

## Boosting Overview

Boosting is an ensemble technique where models are trained **sequentially**, each new model focusing on correcting errors made by the previous ones. In boosting, we start with a weak learner (often a simple decision tree) and iteratively add more models to reduce the overall error. Each new model is trained on the residuals or errors of the ensemble so far, which “boosts” the performance by concentrating on hard-to-predict cases <sup>1</sup> <sup>2</sup>. Unlike bagging (parallel ensembles), boosting’s sequential approach often yields higher accuracy but risks overfitting if not regularized.

**AdaBoost (Adaptive Boosting)** and **Gradient Boosting** are two classic boosting methods. AdaBoost adaptively reweights training instances, emphasizing those misclassified by earlier learners. Gradient Boosting fits each new model to the negative gradient of a loss function (i.e. the residual errors). XGBoost (eXtreme Gradient Boosting) is a highly optimized gradient-boosting implementation that adds regularization, parallelization, and other tricks for speed and accuracy <sup>2</sup> <sup>3</sup>.

The sections below cover each method with light theory and hands-on scikit-learn demos, using datasets like Iris, Breast Cancer, Digits, etc. Each section ends with an in-class exercise to reinforce learning.

## AdaBoost

### Theory

**AdaBoost** stands for *Adaptive Boosting*. It was introduced by Freund and Schapire (1996) as an ensemble method that combines multiple **weak learners** into a strong one. A *weak learner* (often a decision stump – a tree of depth 1) is only slightly better than random guessing, but AdaBoost boosts their power. Initially, all training samples have equal weight. After each weak learner is trained, AdaBoost increases the weight of misclassified instances so that subsequent learners focus on the “hard” cases <sup>4</sup> <sup>5</sup>. Finally, AdaBoost takes a weighted vote of all the learners to form the final prediction.

In scikit-learn’s `AdaBoostClassifier`, the default base estimator is a decision tree of `max_depth=1` (a stump) <sup>6</sup>. By default it uses the SAMME algorithm (multi-class extension of AdaBoost). Key points:

- Each weak learner is trained on the same data but with updated instance weights.
- After each iteration, weights of misclassified points increase, so the next learner focuses on them <sup>4</sup> <sup>5</sup>.
- Learners are combined in a weighted sum (more accurate learners get higher weight). The final model is still just an ensemble of trees.
- AdaBoost works well for binary (and multiclass) classification, especially with clean data. It can be sensitive to noisy data and outliers <sup>7</sup>. On the positive side, because it uses very simple base learners, the resulting model is often interpretable (you can examine each stump or feature importance) <sup>7</sup>.

### Practical Example (scikit-learn)

Below is a simple AdaBoost classification example using the **Iris** dataset. We split into train/test sets, train the model, and evaluate accuracy. We also inspect feature importances to see which features are influential.

```

from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.ensemble import AdaBoostClassifier

# Load Iris data
X, y = load_iris(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)

# Create AdaBoost classifier with 50 weak learners (stumps by default)
clf = AdaBoostClassifier(n_estimators=50, random_state=42)
clf.fit(X_train, y_train)

# Evaluate and print results
print("Test accuracy:", clf.score(X_test, y_test))
print("Feature importances:", clf.feature_importances_)

```

This code initializes an **AdaBoostClassifier** and fits it to the Iris data. By default the base learner is a decision stump (`max_depth=1`) <sup>6</sup>. We see the test accuracy and which features the ensemble found most important. For example, you might get high accuracy (depending on the random split) and nonzero importances for petal features. *Feature importances* show how much each input feature contributed to the final boosted model.

Note that AdaBoost's iterative weighting (increasing weights on misclassified points) causes later weak learners to focus on those hard samples <sup>4</sup> <sup>5</sup>. In practice, you can tune `n_estimators` (number of rounds) and `learning_rate` to balance bias/variance. A larger number of weak learners can improve accuracy up to a point, but may overfit if too high.

