

XGBoost：理论与实践

2025 年 9 月 10 日

1 引言

XGBoost 是高效、可扩展的梯度提升树（GBDT）实现。它通过二阶近似优化、稀疏感知的分裂查找、学习率（shrinkage）与行/列子采样、以及显式结构正则化等机制，兼顾训练效率与泛化性能。

2 原理与公式

梯度提升以逐步加法模型 $F_M(\mathbf{x}) = \sum_{m=1}^M f_m(\mathbf{x})$ 拟合弱学习器（浅树），XGBoost 最小化带正则的目标：

$$\mathcal{L} = \sum_{i=1}^n \ell(y_i, \hat{y}_i) + \sum_{m=1}^M \Omega(f_m), \quad \Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2, \quad (1)$$

其中 T 为叶子数、 w 为叶子打分。对损失在当前预测处做二阶泰勒展开，得到节点上梯度 g_i 与海森 h_i 的求和，左右子集 L, R 的分裂增益为：

$$\text{Gain} = \frac{1}{2} \left(\frac{(\sum_{i \in L} g_i)^2}{\sum_{i \in L} h_i + \lambda} + \frac{(\sum_{i \in R} g_i)^2}{\sum_{i \in R} h_i + \lambda} - \frac{(\sum_{i \in L \cup R} g_i)^2}{\sum_{i \in L \cup R} h_i + \lambda} \right) - \gamma. \quad (2)$$

通过 λ, γ 的结构正则、学习率（shrinkage）、列/行子采样、深度与叶子约束来控制模型复杂度与过拟合。

3 应用与技巧

- **特征与缩放**：支持数值与独热编码后的类别特征；树模型通常无需标准化。
- **关键超参**：`n_estimators`、`max_depth`、`learning_rate`、`subsample`、`colsample_bytree`、`reg_alpha`/`reg_lambda`。
- **早停**：使用 `eval_set` 与 `early_stopping_rounds` 在验证集上早停。

- 类不平衡：设置 `scale_pos_weight` 或采用分层采样。
- 解释性：内置重要性可作初筛；更稳健可用置换重要性或 SHAP。

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 XGBoost 配图

```
1 python gen_xgboost_figures.py
```

