

Lecture 16 - Mar 4

Boosting

AdaBoost.M1

Reading

- [Elements of Statistical Learning](#)
 - Ch 10: Boosting and Additive Trees
- QuantDare (2016). What is the difference between Bagging and Boosting?
<https://quandare.com/what-is-the-difference-between-bagging-and-boosting/>

Motivation: Combine the outputs of many
"weak" classifiers to form a powerful ensemble



better than
random guessing

Upcoming Deadlines

Homework 4 (Mar 15)

Boosting

Motivation: Combine the outputs of many "weak" classifiers to form a powerful "committee"

AdaBoost.M1 (Freund + Schapire 1997)

Let $Y \in \{-1, 1\}$ with predictors X + classifier $G(x) \in \{-1, 1\}$

Error rate on training sample:

$$\bar{err} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{y_i \neq G(x_i)\}}$$

and the error rate on future predictions is $E_{xy} [\mathbb{1}_{\{y \neq G(x)\}}]$

A weak classifier has error rate slightly better than random guessing

Boosting: sequentially apply weak classification algorithm to repeatedly modified version of the data to produce weak classifiers $G_m(x)$, $m=1, \dots, M$

Combine predictions through a weighted majority vote

$$G(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m G_m(x) \right)$$

"Inaction algorithm. They

$$G(x) = \text{sign}(\sum_{m=1}^M$$

with weights computed by the boosting algorithm. They give higher influence to more accurate classifiers

Data modifications at each boosting step apply weights w_1, \dots, w_N to training observations (x_i, y_i) . Initially, all $w_i = \frac{1}{N}$. In successive iterations, observation weights are modified:

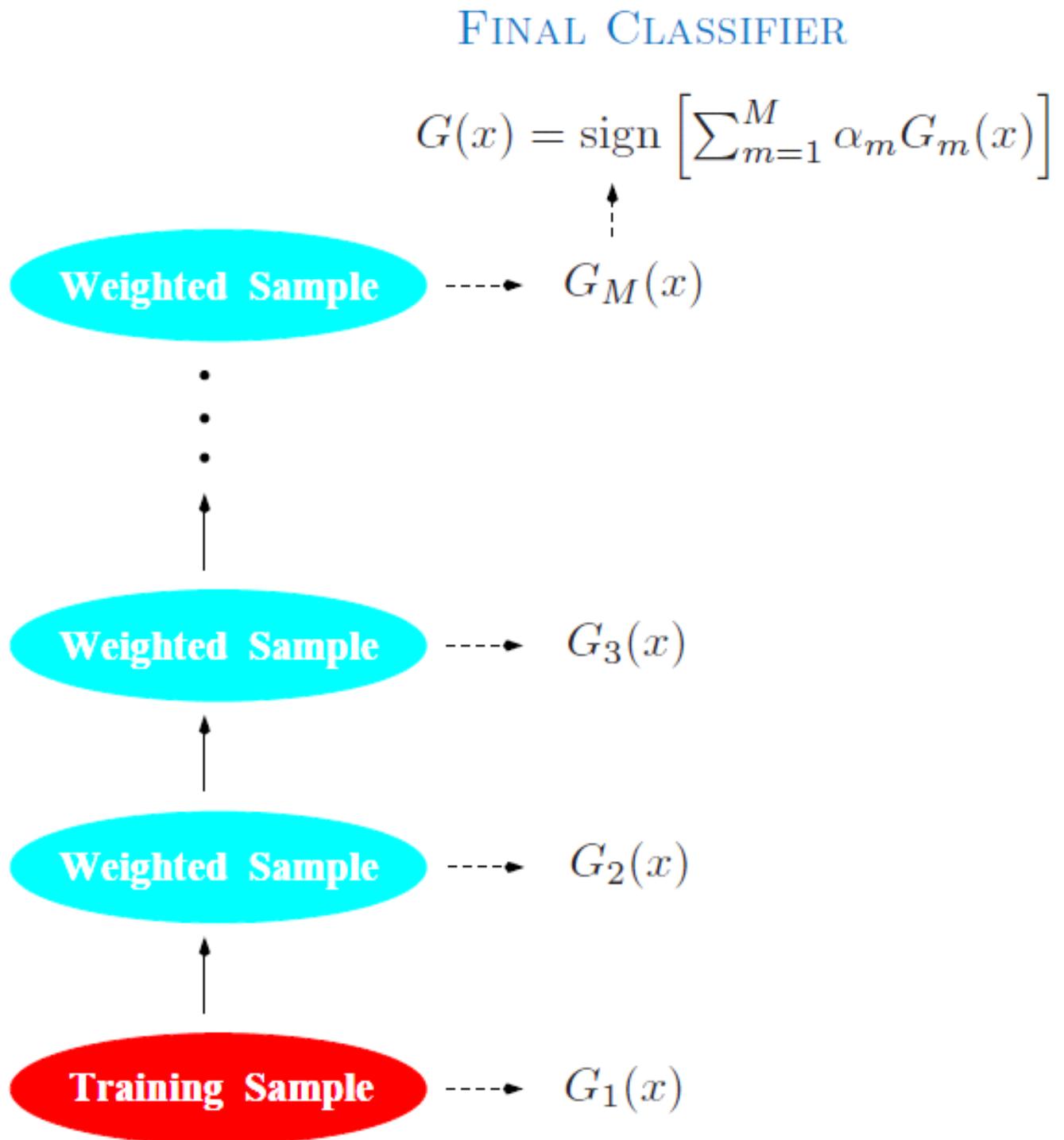
misclassified observations $\rightarrow w_i$ increases

correctly classified obs $\rightarrow w_i$ decreases

\Rightarrow difficult points get more influence in new models fitted with weights