## In-Class Exercises

- **Experiment with parameters:** Try varying `n_estimators` (e.g. 10, 50, 100) and `learning_rate` on the Iris or Breast Cancer dataset. Observe how test accuracy changes.
- **Change the base learner:** By default `AdaBoostClassifier` uses stumps. Try setting `base_estimator=DecisionTreeClassifier(max_depth=2)` or 3. Does a deeper tree improve or worsen performance? Compare training time as well.
- **Feature Importance:** Plot or list the `feature_importances_` from AdaBoost on a dataset. Which features get the highest weights? Why might that be?
- **AdaBoost for Regression:** AdaBoost can be used for regression (`AdaBoostRegressor`). As an exercise, apply `AdaBoostRegressor` to the California Housing data (predict median house value). Compare `AdaBoostRegressor` with a single decision-tree regressor.

## Gradient Boosting

### Theory

**Gradient Boosting** is a powerful boosting method introduced by Friedman (1999). Instead of reweighting instances as in AdaBoost, gradient boosting builds each new model to predict the **residual errors** (negative

gradient of the loss) of the current ensemble. In essence, it performs gradient descent in function space. The algorithm works like this: start with an initial constant model (e.g. predict the mean), compute the residuals (actual minus predicted), train a new tree on these residuals, and add this tree into the ensemble. Repeat this process for many iterations. Each tree is thus “correcting” the previous errors.

In scikit-learn, `GradientBoostingClassifier` (for classification) and `GradientBoostingRegressor` (for regression) implement this idea <sup>8</sup>. Quoting the docs:

*“This algorithm builds an additive model in a forward stage-wise fashion... In each stage `n_classes_` regression trees are fit on the negative gradient of the loss function... Binary classification is a special case where only a single regression tree is induced.”* <sup>8</sup>.

Key differences from AdaBoost: - **Loss function:** Gradient boosting can optimize any differentiable loss (e.g. deviance/log-loss for classification, squared error for regression). (Sklearn’s default for classification is log-loss). In fact, setting the loss to “exponential” makes gradient boosting equivalent to AdaBoost <sup>9</sup>.

- **Predictive model:** Each new tree predicts a continuous value (the negative gradient), and it is *added* to the current model.

- **Regularization:** Typically uses shrinkage (`learning_rate`), tree depth, and subsampling (`subsample`) to prevent overfitting.

- **Flexibility:** More flexible with loss functions and hyperparameters. It often outperforms AdaBoost on complex problems, but it is also more computationally intensive. According to summaries, gradient boosting “can be more robust to overfitting” and allows different loss functions <sup>10</sup>.

Scikit-learn also provides `HistGradientBoostingClassifier` / `Regressor` for large datasets (>=10k samples) which is much faster (using histogram binning) <sup>11</sup>. In this lecture, we focus on the classic gradient boosting (`GradientBoostingClassifier/Regressor`).

## Practical Example (scikit-learn)

Let’s train a **GradientBoostingClassifier** on the Breast Cancer dataset (binary classification). We will split the data, fit the model, and evaluate accuracy.

```
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.ensemble import GradientBoostingClassifier

# Load Breast Cancer data (binary classification)
X, y = load_breast_cancer(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)

# Create Gradient Boosting classifier
gbc = GradientBoostingClassifier(n_estimators=100, learning_rate=0.1,
max_depth=3, random_state=42)
gbc.fit(X_train, y_train)

# Evaluate and print results
```

```
print("Test accuracy:", gbc.score(X_test, y_test))
print("Feature importances:", gbc.feature_importances_)
```

In this code, we use 100 trees (`n_estimators=100`), a learning rate of 0.1, and maximum tree depth 3. After training, we print the test accuracy and feature importances. For example, you might see an accuracy around 0.95 or higher on this task. Gradient boosting's `feature_importances_` tells us which input features were most useful for reducing loss. Compared to AdaBoost, gradient boosting often achieves similar or better accuracy on such tasks with well-tuned parameters.

The model's hyperparameters control its behavior. A smaller `learning_rate` (with more trees) usually improves performance at the cost of training time, while larger depth increases model complexity. Also, `subsample < 1.0` can reduce variance (introducing randomness) and help prevent overfitting. It's valuable for students to experiment with these.