Listing 2: `gen_xgboost_figures.py` 源码

```
1 """
2 Figure generator for the XGBoost 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 - xgboost (optional; falls back to scikit-learn GradientBoosting if
   missing)
11
12 Install (if needed):
13     pip install numpy matplotlib scikit-learn xgboost
14
15 This script avoids newer or experimental APIs for broader compatibility
   .
16 """
17 from __future__ import annotations
18
19 import os
20 import numpy as np
21 import matplotlib.pyplot as plt
22 from matplotlib.colors import ListedColormap
23
24 try:
25     import xgboost as xgb # type: ignore
26     HAS_XGB = True
27 except Exception:
```

```
28     xgb = None
29     HAS_XGB = False
30
31 from sklearn.datasets import make_moons, make_classification
32 from sklearn.model_selection import train_test_split
33 from sklearn.metrics import log_loss
34
35 try:
36     from sklearn.ensemble import GradientBoostingClassifier
37 except Exception:
38     GradientBoostingClassifier = None # type: ignore
39
40
41 def _ensure_figures_dir(path: str | None = None) -> str:
42     """Create figures directory under this chapter regardless of CWD.
43     """
44     if path is None:
45         base = os.path.dirname(os.path.abspath(__file__))
46         path = os.path.join(base, "figures")
47     os.makedirs(path, exist_ok=True)
48     return path
49
50 def _plot_decision_boundary(ax, clf, X, y, title: str):
51     x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
52     y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5
53     xx, yy = np.meshgrid(
54         np.linspace(x_min, x_max, 400), np.linspace(y_min, y_max, 400)
55     )
56     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
57     cmap_light = ListedColormap(["#FFEEEE", "#EEEEFF"])
58     cmap_bold = ListedColormap(["#E74C3C", "#3498DB"])
59     ax.contourf(xx, yy, Z, cmap=cmap_light, alpha=0.8, levels=np.unique(
60         Z).size)
61     ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold, edgecolors="k", s
62               =20)
63     ax.set_title(title)
64     ax.set_xlabel("Feature 1")
65     ax.set_ylabel("Feature 2")
66
67 def _xgb_classifier(**kwargs):
68     if HAS_XGB:
```

```

68     params = dict(
69         n_estimators=200,
70         max_depth=3,
71         learning_rate=0.1,
72         subsample=1.0,
73         colsample_bytree=1.0,
74         objective="binary:logistic",
75         tree_method="hist",
76         random_state=0,
77         n_jobs=0,
78     )
79     params.update(kwargs)
80     return xgb.XGBClassifier(**params)
81 else:
82     if GradientBoostingClassifier is None:
83         raise RuntimeError("Neither xgboost nor
84                             GradientBoostingClassifier available.")
85     params = dict(
86         n_estimators=kwargs.get("n_estimators", 200),
87         max_depth=kwargs.get("max_depth", 3),
88         learning_rate=kwargs.get("learning_rate", 0.1),
89         random_state=0,
90     )
91     return GradientBoostingClassifier(**params)
92
93 def fig_xgb_decision_boundary_2class(out_dir: str) -> str:
94     np.random.seed(0)
95     X, y = make_moons(n_samples=500, noise=0.3, random_state=0)
96     clf = _xgb_classifier()
97     clf.fit(X, y)
98
99     fig, ax = plt.subplots(figsize=(5.5, 4.5), dpi=150)
100    title = "XGBoost boundary (max_depth=3, lr=0.1)" if HAS_XGB else "
101            GBDT boundary (fallback)"
102    _plot_decision_boundary(ax, clf, X, y, title)
103    out_path = os.path.join(out_dir, "xgb_decision_boundary_2class.png"
104                            )
105    fig.tight_layout()
106    fig.savefig(out_path)
107    plt.close(fig)
108    return out_path