\Rightarrow classifiers forced to focus on difficult points

Boosting Diagram



AdaBoost.M1 Algorithm

(1) Set $w_i = \frac{1}{N}, i=1, \dots, N$

(2) For $M=1$ to M :

(a) fit classifier $G_m(x)$ to data with w_i 's

(b) Compute

$$\text{err}_m = \frac{\sum_{i=1}^N w_i \mathbb{1}_{\{y_i \neq G_m(x)\}}}{\sum_{i=1}^N w_i}$$

fraction of weight misclassified

(c) Compute

$$\alpha_m = \log\left(\frac{1 - \text{err}_m}{\text{err}_m}\right)$$

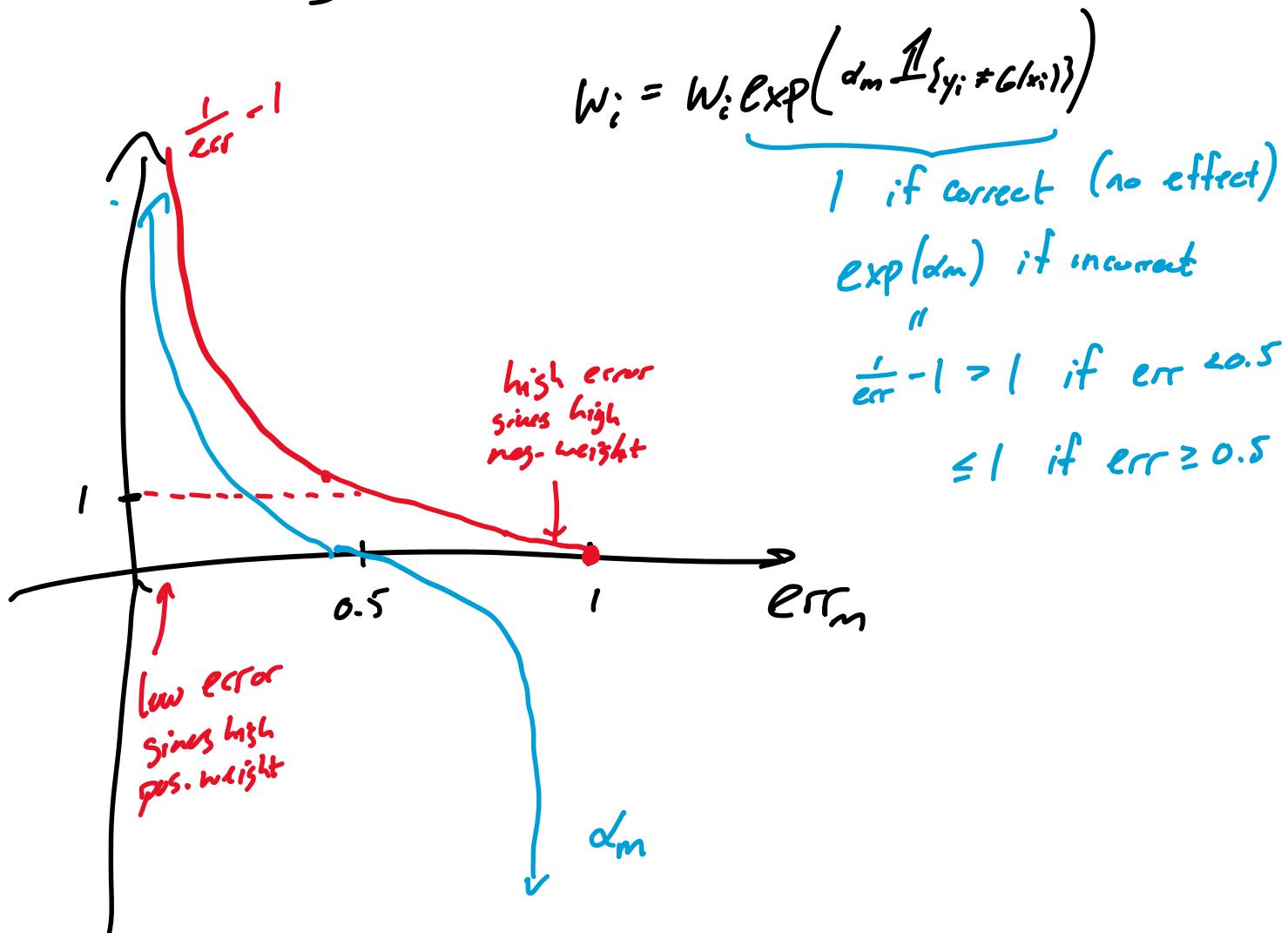
(d) Set $w_i = w_i \exp(\alpha_m \mathbb{1}_{\{y_i \neq G_m(x)\}})$