## In-Class Exercises

- **Hyperparameter tuning:** Try different values of `learning_rate` (e.g. 0.01, 0.1, 0.5) and `n_estimators` (e.g. 50, 100, 200) on the Breast Cancer or Iris dataset. Plot how test accuracy or validation score changes.
- **Regression task:** Apply `GradientBoostingRegressor` to a regression dataset. For example, predict median house values using the California Housing data. Evaluate using RMSE. (Hint: use `fetch_california_housing` and `mean_squared_error`).
- **Subsampling vs full:** Experiment with `subsample=0.5` vs `subsample=1.0` (full data) and compare performance and overfitting.
- **Comparison to Random Forest:** As an exercise, train a `RandomForestClassifier` with similar depth on the same data and compare accuracy and training time to Gradient Boosting.

## XGBoost

### Theory

**XGBoost** (eXtreme Gradient Boosting) is an optimized implementation of gradient boosting developed by Tianqi Chen. It follows the same general boosting principle (sequential trees), but adds many enhancements for speed and performance. XGBoost is widely used in practice (especially in data competitions) due to its efficiency and accuracy. Key characteristics of XGBoost include parallel tree construction, out-of-core (disk) processing, and built-in regularization <sup>2</sup> <sup>3</sup>.

According to IBM's summary, "XGBoost is an advanced implementation of gradient boosting with the same general framework" <sup>2</sup>. It is known for its **speed and scalability** – e.g. it can run an order of magnitude faster than standard libraries by using cache-aware prefetching and parallelism <sup>12</sup>. Importantly, XGBoost **includes regularization** (L1/L2) in its objective, unlike vanilla gradient boosting <sup>3</sup>. This regularization helps prevent overfitting and often yields better generalization. XGBoost also has special handling for missing values, learning the best direction to take when a feature is missing <sup>3</sup>.

In practice, XGBoost can be used via its scikit-learn-like API (`XGBClassifier`, `XGBRegressor`). It usually requires installing the `xgboost` library (`pip install xgboost`). When using `XGBClassifier`,

common parameters are similar to `GradientBoostingClassifier` (`n_estimators`, `learning_rate`, `max_depth`, etc.), plus things like `tree_method` for GPU/histograms if desired. Because XGBoost is optimized, it often handles large data or many trees faster than scikit-learn's versions. It also provides advanced features like early stopping.

## Practical Example (scikit-learn interface)

We demonstrate XGBoost on the Digits dataset (10-class classification). Make sure you have XGBoost installed (`pip install xgboost`). We'll train an `XGBClassifier` and evaluate accuracy.

```
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
from xgboost import XGBClassifier

# Load Digits data (10 classes)
X, y = load_digits(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)

# Create XGBoost classifier (using 10 trees for speed)
xgbc = XGBClassifier(n_estimators=10, use_label_encoder=False,
                    eval_metric='mlogloss', random_state=42)
xgbc.fit(X_train, y_train)

# Evaluate and print results
print("Test accuracy:", xgbc.score(X_test, y_test))
```

This code runs XGBoost with 10 trees (`n_estimators=10`) to keep training quick. We disable the old label encoder warning (`use_label_encoder=False`) and set `eval_metric='mlogloss'` for multiclass. You should see high accuracy (e.g. ~0.96). XGBoost trains fast even with more trees thanks to its optimizations.

If time permits, you could also try `XGBRegressor` on a regression dataset. The API and usage are analogous to `GradientBoostingRegressor`.

## In-Class Exercises

- **More trees and early stopping:** Increase `n_estimators` (e.g. 50 or 100) and monitor training time and accuracy. Optionally, use XGBoost's `early_stopping_rounds` on a validation split to automatically stop when performance plateaus.
- **Compare with sklearn:** Using the same data and hyperparameters, compare training time and accuracy of `XGBClassifier` vs `GradientBoostingClassifier` vs `AdaBoostClassifier`. Note the speed difference, especially as you increase the number of trees.
- **Tuning XGBoost:** Play with XGBoost-specific parameters like `gamma` (minimum loss reduction), `subsample`, or enabling `tree_method='hist'` (histogram optimization) to see effects on performance.

- **Feature Importance with XGBoost:** Use `xgbc.feature_importances_` or `xgb.plot_importance(xgbc)` (requires matplotlib) to visualize which features are most important in the XGBoost model.