```

```
108
109 def fig_xgb_learning_rate_compare(out_dir: str) -> str:
110     np.random.seed(1)
111     X, y = make_moons(n_samples=550, noise=0.3, random_state=1)
112     models = [
113         (_xgb_classifier(learning_rate=0.05, n_estimators=400), "
114             learning_rate=0.05, n_estimators=400"),
115         (_xgb_classifier(learning_rate=0.3, n_estimators=150), "
116             learning_rate=0.3, n_estimators=150"),
117     ]
118     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=
119         True, sharey=True)
120     for ax, (m, title) in zip(axes, models):
121         m.fit(X, y)
122         label = ("XGBoost: " if HAS_XGB else "GBDT: ") + title
123         _plot_decision_boundary(ax, m, X, y, label)
124     fig.suptitle("Effect of learning_rate with trees budget")
125     out_path = os.path.join(out_dir, "xgb_learning_rate_compare.png")
126     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
127     fig.savefig(out_path)
128     plt.close(fig)
129     return out_path
130
131
132 def fig_xgb_max_depth_compare(out_dir: str) -> str:
133     np.random.seed(2)
134     X, y = make_moons(n_samples=600, noise=0.32, random_state=2)
135     models = [
136         (_xgb_classifier(max_depth=2, n_estimators=250), "max_depth=2"),
137         ,
138         (_xgb_classifier(max_depth=4, n_estimators=250), "max_depth=4"),
139         ,
140         (_xgb_classifier(max_depth=8, n_estimators=250), "max_depth=8"),
141         ,
142     ]
143     fig, axes = plt.subplots(1, 3, figsize=(12.5, 4.2), dpi=150, sharex
144         =True, sharey=True)
145     for ax, (m, title) in zip(axes, models):
146         m.fit(X, y)
147         label = ("XGBoost: " if HAS_XGB else "GBDT: ") + title
148         _plot_decision_boundary(ax, m, X, y, label)
149     fig.suptitle("Effect of max_depth")
150     out_path = os.path.join(out_dir, "xgb_max_depth_compare.png")
```

```
144     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
145     fig.savefig(out_path)
146     plt.close(fig)
147     return out_path
148
149
150 def fig_xgb_feature_importances(out_dir: str) -> str:
151     X, y = make_classification(
152         n_samples=800,
153         n_features=10,
154         n_informative=4,
155         n_redundant=3,
156         n_repeated=0,
157         random_state=7,
158         shuffle=True,
159     )
160     clf = _xgb_classifier(n_estimators=300, max_depth=4, learning_rate
161                          =0.1)
162     clf.fit(X, y)
163     importances = getattr(clf, "feature_importances_", None)
164     if importances is None:
165         # Fallback: uniform zeros to avoid crash
166         importances = np.zeros(X.shape[1], dtype=float)
167
168     fig, ax = plt.subplots(figsize=(7.0, 4.0), dpi=160)
169     idx = np.arange(importances.size)
170     ax.bar(idx, importances, color="#F39C12")
171     ax.set_xticks(idx)
172     ax.set_xticklabels([f"f{i}" for i in idx])
173     ax.set_ylabel("importance")
174     title = "XGBoost feature importances" if HAS_XGB else "GBDT feature
175             importances"
176     ax.set_title(title)
177     ax.set_ylim(0, max(0.25, float(importances.max()) + 0.05))
178     for i, v in enumerate(importances):
179         ax.text(i, v + 0.01, f"{v:.2f}", ha="center", va="bottom",
180                fontsize=8)
181     out_path = os.path.join(out_dir, "xgb_feature_importances.png")
182     fig.tight_layout()
183     fig.savefig(out_path)
184     plt.close(fig)
185     return out_path
```

```

184
185 def fig_xgb_eval_logloss_curve(out_dir: str) -> str:
186     np.random.seed(3)
187     X, y = make_classification(
188         n_samples=1200,
189         n_features=15,
190         n_informative=5,
191         n_redundant=5,
192         random_state=3,
193     )
194     X_train, X_val, y_train, y_val = train_test_split(X, y, test_size
195         =0.3, random_state=3)
196
197     if HAS_XGB:
198         clf = _xgb_classifier(n_estimators=300, learning_rate=0.1,
199             max_depth=4)
200         clf.fit(X_train, y_train, eval_set=[(X_train, y_train), (X_val,
201             y_val)], eval_metric="logloss", verbose=False)
202         res = clf.evals_result()
203         tr = np.array(res.get("validation_0", {}).get("logloss", []),
204             dtype=float)
205         va = np.array(res.get("validation_1", {}).get("logloss", []),
206             dtype=float)
207     else:
208         # Fallback using staged decision on GradientBoosting
209         clf = _xgb_classifier(n_estimators=300, learning_rate=0.1,
210             max_depth=3)
211         clf.fit(X_train, y_train)
212         tr_list, va_list = [], []
213         # GradientBoostingClassifier provides staged_predict_proba
214         if hasattr(clf, "staged_predict_proba"):
215             for y_tr_prob, y_va_prob in zip(clf.staged_predict_proba(
216                 X_train), clf.staged_predict_proba(X_val)):
217                 tr_list.append(log_loss(y_train, y_tr_prob))
218                 va_list.append(log_loss(y_val, y_va_prob))
219         else:
220             # Last resort: single-point curves
221             y_tr_prob = clf.predict_proba(X_train)
222             y_va_prob = clf.predict_proba(X_val)
223             tr_list = [log_loss(y_train, y_tr_prob)]
224             va_list = [log_loss(y_val, y_va_prob)]
225     tr, va = np.array(tr_list), np.array(va_list)

