



Machine Learning for Econometrics

Lecture 4: Boosting

Yi He

November 22, 2022

Plan for Today

1. Introduction
2. Least Squares Boosting
3. AdaBoost

Introduction

Example: Predicting Stock Direction

Let the target $Y_t = \text{sign}(r_t) \in \{-1, 1\}$ be the sign of the stock return at time t , and the feature $X_t = r_{t-1}$ to be the return at time $t - 1$.

- Momentum Rule

$$g(X_t) = \begin{cases} 1 & X_t > 0 \\ -1 & X_t < 0 \end{cases}$$

- Reversion Rule

$$-g(X_t) = \begin{cases} -1 & X_t > 0 \\ 1 & X_t < 0 \end{cases}$$

- Error rates for daily returns on S&P 500 index between 2012 and 2021 ($n = 2516$): Momentum = 52%, Reversion = 48%.

Decision Stumps

- Target $Y \in \{-1, 1\}$ and a univariate feature $X \in \mathbb{R}$.
- Consider a simple threshold rule for decision:

$$g(X) = \begin{cases} 1 & X > \theta \\ -1 & X < \theta \end{cases} \text{ or } -g(X) = \begin{cases} -1 & X > \theta \\ 1 & X < \theta \end{cases}$$

- Decision stumps:

$$\mathcal{G} = \{g(x) = \text{sign}(\theta - x) \cdot b : \theta \in \mathbb{R}, b = \pm 1\},$$

- For multivariate feature $X \in \mathbb{R}^d$:

$$\mathcal{G} = \left\{ g(x) = \text{sign}(\theta - x_j) \cdot b : \underbrace{j \in \{1, \dots, d\}}_{\text{pick one feature}}, \theta \in \mathbb{R}, b = \pm 1 \right\}.$$

Decision Stumps As Weak Classifiers

- Take any decision stump $g \in \mathcal{G}$, then $-g \in \mathcal{G}$
- $g \in \mathcal{G}$ or $-g \in \mathcal{G}$ is no worse than a random guess but shows only weak predictive performance in general:

$$\min \{\mathbb{P}(g(X) \neq Y), \mathbb{P}(-g(X) \neq Y)\} \leq \frac{1}{2}$$

because

$$\begin{aligned} & \mathbb{P}(g(X) \neq Y) + \mathbb{P}(-g(X) \neq Y) \\ &= \mathbb{P}(g(X) \neq Y) + \mathbb{P}(g(X) = Y) = 1. \end{aligned}$$

- In our stock example: reversion rule = 48% < 50%.

Does Ensembling Work?

Can we combine the performance of many weak classifiers, say, $g_1(x), \dots, g_M(x)$ to produce a powerful committee?

Consider $M = 25$ weak classifiers and fix some value x :

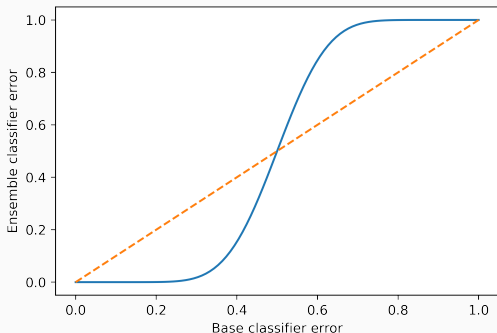
- Each classifier $g_m \in \{-1, 1\}$ has equal base error rate $\varepsilon = 0.35$
- Assume errors made by classifiers are **independent** (hardly true in practice, but let us assume it here for simplicity)
- Ensemble classifier

$$\text{sign}(g_1 + \dots + g_{25}) = \begin{cases} 1 & g_1 + \dots + g_{25} > 0 \\ -1 & g_1 + \dots + g_{25} < 0, \end{cases}$$

where sign denotes the sign function.