(3) Output $G(x) = \text{sign}\left(\sum_{m=1}^M \alpha_m G_m(x)\right)$

This is "discrete AdaBoost"; "Real AdaBoost" works to predict
 $y \in [-1, 1]$.

$$\alpha_m = \log\left(\frac{1 - \text{err}_m}{\text{err}_m}\right) = \log\left(\frac{1}{\text{err}_m} - 1\right)$$

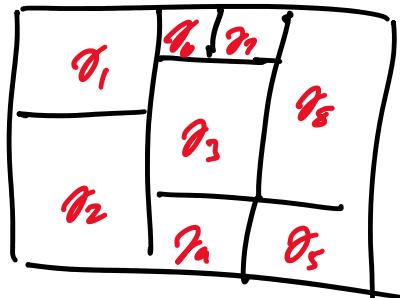
1. 1. ...



Trees split the feature space into rectangles

$$R_1, R_2, \dots, R_J$$

If $x \in R_j$, then $f(x) = g_j$



We can express the tree as

$$T(x; \theta) = \sum_{j=1}^J \gamma_j \mathbb{1}_{\{x \in R_j\}}$$

with parameters $\Theta = \{R_j, g_j\}_{j=1, \dots, J}$

We train T by trying to minimize

$$\arg \min_{\Theta} \sum_{j=1}^J \underbrace{\sum_{x_i \in R_j} L(y_i, g_j)}_{\text{Sum of losses for all } x_i \in R_j}$$

This is hard... so we rely on Suboptimal solutions by

① Find g_i given R_i (\bar{y}_i or model class)

② Find R_i 's by greedy, recursive splitting

$$\hat{\theta} = \operatorname{argmin}_{\theta} \sum_{i=1}^N \tilde{L}(y_i, T(x_i; \theta))$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, T(x_i; \theta))$$

Boosted tree model : $f_M(x) = \sum_{m=1}^M T(x; \theta_m)$

$$\hat{\theta}_M = \underset{\theta_m}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, f_{M-1}(x_i) + T(x_i; \theta_m))$$

$$\hat{\theta}_M = \{R_{jm}, \gamma_{jm}\}$$

frequently, we have to
approx

Given R_{jm} , finding optimal γ_{jm} is easy..-

$$\hat{\gamma}_{jm} = \underset{\gamma_{jm}}{\operatorname{argmin}} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma_{jm})$$

AdaBoost \rightarrow two-class classifier with exponential loss $e^{-y_i f_m(x_i)}$

$$\hat{\theta}_m = \underset{\theta_m}{\operatorname{argmin}} \sum_{i=1}^n w_i^{(m)} \exp(-y_i T(x_i; \theta_m))$$

$$\hat{\gamma}_{jm} = \ln \left(\frac{\sum_{x_i \in R_{jm}} w_i^{(m)} \mathbb{1}_{\{y_i=1\}}}{\sum_{x_i \in R_{jm}} w_i^{(m)} \mathbb{1}_{\{y_i=-1\}}} \right)$$

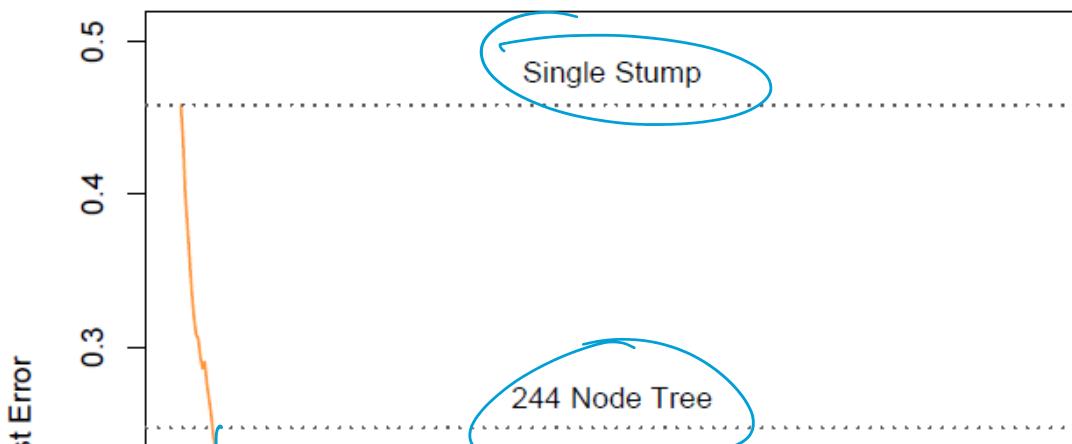
AdaBoost Performance

The power of AdaBoost to dramatically increase the performance of even a very weak classifier is illustrated in Figure 10.2. The features X_1, \dots, X_{10} are standard independent Gaussian, and the deterministic target Y is defined by

$$Y = \begin{cases} 1 & \text{if } \sum_{j=1}^{10} X_j^2 > \chi_{10}^2(0.5), \\ -1 & \text{otherwise.} \end{cases} \quad (10.2)$$

Const.