```

```
220     fig, ax = plt.subplots(figsize=(6.8, 4.2), dpi=160)
221     ax.plot(np.arange(1, len(tr) + 1), tr, label="train logloss", color
222            = "#2ECC71")
223     ax.plot(np.arange(1, len(va) + 1), va, label="valid logloss", color
224            = "#E74C3C")
225     ax.set_xlabel("n_trees")
226     ax.set_ylabel("logloss")
227     ax.set_title("Evaluation curve (logloss vs trees)")
228     ax.legend()
229     ax.grid(True, linestyle=":", alpha=0.4)
230     out_path = os.path.join(out_dir, "xgb_eval_logloss_curve.png")
231     fig.tight_layout()
232     fig.savefig(out_path)
233     plt.close(fig)
234     return out_path
235
236 def main():
237     out_dir = _ensure_figures_dir(None)
238     generators = [
239         fig_xgb_decision_boundary_2class,
240         fig_xgb_learning_rate_compare,
241         fig_xgb_max_depth_compare,
242         fig_xgb_feature_importances,
243         fig_xgb_eval_logloss_curve,
244     ]
245     print("Generating figures into:", os.path.abspath(out_dir))
246     if not HAS_XGB:
247         print("xgboost not found; falling back to
248               GradientBoostingClassifier where possible.")
249     for gen in generators:
250         try:
251             p = gen(out_dir)
252             print("Saved:", p)
253         except Exception as e:
254             print("Failed generating", gen.__name__, ":", e)
255
256 if __name__ == "__main__":
257     main()
```


