import numpy as np
22 import matplotlib.pyplot as plt
23 from matplotlib.colors import ListedColormap
```

```

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-
29         learn"
30     )
31
32 def _ensure_figures_dir(path: str | None = None) -> str:
33     """Create figures directory under this chapter regardless of CWD.
34     """
35     if path is None:
36         base = os.path.dirname(os.path.abspath(__file__))
37         path = os.path.join(base, "figures")
38     os.makedirs(path, exist_ok=True)
39     return path
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
51                 (Z).size)
52     ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold, edgecolors="k", s
53               =20)
54     ax.set_title(title)
55     ax.set_xlabel("Feature 1")
56     ax.set_ylabel("Feature 2")
57
58 def fig_rf_decision_boundary_2class(out_dir: str) -> str:
59     np.random.seed(0)
60     X, y = make_moons(n_samples=500, noise=0.3, random_state=0)
61     clf = RandomForestClassifier(
62         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=2), "max_features=1"),
98         (RandomForestClassifier(max_features="sqrt", n_estimators=150,
            random_state=2), "max_features=sqrt"),
```

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

```
140     return out_path
141
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,
159             random_state=3
160         )
161         rf.fit(X, y)
162         oob_scores.append(rf.oob_score_)
163
164     fig, ax = plt.subplots(figsize=(6.5, 4.0), dpi=160)
165     ax.plot(trees, oob_scores, marker="o", color="#9B59B6")
166     ax.set_xlabel("n_estimators")
167     ax.set_ylabel("OOB score")
168     ax.set_title("Out-of-bag score vs number of trees")
169     ax.grid(True, linestyle=":", alpha=0.4)
170     out_path = os.path.join(out_dir, "rf_oob_curve.png")
171     fig.tight_layout()
172     fig.savefig(out_path)
173     plt.close(fig)
174     return out_path
175
176
177 def main():
178     out_dir = _ensure_figures_dir(None)
179     generators = [
180         fig_rf_decision_boundary_2class,
181         fig_rf_n_estimators_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 结果

