

XGBoost is the gradient-boosting library behind countless winning models, covered here from training and tuning to early stopping.

Data: DMatrix vs the sklearn API

Two front doors to the same engine: the native DMatrix plus `xgb.train`, or `XGBClassifier`/`XGBRegressor` that take plain arrays.

PURPOSE	COMMAND	RESULT
Wrap arrays in the native container.	<pre>dtrain = xgb.DMatrix(X_train, label=y_train)</pre>	
Faster hist data container.	<pre>dtrain = xgb.QuantileDMatrix(X_train, label=y_train)</pre>	
Or skip DMatrix: the sklearn API.	<pre>clf = xgb.XGBClassifier() clf.fit(X_train, y_train)</pre>	
Add a validation set to watch.	<pre>dvalid = xgb.DMatrix(X_valid, label=y_valid)</pre>	
Carry feature names / weights.	<pre>xgb.DMatrix(X, label=y, weight=w, feature_names=cols)</pre>	

Native path: wrap arrays in `xgb.DMatrix` and train with `xgb.train`; or let `XGBClassifier`/`XGBRegressor` take plain arrays directly

Train: classifier and regressor







One engine, two estimators: `XGBClassifier` for labels, `XGBRegressor` for numbers, both `.fit` / `.predict`. The native path is `xgb.train(params, dtrain, num_boost_round)`.

PURPOSE	COMMAND	RESULT
Train a classifier (sklearn API).	<pre>clf = xgb.XGBClassifier(n_estimators=300, tree_method="hist").fit(X_train, y_train)</pre>	
Train a regressor (sklearn API).	<pre>reg = xgb.XGBRegressor(objective="reg:squarederror").fit(X_train, y_train)</pre>	
Train via the native API.	<pre>bst = xgb.train(params, dtrain, num_boost_round=300)</pre>	
Pick the objective (task).	<pre>params = {"objective": "binary:logistic", "eval_metric": "logloss"}</pre>	
Use CPU or GPU.	<pre>xgb.XGBClassifier(tree_method="hist", device="cuda")</pre>	
Same engine, two estimators.	<pre>clf.predict(X) reg.predict(X)</pre>	

`XGBClassifier` / `XGBRegressor` wrap the same boosting engine; `xgb.train(params, dtrain, num_boost_round)` returns the native `Booster`

Key hyperparameters

The three knobs that decide capacity: how many trees (`n_estimators`), how deep each tree (`max_depth`), and how much each tree counts (`learning_rate`). Lower rate plus more trees usually

PURPOSE	COMMAND	RESULT
Number of boosting rounds (trees).	<code>xgb.XGBClassifier(n_estimators=500)</code>	
Tree complexity (depth).	<code>xgb.XGBClassifier(max_depth=6)</code>	
Shrinkage per tree.	<code>xgb.XGBClassifier(learning_rate=0.05)</code>	
Subsample rows / columns.	<code>xgb.XGBClassifier(subsample=0.8, colsample_bytree=0.8)</code>	
Regularize leaf weights.	<code>xgb.XGBClassifier(reg_lambda=1.0, min_child_weight=1)</code>	
Old-style native aliases (avoid in sklearn API).	<code>nativeparams={"eta": 0.05, "lambda": 1.0}</code>	

`n_estimators`, `max_depth`, `learning_rate` are the three capacity knobs; use the underscored sklearn names, not the native `eta` / `lambda` / `alpha` aliases

Early Stopping

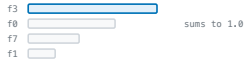




Let a validation set choose the tree count. Set `early_stopping_rounds` in the constructor and pass `eval_set` to fit; training halts when the metric stops improving, and `best_iteration` is recorded.

PURPOSE	COMMAND	RESULT
Turn on early stopping (sklearn API).	<code>xgb.XGBClassifier(n_estimators=1000, early_stopping_rounds=20, eval_metric="logloss")</code>	
Provide the watch set in fit.	<code>clf.fit(X_train, y_train, eval_set=[(X_valid, y_valid)])</code>	
Read where it stopped.	<code>clf.best_iteration</code> and <code>clf.best_score</code>	
Native-API early stopping.	<code>bst = xgb.train(params, dtrain, num_boost_round=1000, evals=[(dvalid, "valid")], early_stopping_rounds=20)</code>	
Callback form (save the best).	<code>xgb.callback.EarlyStopping(rounds=20, save_best=True)</code>	
Deprecated: passing it in fit (sklearn).	<code>clf.fit(..., early_stopping_rounds=20, eval_metric=...)</code>	

Set a generous `n_estimators`, put `early_stopping_rounds` + `eval_metric` in the constructor, pass `eval_set` to fit; read `best_iteration` / `best_score`

Feature Importance (gain)



Which features did the trees actually use? gain is the default and most informative; weight and cover are alternatives.

