# XGBoost: Theory and Practice

September 10, 2025

## 1 Introduction

XGBoost is an efficient and scalable implementation of gradient boosted decision trees (GBDTs). It improves training speed, regularization, and accuracy through second-order optimization, tree sparsity-aware split finding, shrinkage, and subsampling.

## 2 Theory and Formulas

Gradient boosting fits an additive model $F_M(\mathbf{x}) = \sum_{m=1}^{M} f_m(\mathbf{x})$ of shallow trees by stage-wise optimization. XGBoost minimizes a regularized objective

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

where $T$ is the number of leaves and $w$ are leaf scores. Using a second-order Taylor expansion of the loss around current predictions yields per-node sums of gradients $g_i$ and Hessians $h_i$; the split gain for left/right partitions $L, R$ is

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

Regularization via $\lambda, \gamma$, shrinkage (learning rate), column/row subsampling, and maximum depth/leaf constraints control complexity and reduce overfitting.

## 3 Applications and Tips

- **Scaling and types:** handles numeric and one-hot encoded categorical features; no scaling required for trees.

- **Key hyperparameters:** `n_estimators`, `max_depth`, `learning_rate`, `subsample`, `colsample_bytree`, `reg_alpha`/`reg_lambda`.

- **Early stopping:** use validation with `eval_set` and `early_stopping_rounds`.

- **Imbalance:** set `scale_pos_weight` or use stratified sampling.

- **Interpretation:** start with built-in importances; prefer permutation or SHAP for robust insights.