Probability that the committee makes a wrong prediction:

$$\begin{aligned} & \mathbb{P}[\text{sgn}(g_1 + \dots + g_{25}) \neq Y | X = x] \\ &= \mathbb{P}\left(\sum_{m=1}^{25} \mathbb{1}[g_m \neq Y] \geq 13 \mid X = x\right) = \sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1 - \varepsilon)^{25-i} = 0.06 \end{aligned}$$



Least Squares Boosting

Additive Regression Models

Consider the candidate prediction rule in an additive form given by

$$f(x) = \sum_{m=1}^M f_m(x), \quad f_m \in \mathcal{G}, \quad (1)$$

where

- \mathcal{G} is a set of base regression learners (stumps or T-terminal node trees) such that

$$g \in \mathcal{G} \Rightarrow w \cdot g \in \mathcal{G}, \quad \forall w \in \mathbb{R}$$

- M is some positive integer (a hyperparameter)

Let $\text{span}(\mathcal{G})$ denote the set of function generated by (1) over all possible choices of $M \in \{1, 2, 3, \dots\}$.

Least-Squares Regression With Additive Models

Our goal is to estimate an ideal prediction rule $f^* \in \text{span}(\mathcal{G})$ that (nearly) minimizes the **population** risk

$$L_{\mathcal{D}}(f) = \mathbb{E}[\ell(f(X), Y)]$$

for the (half) squared loss function

$$\ell(f(X), Y) = \frac{1}{2} (f(X) - Y)^2$$

- If $\mu(x) = \mathbb{E}[Y|X = x] \in \text{span}(\mathcal{G})$, then it is the ideal rule.
- Unlike neural networks, there is no guarantee of universal approximation property with stumps. [For advanced reader, see Theorem 7 in Alon et al.(2022)]

Fitting Additive Regression Models: First Attempt

Let us take the bases as stumps such that

$$\begin{aligned}f_m(x) &= g(x; j_m, \theta_m, c_{m1}, c_{m2}) \\ &= c_{m1} \mathbf{1}[x_{j_m} < \theta_m] + c_{m2} \mathbf{1}[x_{j_m} \geq \theta_m].\end{aligned}$$

Given a dataset $\{Y_i, X_i : i = 1, \dots, n\}$, the empirical risk is

$$L_S(f) = \frac{1}{2n} \sum_{i=1}^n \left(\sum_{m=1}^M f_m(X_i) - Y_i \right)^2.$$

- Difficult to solve for large M : better optimization/statistical method?
- How to regularize the estimation?

Greedy Stagewise Approach: One At A Time

- The least-squares estimation is easy with $M = 1$, however, and denote the fitted function to be \hat{f}_1 .
- For $M = 2$, the standard algorithm does not store \hat{f}_1 , but optimize f_1 and f_2 simultaneously again. Note that \hat{f}_1 may change.
- To save computational efforts, the stagewise method stores \hat{f}_1 but only fit \hat{f}_2 by minimizing

$$\frac{1}{2n} \sum_{i=1}^n \left(\hat{f}_1(X_i) + f_2(X_i) - Y_i \right)^2 = \frac{1}{2n} \sum_{i=1}^n \left(f_2(X_i) - \tilde{Y}_i^{(1)} \right)^2,$$

where $\tilde{Y}_i^{(1)} = Y_i - \hat{f}_1(X_i)$ are the first-step residuals.

- Repeat this procedure to fit new bases one-at-a-time sequentially: fit \hat{f}_m to the residuals from the previous step $\tilde{Y}_i^{(m-1)} = Y_i - \hat{f}_1(X_i) - \dots - \hat{f}_{m-1}(X_i)$.

Sequential Optimization

- Initialize $F^{(0)}(x)$ with zero or a constant $F^{(0)}(x) = \bar{Y}$ the sample average of the target values.
- For $m = 1$ to M do:
 1. Calculate the residuals $\tilde{Y}_i = Y_i - F^{(m-1)}(X_i)$
 2. Fit a base to the residuals:

$$f_m = \operatorname{argmin}_{f \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \left(f(X_i) - \tilde{Y}_i^{(m-1)} \right)^2.$$

3. Accumulate the base learners:

$$F^{(m)} = F^{(m-1)} + f_m$$

- Output $\hat{f}(x) = F^{(M)}(x)$.