Here $\chi_{10}^2(0.5) = 9.34$ is the median of a chi-squared random variable with 10 degrees of freedom (sum of squares of 10 standard Gaussians). There are 2000 training cases, with approximately 1000 cases in each class, and 10,000 test observations. Here the weak classifier is just a “stump”: a two terminal-node classification tree. Applying this classifier alone to the training data set yields a very poor test set error rate of 45.8%, compared to 50% for random guessing. However, as boosting iterations proceed the error rate steadily decreases, reaching 5.8% after 400 iterations. Thus, boosting this simple very weak classifier reduces its prediction error rate by almost a factor of four. It also outperforms a single large classification tree (error rate 24.7%). Since its introduction, much has been written to explain the success of AdaBoost in producing accurate classifiers. Most of this work has centered on using classification trees as the “base learner” $G(x)$, where improvements are often most dramatic. In fact, Breiman (NIPS Workshop, 1996) referred to AdaBoost with trees as the “best off-the-shelf classifier in the world” (see also Breiman (1998)). This is especially the case for data-mining applications, as discussed more fully in Section 10.7 later in this chapter.



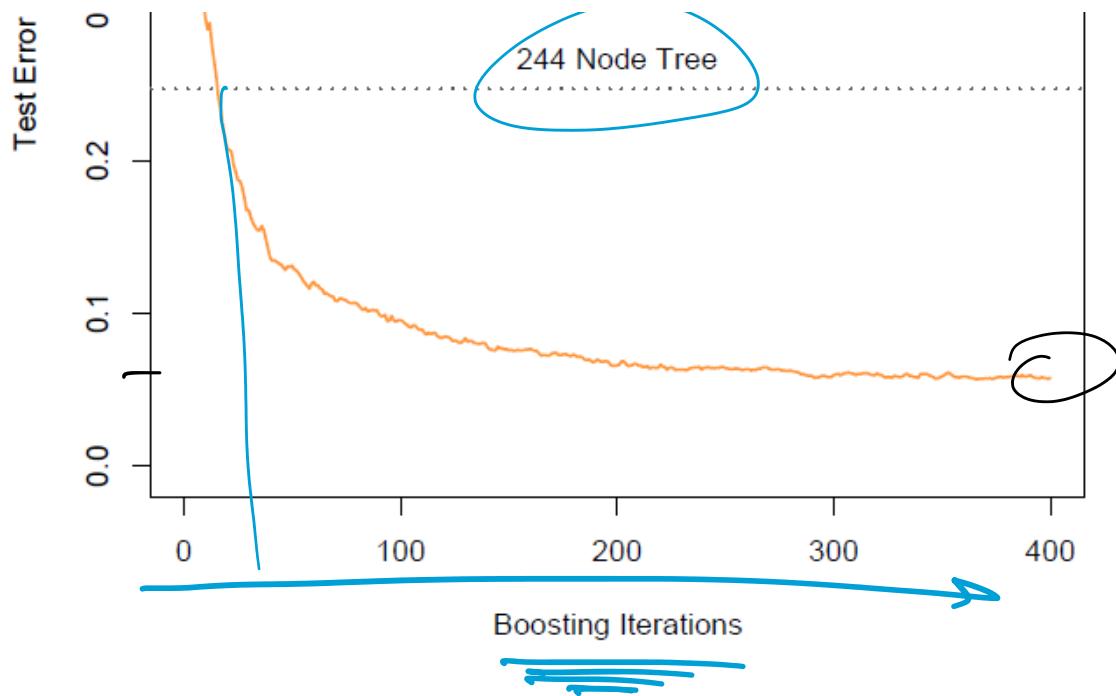


FIGURE 10.2. Simulated data (10.2): test error rate for boosting with stumps, as a function of the number of iterations. Also shown are the test error rate for a single stump, and a 244-node classification tree.