# 4 Python Practice

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

Listing 1: Generate XGBoost figures

```
1  python gen_xgboost_figures.py
```

Listing 2: gen_xgboost_figures.py

```
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      if path is None:
44          base = os.path.dirname(os.path.abspath(__file__))
45          path = os.path.join(base, "figures")
```

```python
46          os.makedirs(path, exist_ok=True)
47          return path
48
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(Z).
              size)
60      ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold, edgecolors="k", s=20)
61      ax.set_title(title)
62      ax.set_xlabel("Feature 1")
63      ax.set_ylabel("Feature 2")
64
65
66  def _xgb_classifier(**kwargs):
67      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 GradientBoostingClassifier
                      available.")
84          params = dict(
85              n_estimators=kwargs.get("n_estimators", 200),
86              max_depth=kwargs.get("max_depth", 3),
87              learning_rate=kwargs.get("learning_rate", 0.1),
88              random_state=0,
89          )
90          return GradientBoostingClassifier(**params)
91
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)
```

```python
 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 "GBDT
              boundary (fallback)"
101       _plot_decision_boundary(ax, clf, X, y, title)
102       out_path = os.path.join(out_dir, "xgb_decision_boundary_2class.png")
103       fig.tight_layout()
104       fig.savefig(out_path)
105       plt.close(fig)
106       return out_path
107
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), "learning_rate
                  =0.05, n_estimators=400"),
114           (_xgb_classifier(learning_rate=0.3, n_estimators=150), "learning_rate
                  =0.3, n_estimators=150"),
115       ]
116       fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=True,
              sharey=True)
117       for ax, (m, title) in zip(axes, models):
118           m.fit(X, y)
119           label = ("XGBoost: " if HAS_XGB else "GBDT: ") + title
120           _plot_decision_boundary(ax, m, X, y, label)
121       fig.suptitle("Effect of learning_rate with trees budget")
122       out_path = os.path.join(out_dir, "xgb_learning_rate_compare.png")
123       fig.tight_layout(rect=[0, 0.03, 1, 0.95])
124       fig.savefig(out_path)
125       plt.close(fig)
126       return out_path
127
128
129   def fig_xgb_max_depth_compare(out_dir: str) -> str:
130       np.random.seed(2)
131       X, y = make_moons(n_samples=600, noise=0.32, random_state=2)
132       models = [
133           (_xgb_classifier(max_depth=2, n_estimators=250), "max_depth=2"),
134           (_xgb_classifier(max_depth=4, n_estimators=250), "max_depth=4"),
135           (_xgb_classifier(max_depth=8, n_estimators=250), "max_depth=8"),
136       ]
137       fig, axes = plt.subplots(1, 3, figsize=(12.5, 4.2), dpi=150, sharex=True,
              sharey=True)
138       for ax, (m, title) in zip(axes, models):
139           m.fit(X, y)
140           label = ("XGBoost: " if HAS_XGB else "GBDT: ") + title
141           _plot_decision_boundary(ax, m, X, y, label)
142       fig.suptitle("Effect of max_depth")
143       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)
```

```python
147        return out_path


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=0.1)
161        clf.fit(X, y)
162        importances = getattr(clf, "feature_importances_", None)
163        if importances is None:
164            # Fallback: uniform zeros to avoid crash
165            importances = np.zeros(X.shape[1], dtype=float)
166
167        fig, ax = plt.subplots(figsize=(7.0, 4.0), dpi=160)
168        idx = np.arange(importances.size)
169        ax.bar(idx, importances, color="#F39C12")
170        ax.set_xticks(idx)
171        ax.set_xticklabels([f"f{i}" for i in idx])
172        ax.set_ylabel("importance")
173        title = "XGBoost feature importances" if HAS_XGB else "GBDT feature
            importances"
174        ax.set_title(title)
175        ax.set_ylim(0, max(0.25, float(importances.max()) + 0.05))
176        for i, v in enumerate(importances):
177            ax.text(i, v + 0.01, f"{v:.2f}", ha="center", va="bottom", fontsize=8)
178        out_path = os.path.join(out_dir, "xgb_feature_importances.png")
179        fig.tight_layout()
180        fig.savefig(out_path)
181        plt.close(fig)
182        return out_path


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=0.3,
            random_state=3)
195
196        if HAS_XGB:
197            clf = _xgb_classifier(n_estimators=300, learning_rate=0.1, max_depth
                =4)
```

```python
198            clf.fit(X_train, y_train, eval_set=[(X_train, y_train), (X_val, y_val)
                   ], eval_metric="logloss", verbose=False)
199            res = clf.evals_result()
200            tr = np.array(res.get("validation_0", {}).get("logloss", []), dtype=
                   float)
201            va = np.array(res.get("validation_1", {}).get("logloss", []), dtype=
                   float)
202        else:
203            # Fallback using staged decision on GradientBoosting
204            clf = _xgb_classifier(n_estimators=300, learning_rate=0.1, max_depth
                   =3)
205            clf.fit(X_train, y_train)
206            tr_list, va_list = [], []
207            # GradientBoostingClassifier provides staged_predict_proba
208            if hasattr(clf, "staged_predict_proba"):
209                for y_tr_prob, y_va_prob in zip(clf.staged_predict_proba(X_train),
                       clf.staged_predict_proba(X_val)):
210                    tr_list.append(log_loss(y_train, y_tr_prob))
211                    va_list.append(log_loss(y_val, y_va_prob))
212            else:
213                # Last resort: single-point curves
214                y_tr_prob = clf.predict_proba(X_train)
215                y_va_prob = clf.predict_proba(X_val)
216                tr_list = [log_loss(y_train, y_tr_prob)]
217                va_list = [log_loss(y_val, y_va_prob)]
218            tr, va = np.array(tr_list), np.array(va_list)
219
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="#2
                   ECC71")
222        ax.plot(np.arange(1, len(va) + 1), va, label="valid logloss", color="#
                   E74C3C")
223        ax.set_xlabel("n_trees")
224        ax.set_ylabel("logloss")
225        ax.set_title("Evaluation curve (logloss vs trees)")
226        ax.legend()
227        ax.grid(True, linestyle=":", alpha=0.4)
228        out_path = os.path.join(out_dir, "xgb_eval_logloss_curve.png")
229        fig.tight_layout()
230        fig.savefig(out_path)
231        plt.close(fig)
232        return out_path
233
234
235 def main():
236     out_dir = _ensure_figures_dir(None)
237     generators = [
238         fig_xgb_decision_boundary_2class,
239         fig_xgb_learning_rate_compare,
240         fig_xgb_max_depth_compare,
241         fig_xgb_feature_importances,
242         fig_xgb_eval_logloss_curve,
243     ]
244     print("Generating figures into:", os.path.abspath(out_dir))
```

