Boosting

Statistical Machine Learning (ENGG*6600*08)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Fall 2023



- Boosting is a <u>meta algorithm</u> which can be used with any model (classifier, regression, etc).
- For binary classification, for example, if we use boosting with a <u>classifier</u> even slightly better than flipping a coin, we will have a strong <u>classifier</u> (we will explain the reason later). Thus, we can say boosting makes the estimation or classification very strong. In other words, boosting addresses the question whether a strong classifier can be obtained from a set of weak classifiers [1, 2].

- The idea of boosting is to learn <u>k models in a hierarchy</u> where every model gives more attention (larger weight) to the instances misclassified (or estimated very badly) by the previous model.
- Finally, the overall estimation or classification is a weighted summation (average) of the k estimations. For an instance x, we have:

$$\widehat{f}(\mathbf{x}) = \sum_{j=1}^{k} \alpha_j h_j(\mathbf{x}).$$
(1)

If the model is <u>classifier</u>, we should probably use sign function:

$$\widehat{f}(\mathbf{x}) = \sup_{j=1}^{k} (\sum_{j=1}^{k} \alpha_j h_j(\mathbf{x})), \quad (2)$$

which is equivalent to **majority voting** among the trained classifiers.

training instances



- Different methods have been proposed for boosting, one of the most well-known ones is AdaBoost (Adaptive Boosting) (1996) [3].
- The algorithm of AdaBoost for binary classification is shown below.

1 Initialize
$$w_i = 1/n, \forall i \in \{1, \dots, n\}$$

2 for j from 1 to k do
3 $h_j(\mathbf{x}) = \arg \min L_j$
4 $\alpha_j = \log(\frac{1-L_j}{L_j})$
5 $w_i = w_i \exp(\alpha_j \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)))$

Algorithm : The AdaBoost Algorithm

- In this algorithm, L_j is the cost function minimized in the j-th model h_j, the I(.) is the indicator function which is one and zero if its condition is and is not satisfied, respectively, and w_i is the weight associated to the i-th instance for weighting it as the input to the next layer of boosting.
- Note that the cost in AdaBoost is:

$$(3)$$

which makes sense because it gets larger if the observations of more instances are estimated incorrectly.

1 Initialize
$$w_i = 1/n, \forall i \in \{1, \dots, n\}$$

2 for j from 1 to k do
3 $h_j(\mathbf{x}) = \arg \min L_j$
4 $\alpha_j = \log(\frac{1-L_j}{L_j})$
5 $w_i = w_i \exp(\alpha_j \widehat{\mathbb{I}(y_i \neq h_j(\mathbf{x}_i))})$

Algorithm : The AdaBoost Algorithm

- Here, we can have several cases which help us understand the interpretation of the AdaBoost algorithm:
- If an instance is correctly classified, the I(y_i ≠ h_i(x_i)) is zero and thus the w_i will be still w_i without any change. This makes sense because the correctly classified instance should not gain a significant weight in the next layer of boosting.



Algorithm : The AdaBoost Algorithm

- If an instance is misclassified, the $\mathbb{I}(y_i \neq h_j(\mathbf{x}_i))$ is one. In this case, we can have two sub-cases:
 - If the classifier which classified that instance was a bad classifier, its cost would be like flipping a coin, i.e., $L_j = 0.5$. Therefore, we will have $\alpha_j = \log(1) = 0$ and again the w_i will still be w_i without any change. This makes sense because we cannot trust the bad classifier whether the instance is correctly or incorrectly classified and thus we should not make any decision based on that.
 - If the classifier which classified that instance was a good classifier, then we have $L_j \neq 0.5$ and as we also have $\mathbb{I}(y_i \neq h_j(\mathbf{x}_i)) = 1$, the weight will change as
 - $(\overline{w_i := w_i} \exp(\alpha_i))$. This also is intuitive because the previous model in the boosting was a good classifier and we can trust it and that good classifier could not classify the instance correctly. Therefore, we should notice that instance more in the next model in the boosting hierarchy.