Bagging vs. Boosting

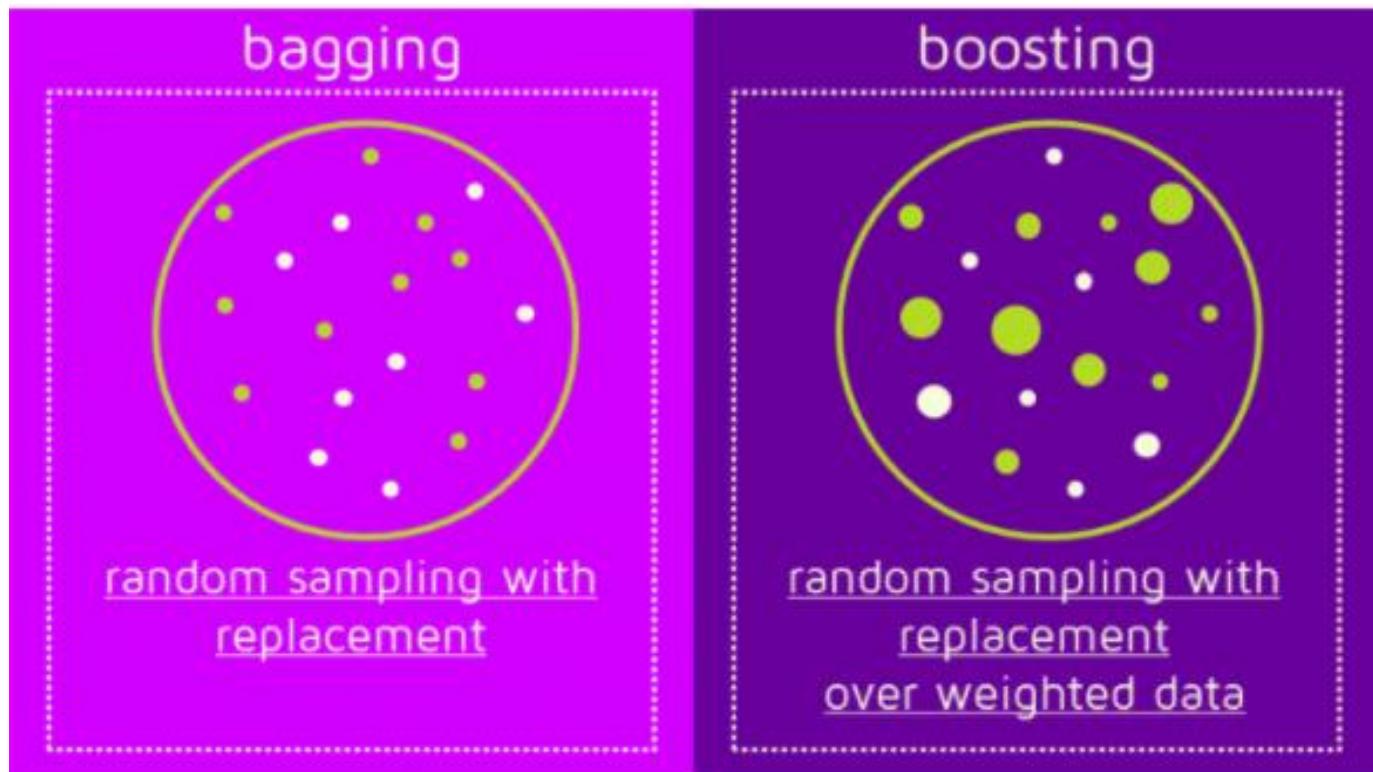
(from the QuantDare blog post referenced above)

The main causes of error in learning are due to **noise, bias, and variance**. Ensemble helps to minimize these factors.

These methods are designed to improve the stability and the accuracy of ML algorithms.

Combinations of classifiers **decrease variance**, especially in the case of unstable classifiers, and may produce a more reliable classification than a single classifier.

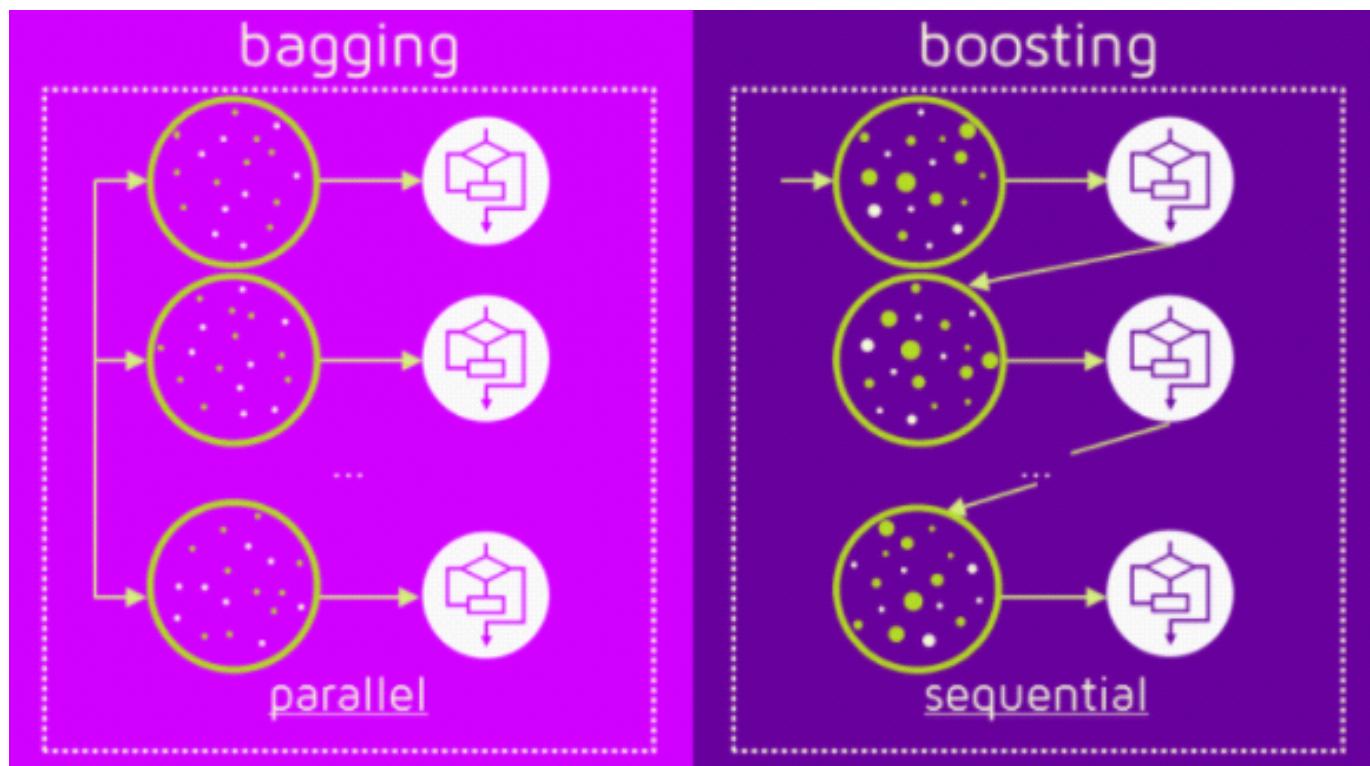
To use Bagging or Boosting you must select a base learner algorithm. For example, if we choose a classification tree, they use a pool of trees.