6

```
245    if not HAS_XGB:
246        print("xgboost not found; falling back to GradientBoostingClassifier
               where possible.")
247    for gen in generators:
248        try:
249            p = gen(out_dir)
250            print("Saved:", p)
251        except Exception as e:
252            print("Failed generating", gen.__name__, ":", e)
253
254
255 if __name__ == "__main__":
256     main()
```
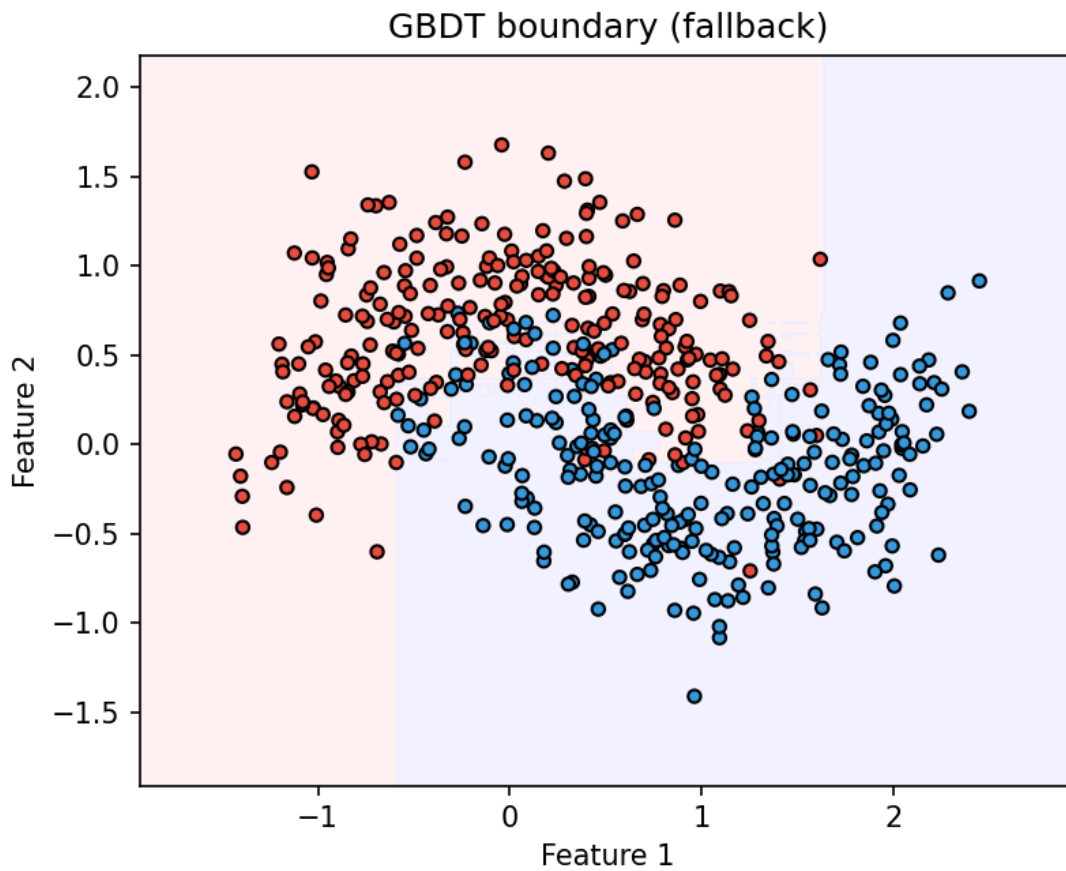
# 5   Result



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

Effect of learning_rate with trees budget



Figure 2: Learning rate effect with a budget of trees.
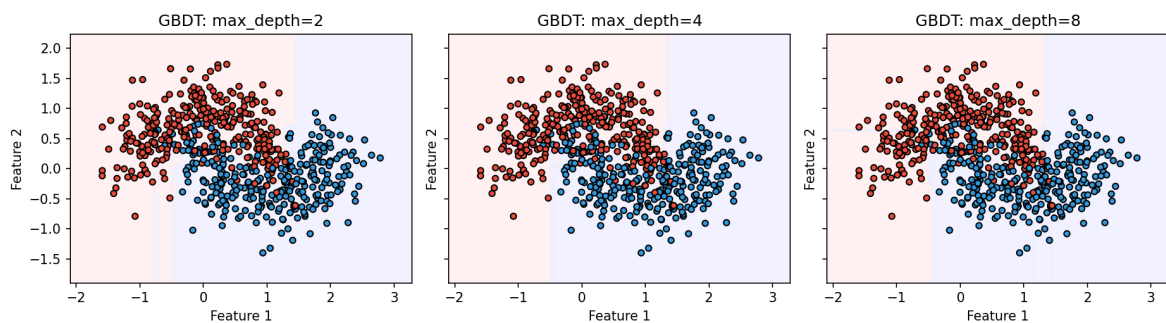
Effect of max_depth