5 结果

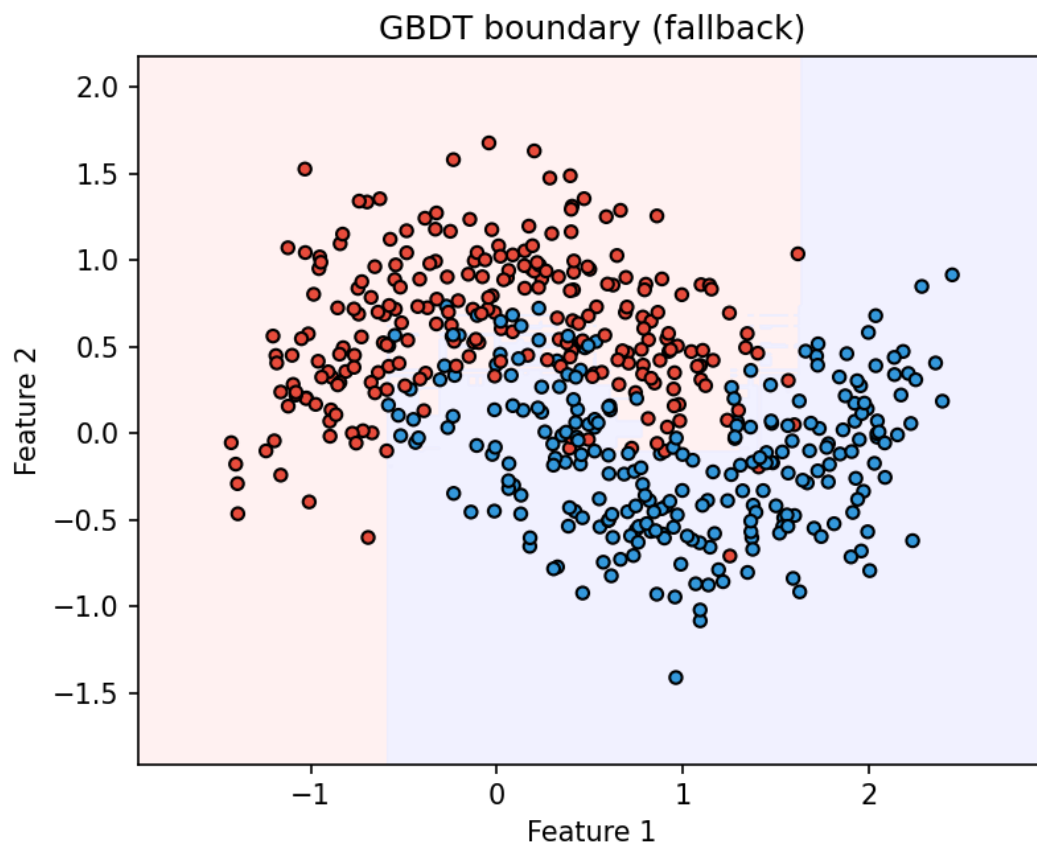


图 1: XGBoost 在两类数据上的决策边界。

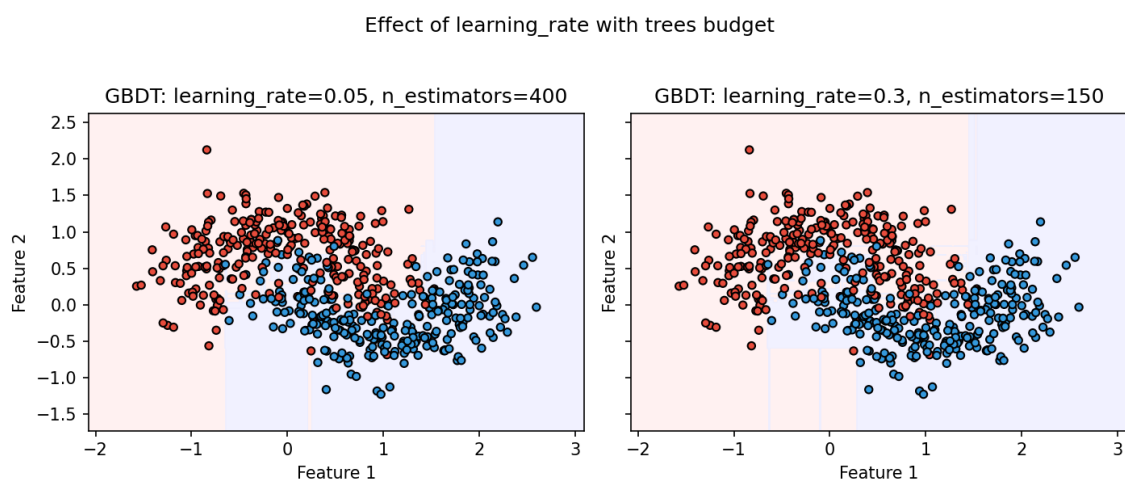


图 2: 在固定树预算下，不同学习率的效果对比。

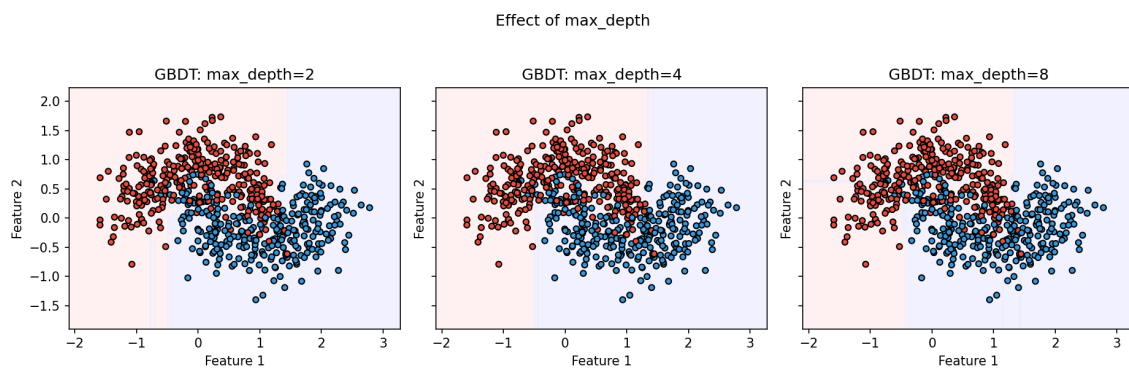


图 3: 不同 max_depth 下的决策边界对比。