N new training data sets are produced by **bootstrap samples** from the original set.

Bagging: any element has the **same probability** to appear in a new dataset.

Boosting: observations are **weighted** so some are more likely to be in new datasets

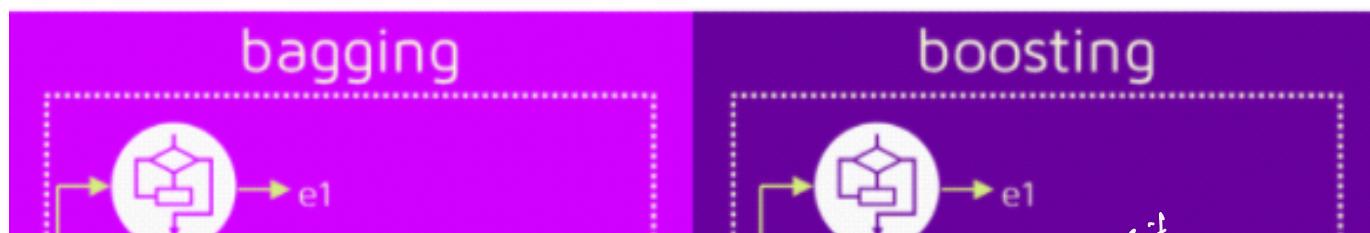


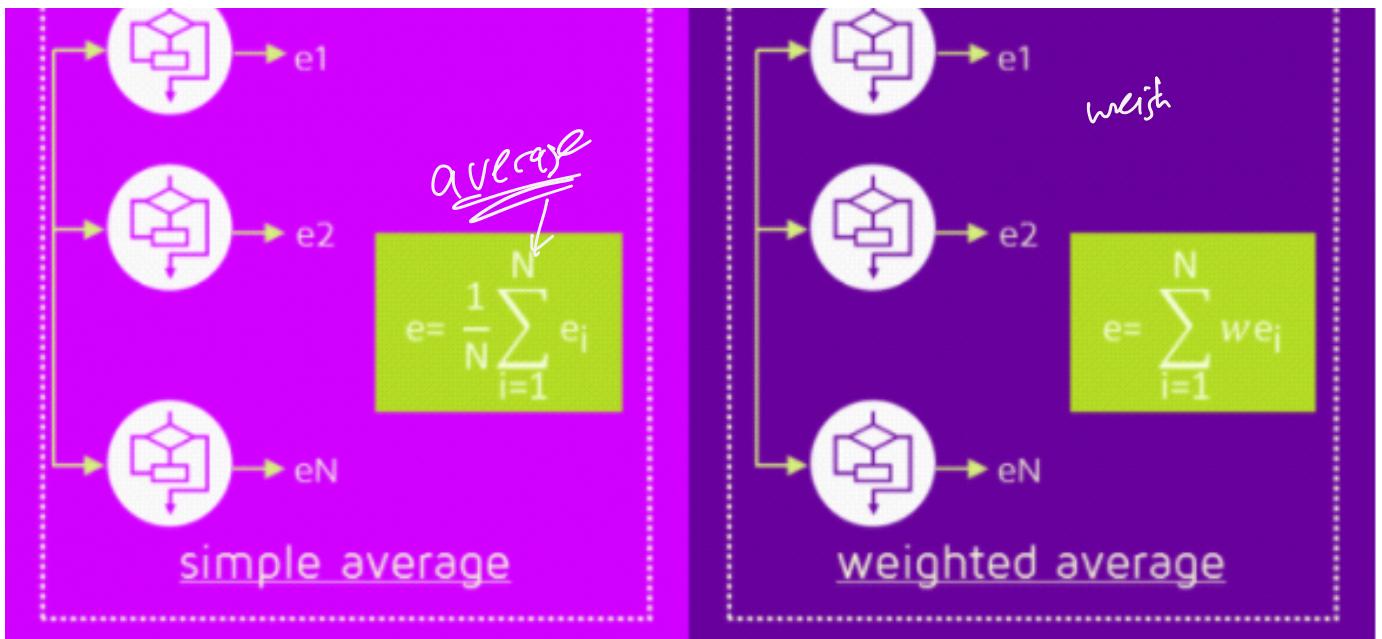
Boosting: each classifier is trained on data, taking into account the previous classifiers' success. After each training step, the weights are redistributed.