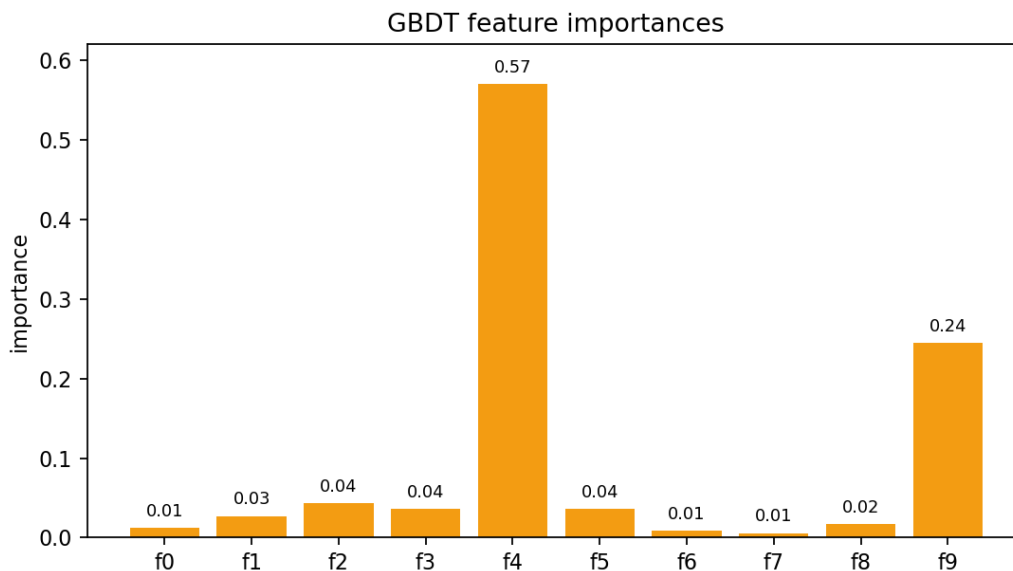Figure 3: Decision boundaries under different max_depth.
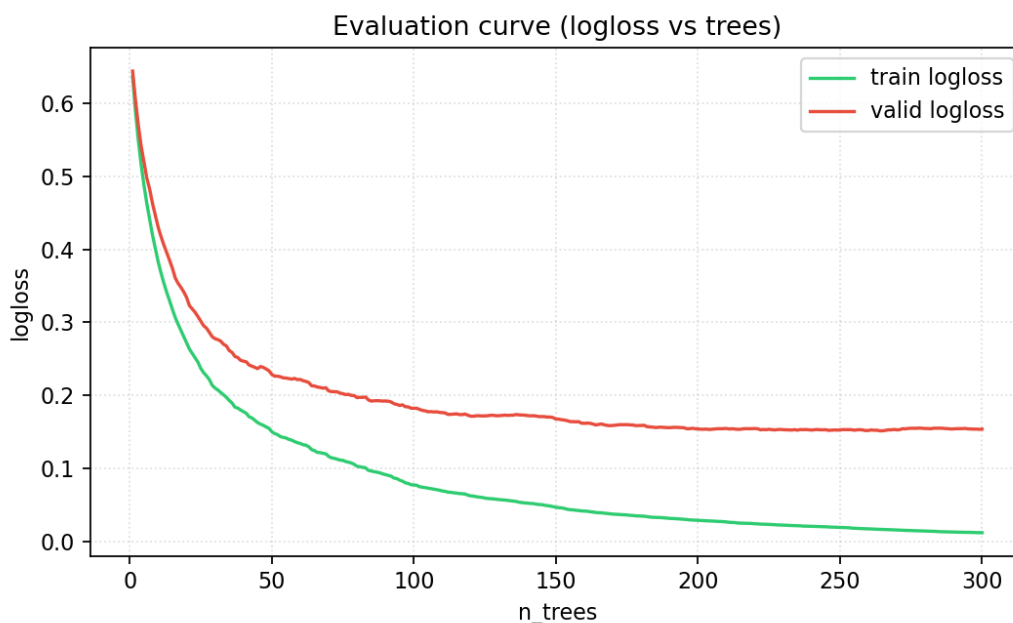
Figure 4: Feature importances from XGBoost.



Figure 5: Training/validation logloss vs number of trees.

# 6 Summary

XGBoost combines efficient tree boosting with strong regularization and advanced split finding. With careful tuning of depth, learning rate, and sampling, it delivers state-of-the-art performance on many tabular tasks.