Regularization: Early Stopping

- Boosting *forever* can overfit (although slowly)!

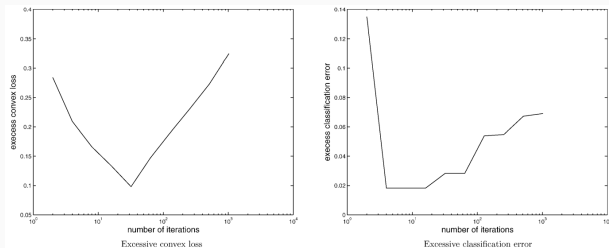


Figure 3 from Zhang and Yu (2005)

- A simple strategy is early-stopping: choose M based on a validation set

Regularization: Shrinkage

- Initialize $F^{(0)}(x)$ with zero or a constant $F^{(0)}(x) = \bar{Y}$ the sample average of the target values.
- For $m = 1$ to M do:
 1. Calculate the residuals $\tilde{Y}_i = Y_i - F^{(m-1)}(X_i)$
 2. Fit a base to the residuals:

$$g_m = \operatorname{argmin}_{f \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \left(f(X_i) - \tilde{Y}_i^{(m-1)} \right)^2.$$

3. Set $f_m = \eta \cdot g_m \in \mathcal{G}$. Accumulate the base learners:

$$F^{(m)} = F^{(m-1)} + f_m = F^{(m-1)} + \eta \cdot g_m$$

- Output $\hat{f}(x) = F^{(M)}(x)$.

The hyperparameter $\eta \in (0, 1)$ is called ‘learning rate’ or ‘shrinkage parameter’.

AdaBoost

- Binary target $Y \in \{-1, 1\}$ and features $X \in \mathbb{R}^d$
- Recall the softmax representation

$$\begin{bmatrix} \mathbb{P}(Y = 1|X = x) \\ \mathbb{P}(Y = -1|X = x) \end{bmatrix} = \begin{bmatrix} \frac{\exp(a_1(x))}{\exp(a_1(x)) + \exp(a_2(x))} \\ \frac{\exp(a_2(x))}{\exp(a_1(x)) + \exp(a_2(x))} \end{bmatrix} = \begin{bmatrix} \sigma(a(x)) \\ \sigma(-a(x)) \end{bmatrix}$$

where

$$a(x) = a_1(x) - a_2(x) = \log \frac{\mathbb{P}(Y = 1|X = x)}{\mathbb{P}(Y = -1|X = x)}$$

is called the **log odds**, and $\sigma(a) = \frac{1}{1 + \exp(-a)}$ is the sigmoid function (Lecture 2).

- The Bayes classifier $C^{\text{Bayes}}(x) = \text{sign}(a(x))$.

Additive Logistic Model

Suppose we want to estimate $a(x)$ with some additive models

$$f(x) = \sum_{m=1}^M f_m(x), \quad f_m \in \mathcal{B}.$$

such that the resulted Bayes decision rule

$$h(x) = \text{sign} \left(\sum_{m=1}^M f_m(x) \right) = \text{sign} \left(\sum_{m=1}^M \alpha_m g_m(x) \right)$$

is a **weighted** majority vote. Each $g_m(x) \in \{-1, 1\}$ is a 'voter', and $\alpha_m > 0$ is its voting weight.

\Rightarrow Take \mathcal{B} to be stumps such that

$$\mathcal{B} = \{f(x) = \alpha g(x) : \alpha > 0, g \in \mathcal{G}\},$$

where \mathcal{G} denotes the set of base **classifiers** so that $g(x) \in \{-1, 1\}$.