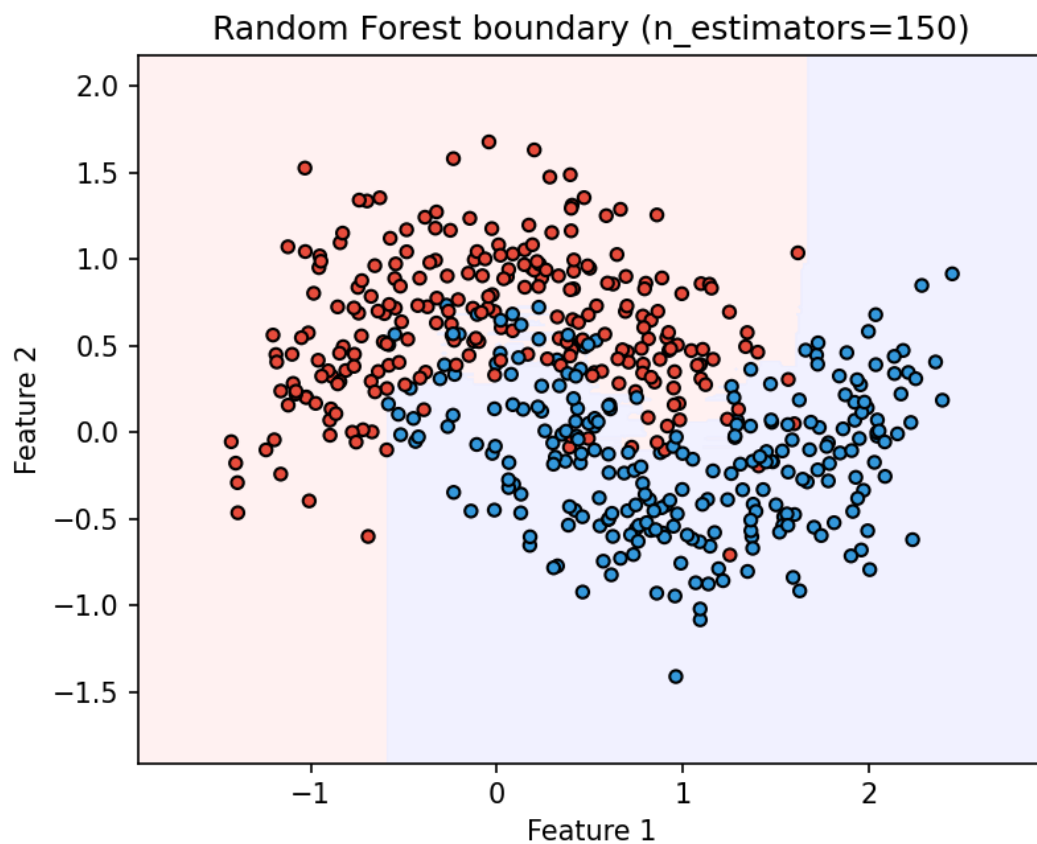


图 1: 随机森林在两类数据上的决策边界。

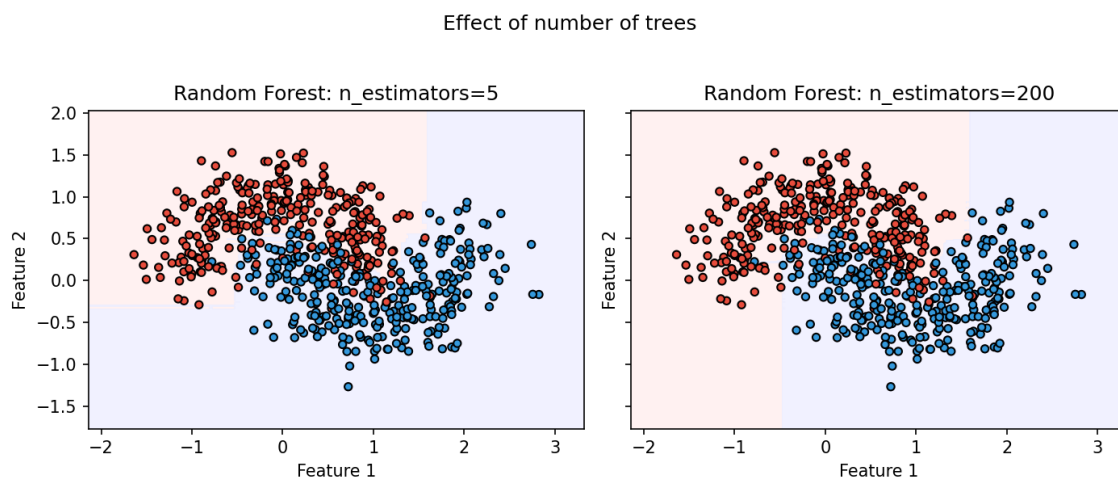


图 2: 树数量影响: 少树 vs 多树的对比。

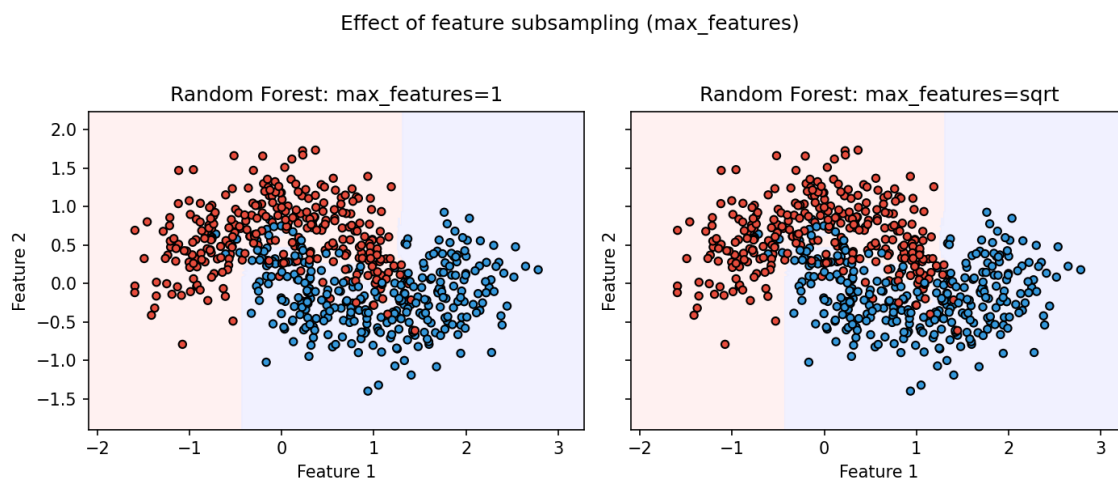


图 3: 不同 $max_features$ 下的决策边界对比。