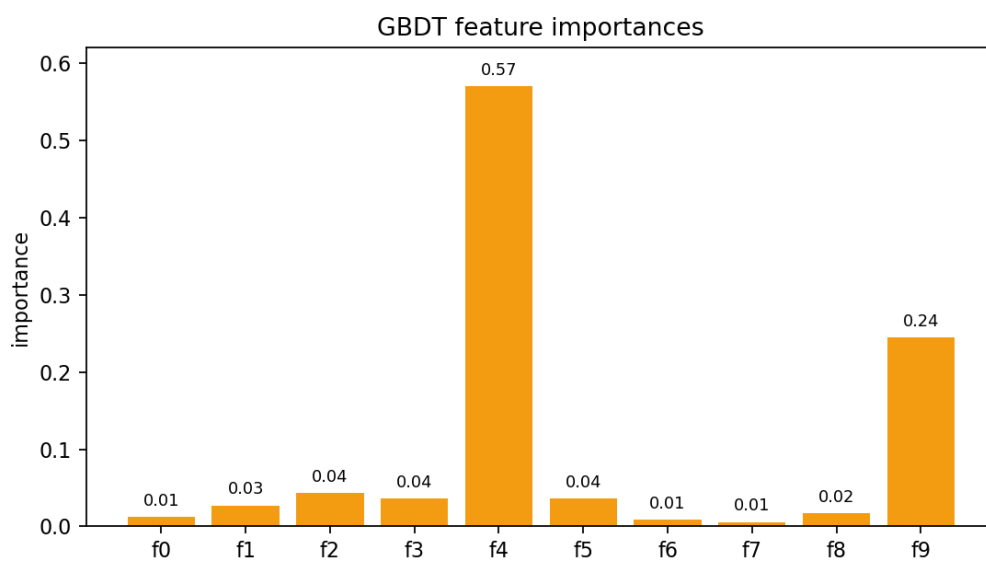


图 4: XGBoost 的特征重要性。

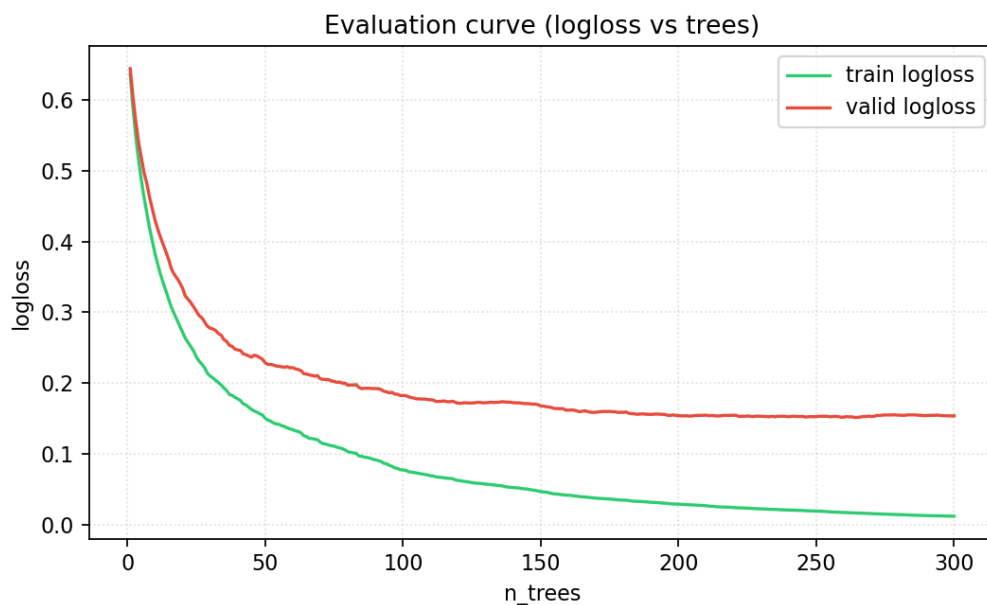


图 5: 训练/验证 logloss 随树数变化的曲线。

6 总结

XGBoost 在高效的树提升框架上引入结构正则、二阶近似与稀疏感知的分裂查找，并结合学习率与采样策略，在表格数据任务中具有很强的基准表现。通过调节深度、学习率与采样比例，可在偏差-方差与计算开销之间取得良好平衡。