Exponential Loss

(Tute) The true log odds minimize the expected **exponential loss**:

$$a(x) = \operatorname{argmin}_{f: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E}[\ell_{\text{exp}}(f(X), Y) \mid X = x],$$

with

$$\ell_{\text{exp}}(f(x), y) = \phi(f(x) \cdot y), \quad \phi(x) = \exp\left(-\frac{1}{2}x\right).$$

- The exponential loss is similar to hinge loss:

$$\ell_{\text{Hinge}}(f(x), y) = \phi(f(x) \cdot y), \quad \phi(x) = \max\{0, 1 - x\}.$$

- The exponential loss is similar to the zero-one loss

$$\ell_{0-1}(f(x), y) = \phi(f(x) \cdot y), \quad \phi(x) = \mathbb{1}[x < 0].$$

- However, the function $\phi(x) = \exp\left(-\frac{1}{2}x\right)$ is strictly convex and identifies $a(x)$.

Factorizing the Empirical Exponential Risk

Again, it is hard to optimize the empirical (exponential) risk over all base functions f_1, \dots, f_M simultaneously for large M .

AdaBoost fits f_m **sequentially**:

- The empirical exponential risk function at m -th step

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \exp \left(-\frac{1}{2} \sum_{j=1}^{m-1} f_j(X_i) Y_i - \frac{1}{2} f_m(X_i) Y_i \right) \\ &= \sum_{i=1}^n \frac{1}{n} \exp \left(-\frac{1}{2} \sum_{j=1}^{m-1} f_j(X_i) Y_i \right) \cdot \exp \left(-\frac{1}{2} f_m(X_i) Y_i \right) \\ &= \sum_{i=1}^n \tilde{w}_i^{(m)} \cdot \ell_{\exp}(f_m(X_i), Y_i) \end{aligned}$$

- We are minimizing *weighted* exponential error at each step.

- @UvA
- The weights $\tilde{w}_i^{(m)}$ are known from the previous steps.

Normalizing the Weights

- Normalizing the weights $\tilde{w}_i^{(m)}$ does not change the solution:

$$w_i^{(m)} = \frac{\tilde{w}_i^{(m)}}{\sum_{i=1}^n \tilde{w}_i^{(m)}}$$

- Normalized weights are easier to interpret as probabilities.
- Henceforth we minimize the error function at m -th step:

$$f_m = \operatorname{argmin}_{f \in \mathcal{B}} \sum_{i=1}^n w_i^{(m)} \cdot \ell_{\exp}(f(X_i), Y_i)$$

with the **normalized** weights such that $\sum_{i=1}^n w_i^{(m)} = 1$.

AdaBoost: Solving the Base Learner

Now consider any candidate base $f = \alpha g$, $g \in \mathcal{G}$ and $\alpha \geq 0$.

Because $g(X_i), Y_i \in \{-1, 1\}$, $g(X_i)Y_i = \begin{cases} -1 & g(X_i) \neq Y_i \\ 1 & g(X_i) = Y_i, \end{cases}$

$$\begin{aligned} & \ell_{\text{exp}}(f(X_i), Y_i) \\ &= \exp\left(-\frac{1}{2}\alpha g(X_i)Y_i\right) = \begin{cases} \exp\left(\frac{1}{2}\alpha\right) & g(X_i) \neq Y_i \\ \exp\left(-\frac{1}{2}\alpha\right) & g(X_i) = Y_i \end{cases} \\ &= \exp\left(\frac{1}{2}\alpha\right) \mathbb{1}[g(X_i) \neq Y_i] + \exp\left(-\frac{1}{2}\alpha\right) \underbrace{\mathbb{1}[g(X_i) = Y_i]}_{1 - \mathbb{1}[g(X_i) \neq Y_i]} \\ &= \left(\underbrace{\exp\left(\frac{1}{2}\alpha\right) - \exp\left(-\frac{1}{2}\alpha\right)}_{>0, \Leftrightarrow \alpha > 0} \right) \mathbb{1}[g(X_i) \neq Y_i] + \exp\left(-\frac{1}{2}\alpha\right) \end{aligned}$$