• Additive models [4] can be used to explain why boosting works [5, 6]. In an additive model, we map the data as $x \mapsto \phi_j(x), \forall j \in \{1, \dots, k\}$ and then add them using some weights β_j 's:

$$\phi(\mathbf{x}) = \sum_{j=1}^{k} \beta_j \, \phi_j(\mathbf{x}). \tag{4}$$

- A well-known example of the additive model is Radial Basis Function (RBF) neural network (here with k hidden nodes) which uses Gaussian mappings [7, 8].
- Now, consider a cost function for an instance as: $L(y, h(x)) := \exp(-y h(x)), \quad (5)$

where y is the observation or label for x and h(x) is the model's estimation of y. This cost is intuitive because when the instance is misclassified, the signs of y and h(x) will be different (so y h(x) < 0) and the cost will be large, while in case of correct classification, the signs are similar (so y h(x) > 0) and the cost is small.

• If we add up the cost over the *n* training instances, we have:

$$L_t(y, h(\mathbf{x})) := \sum_{i=1}^n \exp(-y_i h(\mathbf{x}_i)),$$
(6)

where L_t denotes the total cost.

• In Eq. (4), $\phi(\mathbf{x}) = \sum_{j=1}^{k} \beta_j \phi_j(\mathbf{x})$, if we rename the mapping to $h(\mathbf{x})$, which is the model used in boosting, we will have:

$$\bigstar \qquad \bigstar \qquad \begin{pmatrix} h(\mathbf{x}) \\ j=1 \end{pmatrix} = \sum_{j=1}^{k} \beta_j h_j(\mathbf{x}).$$
(7)

• We can write this expression as a **forward stage-wise additive model** [5, 6] in which we work with the models one by one where we add up the previously worked models:

$$\begin{cases} \underbrace{f_{q-1}(\mathbf{x})}_{j=1} = \sum_{j=1}^{q-1} \beta_j h_j(\mathbf{x}), \\ \underbrace{f_q(\mathbf{x})}_{j=1} = \underbrace{f_{q-1}(\mathbf{x})}_{q-1} + \underbrace{\beta_q h_q(\mathbf{x})}_{q-1}, \\ \underbrace{q \leq k}_{j=1}, \end{cases}$$
(8)

where $h(\mathbf{x}) = f_k(\mathbf{x})$. • Therefore, minimizing the cost, i.e., Eq. (6), $L_t(y, h(\mathbf{x})) := \sum_{i=1}^n \exp(-y_i h(\mathbf{x}_i))$, for the <u>j-th model</u> in the <u>additive manner</u> is: $\min_{\substack{j:j,h_j \ i=1}} \sum_{i=1}^n \exp(-y_i f_{j-1}(\mathbf{x}_i) + \beta_j h_j(\mathbf{x}_i))) = \min_{\beta_j,h_j} \sum_{i=1}^n \exp(-y_i f_{j-1}(\mathbf{x}_i)) \exp(-y_i \beta_j h_j(\mathbf{x}_i)).$

Boosting

We had:

$$\min_{\beta_j,h_j} \sum_{i=1}^n \exp\left(-y_i \left[f_{j-1}(\boldsymbol{x}_i) + \beta_j h_j(\boldsymbol{x}_i)\right]\right) = \min_{\beta_j,h_j} \sum_{i=1}^n \exp(-y_i f_{j-1}(\boldsymbol{x}_i)) \exp(-y_i \beta_j h_j(\boldsymbol{x}_i)).$$

• The first term is a constant with respect to β_j and h_j so we name it by w_i :

$$\underline{w_i} := \underline{\exp(-y_i f_{j-1}(\boldsymbol{x}_i))}. \tag{10}$$

 $\mathcal{U}_{\mathcal{U}}$

Thus:

• As in binary AdaBoost, we have
$$\pm 1$$
 for y_i and h_j , we can say:

$$\begin{array}{c}
\underset{\beta_j,h_j}{\underset{i=1}{n}} \tilde{w}_i \exp(-y_i \beta_j h_j(\mathbf{x}_i)). \\
\underset{\beta_j,h_j}{\underset{i=1}{n}} \exp(-\beta_j) \sum_{i=1}^n w_i \mathbb{I}(y_i = h_j(\mathbf{x}_i)) + \exp(\beta_j) \sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)) \\
\overset{(a)}{\underset{\beta_j,h_j}{\underset{i=1}{n}}} \exp(-\beta_j) \sum_{i=1}^n w_i - \exp(-\beta_j) \sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)) + \exp(\beta_j) \sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)), \\
\end{aligned}$$
where (a) is because:

$$\begin{array}{c}
\underset{i=1}{\underset{j=1}{n}} w_i \mathbb{I}(y_i = h_j(\mathbf{x}_i)) = \sum_{i=1}^n w_i - \sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)). \\
\end{array}$$

• We had:

$$\underbrace{\min_{\beta_j,h_j} \exp(-\beta_j) \sum_{i=1}^n w_i \bigoplus_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)) + \underbrace{\exp(\beta_j) \sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i))}_{i=1}$$

• For the sake of minimization, we take the derivative:

$$\underbrace{ \frac{\partial L_t}{\partial \beta_j} = \exp(-\beta_j) \sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)) }_{i=1} \underbrace{ \exp(-\beta_j) \sum_{i=1}^n w_i \mathbb{I}($$

which gives:

$$\Rightarrow (\exp(-\beta_j) + \exp(\beta_j)) \times \underbrace{\sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i))}_{\sum_{i=1}^n w_i} = \exp(-\beta_j)$$

$$\stackrel{(3)}{\Rightarrow} (\exp(-\beta_j) + \exp(\beta_j)) \underbrace{L_j}_{j=1} = \exp(-\beta_j)$$

$$\Rightarrow \underbrace{L_j = \frac{\exp(-\beta_j)}{\exp(-\beta_j) + \exp(\beta_j)}}_{(j=1)} \Rightarrow \underbrace{\exp(2\beta_j) = \frac{1-L_j}{L_j}}_{(j=1)} \Rightarrow 2\beta_j = \underbrace{\log(\frac{1-L_j}{L_j})}_{(11)}$$

where (a) is because of the formula of α_i in the Algorithm.

Theory Based on Additive Models
• According to Eqs. (8), (9), and (10),
$$W_{i} = W_{i}^{c} \mathcal{M}$$

 $\int f_{q-1}(\mathbf{x}) = \sum_{j=1}^{q-1} \beta_{j} h_{j}(\mathbf{x}),$
 $f_{q}(\mathbf{x}) = f_{q-1}(\mathbf{x}) + \beta_{q} h_{q}(\mathbf{x}), \quad q \leq k,$
 $w_{i} := \exp(-y_{i} f_{j-1}(\mathbf{x})).$
we have:
 \mathbf{x}
 $w_{i} := w_{i} \exp(-\underline{y_{i}}\beta_{j} h_{j}(\mathbf{x}_{i})).$ (12)
As we have $y_{i} h_{j}(\mathbf{x}_{i}) = \pm 1$, we can say:
 $-\overline{y_{i}} h_{j}(\mathbf{x}_{i}) = 2 \mathbb{I}(y_{i} \neq h_{j}(\mathbf{x}_{i})) - 1.$ (13)
According Eqs. (11), $\alpha_{j} = 2\beta_{j}$, (12), and (13), we have:
 $\overline{W_{i} := w_{i} \exp(\alpha_{j} \mathbb{I}(y_{i} \neq h(\mathbf{x}_{i})))} \exp(-\beta_{j}),$ (14)

which is equivalent to the formula of w_i in the Algorithm with a factor of $\exp(-\beta_j)$. This factor does not have impact on whether the instance is correctly classified or not.