Misclassified data increases its weights to emphasize the most difficult cases. In this way, subsequent learners will focus on them during their training.

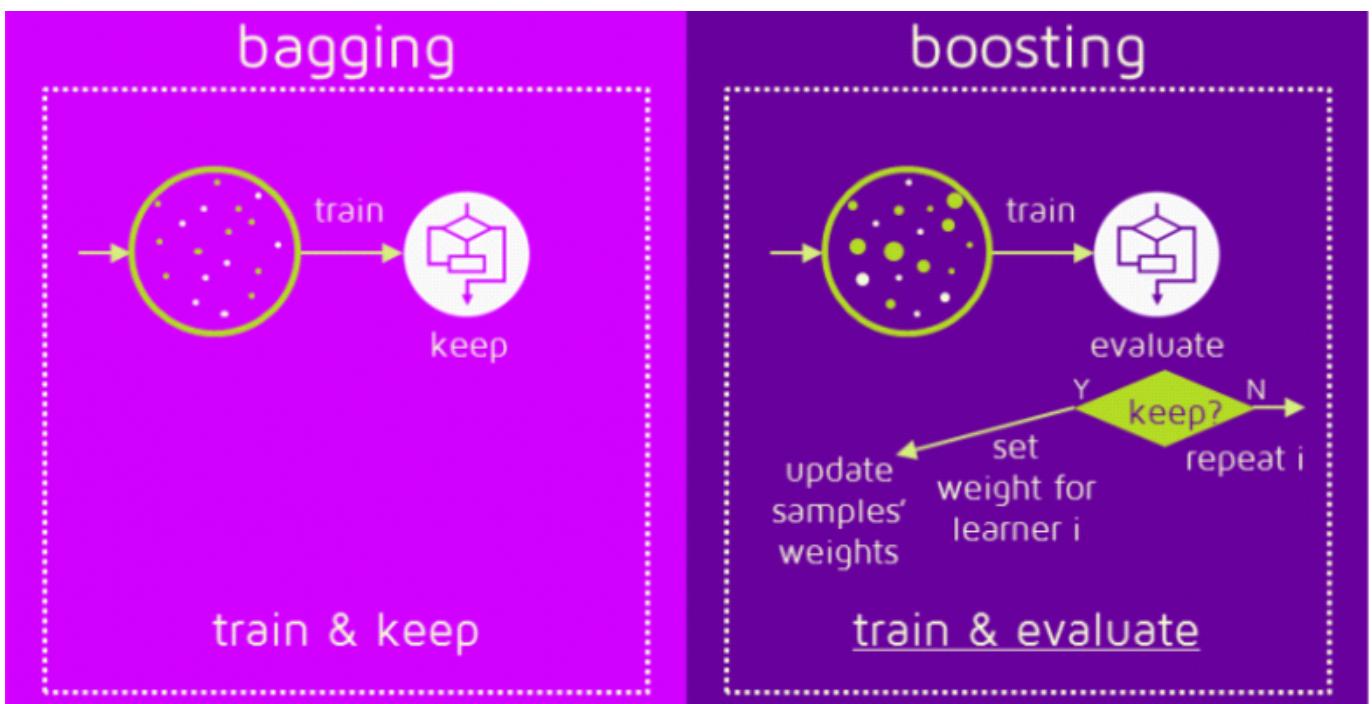
Bagging: result is average the responses of the N classifiers (majority vote).

Boosting: assigns a second set of weights for the N classifiers (weighted vote)





New versions of **AdaBoost**: a model is rejected if error is **worse than random guessing**. Otherwise, repeated until achieving a learner better than a random guess.



In boosting, several alternatives exist to determine the weights to use in the next

training step and in the classification stage.

Most popular: **AdaBoost**, LPBoost, **XGBoost**, **GradientBoost**, BrownBoost

Should I used Bagging or Boosting?

They both **decrease variance**, the resulting in a model with higher stability.

Boosting can generate a combined model with lower errors as it optimizes the advantages and reduces pitfalls of the single model to **reduce underfitting**.
(Bagging doesn't help.)

Bagging is good at **reducing overfitting**.

(Boosting doesn't help ~~much~~)

Similarities	Differences
Both are ensemble methods to get N learners from 1 learner...	... but, while they are built independently for Bagging, Boosting tries to add new models that do well where previous models fail.
Both generate several training data sets by random sampling...	... but only Boosting determines weights for the data to tip the scales in favor of the most difficult cases.
Both make the final decision by averaging the N learners (or taking the majority of them)...	... but it is an equally weighted average for Bagging and a weighted average for Boosting, more weight to those with better performance on training data.
Both are good at reducing variance and provide higher stability...	... but only Boosting tries to reduce bias. On the other hand, Bagging may solve the over-fitting problem, while Boosting can increase it.