PURPOSE	COMMAND	RESULT
Normalized importances (sklearn API).	<code>clf.feature_importances_</code>	
Raw gain per feature (native scores).	<code>clf.get_booster().get_score(importance_type="gain")</code>	
Compare importance types.	<code>get_score(importance_type="weight")</code> and <code>"cover"</code>	
Plot importance directly.	<code>xgb.plot_importance(clf, importance_type="gain", max_num_features=10)</code>	
Draw a single tree.	<code>xgb.plot_tree(clf, num_trees=0)</code>	

`feature_importances_` gives normalized scores; `get_score(importance_type="gain")` gives raw loss reduction per split, usually the best ranking

Persist: save & load a booster



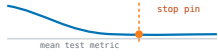


Train once, serve many times: save to portable .json or compact .ubj, reload in another process.

PURPOSE	COMMAND	RESULT
Save the fitted model (sklearn API).	<code>clf.save_model("model.json")</code>	
Save as compact binary (UBJSON).	<code>clf.save_model("model.ubj")</code>	
Load it back into an estimator.	<code>clf2 = xgb.XGBClassifier()</code> <code>clf2.load_model("model.json")</code>	
Save / load a native Booster.	<code>bst.save_model("booster.ubj")</code> <code>xgb.Booster().load_model("booster.ubj")</code>	
Deprecated: pickle / old .bin.	<code>pickle.dump(bst, ...) / ".bin"</code>	
Stable, version-checked formats.	<code>.json = text / .ubj = binary</code>	

`save_model("model.json"/".ubj")` and `load_model` reload a model in another process; avoid pickle and the legacy .bin path

Cross-validation


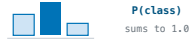



xgb.cv runs k-fold CV and, with early stopping, finds a good num_boost_round honestly.

PURPOSE	COMMAND	RESULT									
K-fold CV on a DMatrix.	<pre>cv = xgb.cv(params, dtrain, num_boost_round=500, nfold=5)</pre>										
Stratified folds for classification.	<pre>xgb.cv(params, dtrain, nfold=5, stratified=True)</pre>										
CV with early stopping.	<pre>xgb.cv(params, dtrain, nfold=5, early_stopping_rounds=20)</pre>										
Read the results table.	<pre>cv # a pandas DataFrame</pre>	<table border="1" data-bbox="1133 596 1349 634"> <thead> <tr> <th>train-mean</th> <th>test-mean</th> <th>test-std</th> </tr> </thead> <tbody> <tr> <td>0.31</td> <td>0.39</td> <td>0.02</td> </tr> <tr> <td>0.28</td> <td>0.37</td> <td>0.02</td> </tr> </tbody> </table>	train-mean	test-mean	test-std	0.31	0.39	0.02	0.28	0.37	0.02
train-mean	test-mean	test-std									
0.31	0.39	0.02									
0.28	0.37	0.02									
Or use sklearn cross-validation.	<pre>from sklearn.model_selection import cross_val_score cross_val_score(clf, X, y, cv=5)</pre>										
Pick num_boost_round honestly.	<pre>best_round = len(cv)</pre>										

xgb.cv returns a tidy per-round train/test DataFrame; add stratified and early_stopping_rounds to pick num_boost_round defensibly

Predict & Explain (SHAP-style)

Get labels, probabilities, or per-feature contributions. predict_contribs returns SHAP values straight from the booster so each prediction breaks down into a base value plus one signed push per feature

PURPOSE	COMMAND	RESULT
Predict labels.	<pre>y_pred = clf.predict(X_test)</pre>	
Predict probabilities.	<pre>proba = clf.predict_proba(X_test)</pre>	
SHAP-style contributions (native).	<pre>bst.predict(dtest, pred_contribs=True)</pre>	
Use the SHAP library directly.	<pre>import shap; shap.TreeExplainer(clf) .shap_values(X_test)</pre>	
Explain one prediction.	<pre>bst.predict(dtest, pred_contribs=True)[i]</pre>	
Score on a test set.	<pre>clf.score(X_test, y_test)</pre>	

predict gives labels, predict_proba gives probabilities; pred_contribs=True returns SHAP values: a base value plus one signed push per feature