• There is an **upper bound on the generalization error** of boosting [9]. In binary boosting, we have ± 1 for y_i and also the sign of $\hat{f}(\mathbf{x}_i)$ is important; therefore, $y_i \hat{f}(\mathbf{x}_i) < 0$ means that we have error for estimating the *i*-th instance. Thus, for an error, we have:

$$\mathbf{f} \qquad \mathbf{y}_i \, \widehat{f}(\mathbf{x}_i) \leq \theta, \tag{15}$$

for a $\theta > 0$. Recall the Eq. (1):

$$\widehat{f}(\mathbf{x}) = \sum_{j=1}^{k} \alpha_j h_j(\mathbf{x}).$$

We can normalize this equation because the sign of it is important:

*

$$\bigstar \qquad \widehat{f}(\mathbf{x}_i) = \frac{\sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i)}{\sum_{j=1}^k \alpha_j}.$$
(16)

According to Eqs. (15) and (16), we have:

$$\underbrace{y_i \hat{f}(\mathbf{x}_i) \leq \theta}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) \leq \theta \sum_{j=1}^k \alpha_j}_{\substack{j=1 \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j \geq 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha_j + \theta \sum_{j=1}^k \alpha_j = 0}_{\substack{i \neq j \\ j=1}} \underbrace{\varphi_j \sum_{j=1}^k \alpha$$



and Eq. (17), we have (take $\underline{a} = 1$ and the exponential term as X in Markov's inequality):

$$\mathbb{P}(y_i \, \hat{f}(\mathbf{x}_i) \le \theta) = \mathbb{P}\left(\exp\left(-y_i \sum_{j=1}^k \alpha_j \, h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j\right) \ge 1\right)$$

$$\stackrel{(18)}{\le} \mathbb{E}\left(\exp\left(-y_i \sum_{j=1}^k \alpha_j \, h_j(\mathbf{x}_i) + \theta \sum_{j=1}^k \alpha_j\right)\right)$$



where (a) is because the expectation is with respect to the data, i.e., x_i and y_i and (b) is according to the definition of expectation.

• Recall the formula of
$$w_i$$
 in the Algorithm:

$$\begin{array}{c} \hline w_i^{(j+1)} = w_i^{(j)} \exp\left(\alpha_j \mathbb{I}(y_i \neq h_j(\mathbf{x}_i))\right), \\ \hline w_i^{(j+1)} = w_i^{(j)} \exp\left(-y_i \alpha_j h_j(\mathbf{x}_i)\right), \\ \hline w_i^{(j+1)} = w_i^{(j)} \exp\left(-y_i \alpha_j h_j(\mathbf{x}_i)\right), \\ \hline w_i^{(j+1)} = \pm 1 \text{ and } h_i(\mathbf{x}_i) = \pm 1. \end{array}$$

• It is not harmful to AdaBoost if we use the normalized weights:

where:





• We found:

$$\Psi \quad \mathbb{P}(y_i \, \hat{f}(\mathbf{x}_i) \le \theta) \le \exp\left(\theta \sum_{j=1}^k \alpha_j\right) \left(\prod_{j=1}^k z_j\right) \sum_{i=1}^n w_i^{(k+1)}.$$

• According to Eqs. (20) and (21),

$$\bigstar \left\{ \begin{array}{c} \underset{i=1}{\checkmark} \\ \begin{array}{c} \underset{i=1}{\checkmark} \\ \end{array} \right\} \begin{array}{c} w_i^{(j+1)} = \frac{w_i^{(j)} \exp\left(-y_i \, \alpha_j \, h_j(\mathbf{x}_i)\right)}{\binom{Z_j}{2}}, \\ \end{array} \right\} \\ \overbrace{\left(\widehat{z_j}\right)}^n := \sum_{i=1}^n w_i^{(j)} \exp\left(-y_i \, \alpha_j \, h_j(\mathbf{x}_i)\right), \end{array} \right\}$$

we have:

$$\sum_{i=1}^{n} w_{i}^{(j+1)} = \frac{\sum_{i=1}^{n} w_{i}^{(j)} \exp\left(-y_{i} \alpha_{j} h_{j}(\mathbf{x}_{i})\right)}{\sum_{i=1}^{n} w_{i}^{(j)} \exp\left(-y_{i} \alpha_{j} h_{j}(\mathbf{x}_{i})\right)} = \underbrace{1.}$$

Therefore:

$$(:) \mathbb{P}(y_i \, \widehat{f}(\mathbf{x}_i) \le \theta) \le \exp\left(\theta \sum_{j=1}^k \alpha_j\right) \left(\prod_{j=1}^k z_j\right).$$
(23)