Boosting Fits an Additive Model

The success of boosting is really not very mysterious. The key lies in expression (10.1). Boosting is a way of fitting an additive expansion in a set of elementary “basis” functions. Here the basis functions are the individual classifiers $G_m(x) \in \{-1, 1\}$. More generally, basis function expansions take the form

Basis expansion models

$$f(x) = \sum_{m=1}^M \beta_m b(x; \gamma_m)$$

... +

β_m

$b(x; \gamma_m)$'s

where $\beta_m, m = 1, 2, \dots, M$ are the expansion coefficients, and $b(x; \gamma) \in \mathbb{R}$ are usually simple functions of the multivariate argument x , characterized by a set of parameters γ . We discuss basis expansions in some detail in Chapter 5.

Typically fit by minimizing loss

$$\min_{\{\beta_m, \gamma_m\}_m} \sum_{i=1}^n L(y_i, \sum_{m=1}^M \beta_m b(x_i; \gamma_m))$$

This can be expensive depending on $L + b$'s, but a feasible alternative when it is feasible to rapidly solve the subproblem of fitting a single basis function

$$\min_{\beta, \gamma} \sum_{i=1}^n L(y_i, \beta b(x_i; \gamma))$$

$$\min_{\beta, \gamma} \sum_{i=1} L(y_i, f_{\text{base}}(x_i))$$

Idea: Sequentially add new basis functions without adjusting parameters + coefficients of those already added

Algorithm 10.2

- (1) Init $f_0(x) = 0$
- (2) For $m = 1$ to M :
 - (a) Compute $(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^n L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma))$
 - (b) Set $\underline{f_m(x)} = \underline{f_{m-1}(x)} + \boxed{\beta_m b(x; \gamma_m)}$

one more
tree

Adds one basis function greedily

AdaBoost.M1 = forward stagewise additive modeling with loss

$$L(y, f(x)) = \exp(-yf(x))$$

(Simple proof in ESL 10.4)

For AdaBoost the basis functions are the individual classifiers $G_m(x) \in \{-1, 1\}$. Using the exponential loss function, one must solve

$$(\beta_m, G_m) = \arg \min_{\beta, G} \sum_{i=1}^N \exp[-y_i(f_{m-1}(x_i) + \beta G(x_i))] \quad \cancel{\text{A}}$$

for the classifier G_m and corresponding coefficient β_m to be added at each step. This can be expressed as

$$(\beta_m, G_m) = \arg \min_{\beta, G} \sum_{i=1}^N w_i^{(m)} \exp(-\beta y_i G(x_i)) \quad (10.9)$$

with $w_i^{(m)} = \exp(-y_i f_{m-1}(x_i))$. Since each $w_i^{(m)}$ depends neither on β nor $G(x)$, it can be regarded as a weight that is applied to each observation. This weight depends on $f_{m-1}(x_i)$, and so the individual weight values change with each iteration m .

The solution to (10.9) can be obtained in two steps. First, for any value of $\beta > 0$, the solution to (10.9) for $G_m(x)$ is

$$G_m = \arg \min_G \sum_{i=1}^N w_i^{(m)} I(y_i \neq G(x_i)), \quad (10.10)$$

which is the classifier that minimizes the weighted error rate in predicting y . This can be easily seen by expressing the criterion in (10.9) as

$$e^{-\beta} \cdot \sum_{y_i=G(x_i)} w_i^{(m)} + e^{\beta} \cdot \sum_{y_i \neq G(x_i)} w_i^{(m)},$$

which in turn can be written as

$$(e^{\beta} - e^{-\beta}) \cdot \sum_{i=1}^N w_i^{(m)} I(y_i \neq G(x_i)) + e^{-\beta} \cdot \sum_{i=1}^N w_i^{(m)}. \quad (10.11)$$

Plugging this G_m into (10.9) and solving for β one obtains

$$\beta_m = \frac{1}{2} \log \frac{1 - \text{err}_m}{\text{err}_m}, \quad (10.12)$$

where err_m is the minimized weighted error rate

$$\text{err}_m = \frac{\sum_{i=1}^N w_i^{(m)} I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i^{(m)}}. \quad (10.13)$$

The approximation is then updated

$$f_m(x) = f_{m-1}(x) + \beta_m G_m(x),$$

which causes the weights for the next iteration to be

$$w_i^{(m+1)} = w_i^{(m)} \cdot e^{-\beta_m y_i G_m(x_i)}. \quad (10.14)$$

Using the fact that $-y_i G_m(x_i) = 2 \cdot I(y_i \neq G_m(x_i)) - 1$, (10.14) becomes

$$w_i^{(m+1)} = w_i^{(m)} \cdot e^{\alpha_m I(y_i \neq G_m(x_i))} \cdot e^{-\beta_m}, \quad (10.15)$$

where $\alpha_m = 2\beta_m$ is the quantity defined at line 2c of AdaBoost.M1 (Algorithm 10.1). The factor $e^{-\beta_m}$ in (10.15) multiplies all weights by the same value, so it has no effect. Thus (10.15) is equivalent to line 2(d) of Algorithm 10.1.

One can view line 2(a) of the AdaBoost.M1 algorithm as a method for approximately solving the minimization in (10.11) and hence (10.10). Hence we conclude that AdaBoost.M1 minimizes the exponential loss criterion (10.8) via a forward-stagewise additive modeling approach.