AdaBoost: Key Steps

Using the normalization condition $\sum_{i=1}^n w_i^{(m)} = 1$:

$$\begin{aligned} & \sum_{i=1}^n w_i^{(m)} \ell_{\exp}(f(X_i), Y_i) \\ &= \left(\underbrace{\exp\left(\frac{1}{2}\alpha\right) - \exp\left(-\frac{1}{2}\alpha\right)}_{>0 \Leftrightarrow \alpha > 0} \right) \underbrace{\sum_{i=1}^n w_i^{(m)} \mathbb{1}[g(X_i) \neq Y_i]}_{J_m(g)} + \exp\left(-\frac{1}{2}\alpha\right) \end{aligned}$$

- Fit a classifier $g_m \in \mathcal{G}$ minimizing the weighed misclassification error $J_m(g)$
- Then solve the first-order condition for α gives that

$$\alpha_m = \log \frac{1 - \epsilon_m}{\epsilon_m}, \text{ with weighted error } \epsilon_m = J_m(g_m)$$

- @UvA • $J_m(g_m) < J_m(-g_m)$ implies that $\epsilon_m < 1/2$ and $\alpha_m > 0$

To reduce computational costs, save the history of $w_i^{(m)}$ and update the data weighting coefficients along the path.

- Using the fact that $Y_i, g_m(X_i) \in \{1, -1\}$,

$$\begin{aligned}\tilde{w}_i^{(m+1)} &= \tilde{w}_i^{(m)} \exp\left(-\frac{1}{2}\alpha_m g_m(X_i) Y_i\right) \\ &\propto w_i^{(m)} \exp\left(-\frac{1}{2}\alpha_m (1 - 2\mathbb{1}(g_m(X_i) \neq Y_i))\right) \\ &\propto w_i^{(m)} \exp(\alpha_m \mathbb{1}(g_m(X_i) \neq Y_i)) \cdot \exp\left(-\frac{1}{2}\alpha_m\right)\end{aligned}$$

- Update and normalize:

$$\begin{aligned}w_i^{(m+1)} &\leftarrow w_i^{(m)} \exp(\alpha_m \mathbb{1}(g_m(X_i) \neq Y_i)), \\ w_i^{(m+1)} &\leftarrow \frac{w_i^{(m+1)}}{\sum_{i=1}^n w_i^{(m+1)}}\end{aligned}$$

Points that are **misclassified** are given **greater** weight when used to train the next classifier in the sequence.

Remark: LogitBoost

- You may consider the logistic loss

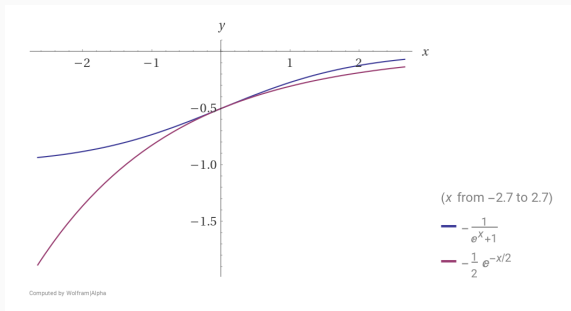
$$\ell_{\text{Logit}}(f(x), y) = \ell_{\text{CE}} \left(\underbrace{\begin{pmatrix} \sigma(f(x)) \\ \sigma(-f(x)) \end{pmatrix}}_{\text{Bivariate Regression Function}}, \begin{pmatrix} \mathbb{1}[y = 1] \\ \mathbb{1}[y = -1] \end{pmatrix} \right) \\ = \phi(f(x) \cdot y)$$

with

$$\phi(x) = \log(1 + \exp(-x)).$$

- ... which gives the 'LogitBoost' algorithm in Friedman et al. (2000).

Comparing the derivative $\phi'(x)$ for AdaBoost (purple) and LogitBoost (blue):



LogitBoost could be less sensitive to outliers

Friedman et al. (2000) conclude in their simulation studies:

It is likely that when the shrinkage parameter [η here] is carefully tuned ..., there would be little performance differential between them.