• On the other hand, according to Eq. (21), $z_j := \sum_{i=1}^n w_i^{(j)} \exp\left(-y_i \alpha_j h_j(\boldsymbol{x}_i)\right)$, we have:

• Recall Eq. (20) for
$$w_i^{j+1}$$
, i.e., $w_i^{(j)} = \frac{w_i^{(j)} \exp(-\alpha_j) \mathbb{I}(y_i = h_j(\mathbf{x}_i))}{z_i} + \frac{w_i^{(j)} \exp(\alpha_j) \mathbb{I}(y_i \neq h_j(\mathbf{x}_i))}{z_i}$. This is in the range [0, 1]

and its summation over error cases can be considered as the probability of error:

$$\sum_{i=1}^{n} (w_i^{(j)}) \mathbb{I}(y_i \neq h_j(\mathbf{x}_i)) = \mathbb{P}(y_i \neq h_j(\mathbf{x}_i)) \stackrel{(a)}{=} \underline{L}_j,$$
(25)

where (a) is because the Eq. (3), $L_j = \frac{\sum_{i=1}^{n} w_i \mathbb{I}(y_i \neq h_j(\mathbf{x}_i))}{\sum_{i=1}^{n} w_i}$, is the cost which is the probability of error. Therefore, the Eq. (24) becomes:

$$z_j = \exp(-\alpha_j) \left(1 - L_j\right) + \exp(\alpha_j) L_j.$$

• Recall the formula of α_j in the Algorithm, i.e., $\alpha_j = \log(\frac{1-L_j}{L_j})$. Scaling it is not harmful to AdaBoost:

$$\alpha_j = \left(\frac{1}{2}\right) \log(\frac{1-L_j}{L_j}). \tag{26}$$

Therefore, we can have:

$$\begin{array}{c} \star \\ & z_{j} = \exp(-\alpha_{j})\left(1 - L_{j}\right) + \exp(\alpha_{j})L_{j} \\ & = \exp(-\frac{1}{2}\log(\frac{1 - L_{j}}{L_{j}}))\left(1 - L_{j}\right) + \exp(\frac{1}{2}\log(\frac{1 - L_{j}}{L_{j}}))L_{j} \\ & = \exp(\log(\sqrt{\frac{L_{j}}{1 - L_{j}}}))\left(1 - L_{j}\right) + \exp(\log(\sqrt{\frac{1 - L_{j}}{L_{j}}}))L_{j} \\ & = \sqrt{\frac{L_{j}}{1 - L_{j}}} + \sqrt{\frac{1 - L_{j}}{L_{j}}} + \sqrt{\frac{1 - L_{j}}{L_{j}}} + \sqrt{\frac{L_{j}(1 - L_{j})}{L_{j}}} \\ & \implies z_{j} = 2\sqrt{L_{j}(1 - L_{j})}. \end{array}$$

$$(27)$$

• We found Eqs. (26), (27), and (23):

$$\underbrace{\mathbf{x}_{j}}_{\alpha_{j}} = \frac{1}{2} \log(\frac{1-L_{j}}{L_{j}}), \quad z_{j} = 2\sqrt{L_{j}(1-L_{j})}, \quad \mathbb{P}(y_{i} \ \widehat{f}(\mathbf{x}_{i}) \le \theta) \le \exp\left(\theta \sum_{j=1}^{k} \widehat{\alpha_{j}}\right) \left(\prod_{j=1}^{k} \widehat{\xi_{j}}\right).$$

• Plugging Eqs. (26) and (27) in Eq. (23) gives: $\mathbb{P}(y_i \ \widehat{f}(\mathbf{x}_i) \le \theta) \le \exp\left(\frac{1}{2}\theta \sum_{j=1}^k \log(\frac{1-L_j}{L_j})\right) \left(2^k \prod_{j=1}^k \sqrt{L_j(1-L_j)}\right)$

$$= 2^{k} \exp\left(\sum_{j=1}^{k} \log\left(\frac{1-L_{j}}{L_{j}}\right)\right) \prod_{j=1}^{k} \sqrt{L_{j}(1-L_{j})}$$