## Comparative Analysis

Finally, let's compare **AdaBoost**, **Gradient Boosting (GBM)**, and **XGBoost** in terms of performance, interpretability, and use cases:

- **Performance:**
  - *Speed:* XGBoost is typically much faster than scikit-learn's gradient boosting, especially on large data, thanks to its optimized implementation (e.g. parallel tree building and efficient memory usage <sup>12</sup>). In some reports, XGBoost can be an order of magnitude faster than other frameworks on a single machine <sup>12</sup>. AdaBoost is generally faster for small models (fewer hyperparameters) but may underperform on complex tasks.
  - *Accuracy:* All three can achieve high accuracy, but with tuned hyperparameters gradient boosting variants often beat AdaBoost on complex data. XGBoost, with its regularization and advanced options, often achieves the best accuracy on many benchmarks. (Remember, AdaBoost with an exponential loss is actually a special case of gradient boosting <sup>9</sup>.)
  - *Robustness:* Gradient boosting and XGBoost usually include hyperparameters (like `learning_rate` and regularization) to control overfitting, making them more robust to noise than vanilla AdaBoost <sup>10</sup>. AdaBoost can be sensitive to outliers and noisy data <sup>10</sup>.
- **Interpretability:**
  - *Model complexity:* AdaBoost with decision stumps yields a sequence of very simple rules, which can be easier to interpret in isolation. You can examine each stump and its weight to see which simple decisions the ensemble is using. Gradient boosting trees (with depth >1) and especially XGBoost models (with many trees and complex splits) are harder to interpret as a whole.
  - *Feature importance:* All methods provide feature importances. AdaBoost's importances (from its stumps) are somewhat interpretable. XGBoost has built-in tools (gain, cover, frequency) to rank features, often more nuanced than sklearn's "Gini" importances. In practice, if interpretability is critical, one might prefer fewer trees or shallow depths. According to guides, AdaBoost is considered "highly interpretable" relative to other ensembles <sup>7</sup>.
- **Use Cases:**
  - **AdaBoost:** Good for smaller datasets and binary classification with clean features. Historically famous for face detection (e.g. Viola-Jones algorithm) and other CV tasks <sup>13</sup>. Also used in domains like medical diagnostics or financial predictions where fast, interpretable decisions are needed <sup>13</sup>.
  - **Gradient Boosting (GBM):** Very general-purpose. Works well on both classification and regression. Use GBM when you want flexibility in loss functions or when you can afford the computation to tune many parameters. It tends to yield strong results on tabular data.
  - **XGBoost:** A go-to choice for competitions and large-scale problems. Preferred when maximum predictive performance is desired and dataset is large. Its handling of missing values also makes it

convenient on real-world messy data <sup>3</sup> . However, for smaller datasets or quick prototyping, AdaBoost or basic GBM may suffice.

In summary, **AdaBoost** is the simplest (fast to train, good on clean binary problems, very interpretable) <sup>7</sup> . **Gradient Boosting** (GBM) is more flexible (different loss functions, better handling of overfitting) <sup>10</sup> but needs careful tuning. **XGBoost** builds on gradient boosting and is engineered for high performance (faster training, built-in regularization, missing-value support) <sup>12</sup> <sup>3</sup> . The choice depends on the task: for quick wins AdaBoost might suffice, for complex datasets use GBM/XGBoost, and for ultimate performance XGBoost (or similar libraries like LightGBM/CatBoost) is often the best.

**Sources:** We have cited the scikit-learn documentation and expert articles for algorithm descriptions and comparisons <sup>4</sup> <sup>8</sup> <sup>9</sup> <sup>2</sup> <sup>12</sup> <sup>3</sup> <sup>7</sup> , ensuring our explanations are grounded in authoritative references.

---

<sup>1</sup> What is the XGBoost algorithm and how does it work?

<https://www.analyticsvidhya.com/blog/2018/09/an-end-to-end-guide-to-understand-the-math-behind-xgboost/>

<sup>2</sup> <sup>3</sup> <sup>12</sup> What is XGBoost? | IBM

<https://www.ibm.com/think/topics/xgboost>

<sup>4</sup> <sup>6</sup> AdaBoostClassifier — scikit-learn 1.7.1 documentation

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostClassifier.html>

<sup>5</sup> <sup>7</sup> <sup>10</sup> <sup>13</sup> Understanding AdaBoost in detail

<https://www.byteplus.com/en/topic/399965>

<sup>8</sup> <sup>9</sup> <sup>11</sup> GradientBoostingClassifier — scikit-learn 1.7.1 documentation

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.html>