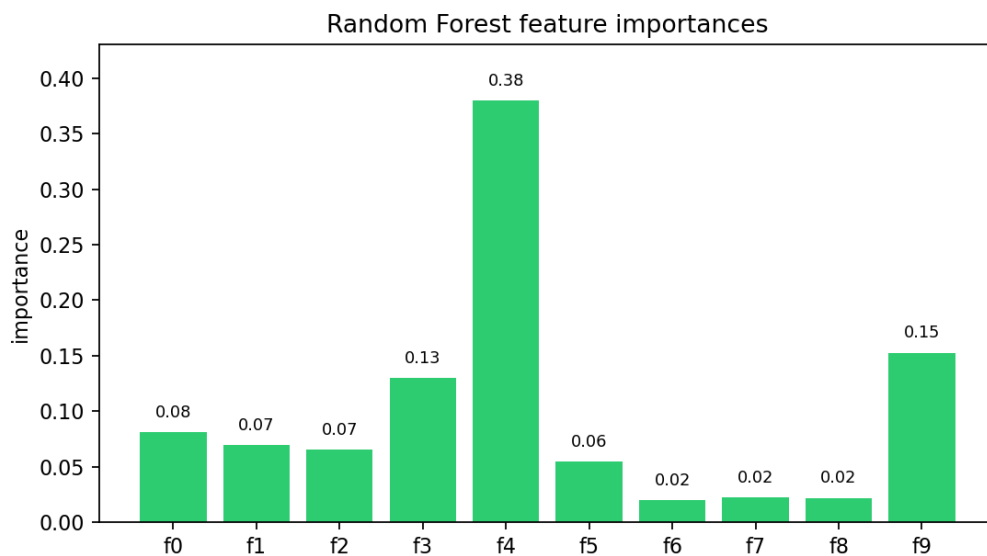


图 4: 随机森林的特征重要性可视化。

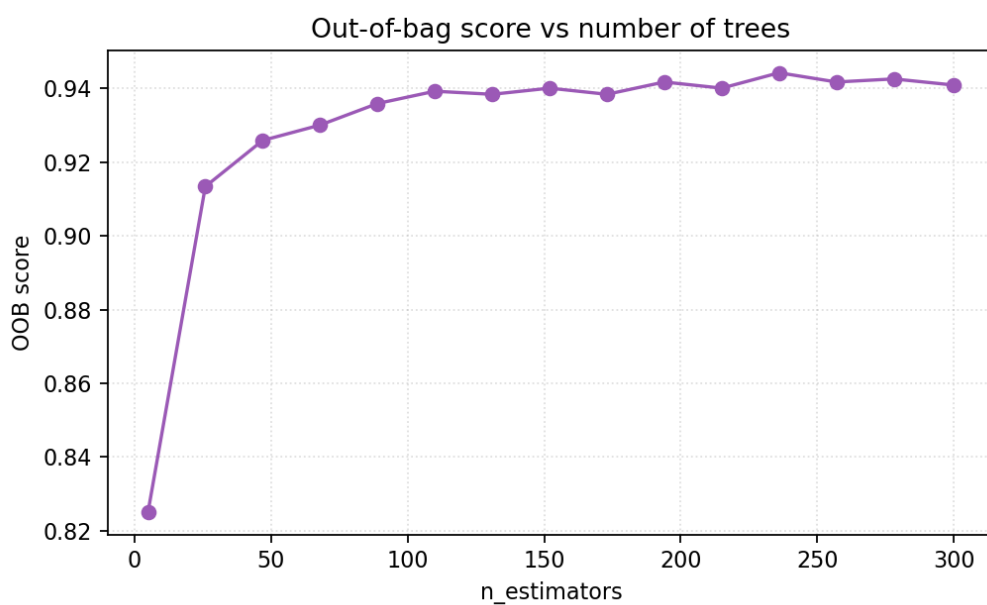


图 5: 袋外（OOB）分数随树数量变化的曲线。

6 总结

随机森林通过“装袋 + 特征子采样”有效降低方差，提供 OOB 内部验证与可用的特征重要性，是工程中可靠的通用基线模型。通过合适的树数与正则化参数，可在准确率与效率之间取得良好平衡。