$$= 2^{k} \prod_{j=1}^{k} \exp\left(\log\left(\frac{1-L_{j}}{L_{j}}\right)^{\theta/2}\right) \prod_{j=1}^{k} \sqrt{L_{j}(1-L_{j})}$$

$$= 2^{k} \prod_{j=1}^{k} \left(\frac{1-L_{j}}{L_{j}}\right)^{\theta/2} \prod_{j=1}^{k} \sqrt{L_{j}(1-L_{j})} = 2^{k} \prod_{j=1}^{k} \sqrt{\frac{1-L_{j}}{L_{j}}} U_{j}(1-L_{j}),$$

which simplifies to the upper bound on the generalization error of AdaBoost [9]:

$$\mathbb{P}(y_i \, \widehat{f}(\mathbf{x}_i) \le \theta) \le 2^k \prod_{j=1}^k \sqrt{L_j^{1-\theta}(1-L_j)^{1+\theta}}$$
(28)

• We found the upper bound on the generalization error of AdaBoost [9]:

$$\bigstar \quad \mathbb{P}\big(y_i \, \widehat{f}(\mathbf{x}_i) \leq \theta\big) \leq 2^k \prod_{j=1}^k \sqrt{L_j^{1-\theta} (1-L_j)^{1+\theta}}, \quad \checkmark$$

where $\mathbb{P}(y_i \hat{f}(\mathbf{x}_i) \leq \theta)$ is the probability that the generalization (true) error for the *i*-th instance is less than $\theta > 0$. • According to Eq. (3), $L_j = \frac{\sum_{i=1}^n w_i \mathbb{I}(y_i \neq h_j(x_i))}{\sum_{i=1}^n w_i}$, we have $L_j \in [0, 1]$. If we have: $L_i \leq 0.5 - \xi$ where $\xi \in (0, 0.5)$, the Eq. (28) becomes: 2(1-Li $\sqrt[\kappa]{L_j^{1- heta}(1-L_j)^{1+ heta}} = \prod_{j=1}^{\kappa} 2 \sqrt{L_j^{1- heta}(1-L_j)^{1+ heta}}$ $\mathbb{P}(y_i \, \widehat{f}(\boldsymbol{x}_i) \leq \theta)$ $(-L_j)^{1+ heta} = \prod_{j=1}^{n-1}$ 1/2¹⁻⁰ $^{-\theta} 2^{1+\theta} (1-L_j)^{1+\theta}$ $(1-\theta)(2(1-L_i))^{1+\theta}$ (30)

- We found: $\left(\mathbb{P}(y_i \, \widehat{f}(\mathbf{x}_i)) \leq \left(\sqrt{(1-2\xi)^{1-\theta}(1+2\xi)^{1+\theta}}\right)\right)$ which is a very good upper bound because $(\theta < \xi, \phi)$ have $\sqrt{(1-2\xi)^{1-\theta}(1+2\xi)^{1+\theta}} < 1$; thus, the probability of error, $\mathbb{P}(y_i \hat{f}(\mathbf{x}_i) \leq \theta)$, decreases exponentially with k which is the number of models used in boosting. This shows that boosting helps us reduce the generalization error and thus helps us avoid overfitting. In other words, because of the bound on generalization error, boosting overfits verv hardly. • If ξ is a very small positive number, the $L_i \leq 0.5 - \xi$ is a little smaller than 0.5, i.e., $L_i \lesssim 0.5$. As we are discussing binary classification in boosting, $L_i = 0.5$ means random classification by flipping a coin. Therefore, for having the great bound of Eq. (30), having weak base models (a little better than random decision) suffices. This shows the effectiveness of boosting.
 - Note that a very small ξ means a very small θ because of $\theta < \xi$; therefore, it means a very small probability of error because of $\mathbb{P}(y_i \hat{f}(\mathbf{x}_i) \le \theta)$.
 - It is noteworthy that both <u>boosting and bagging can</u> be seen as <u>ensemble learning</u> [10] (or majority voting) methods which use <u>model averaging</u> [11, 12] and are very effective in learning theory.
 - Moreover, both <u>boosting and bagging</u> reduce the <u>variance of estimation</u> [13, 9], especially for the models with high variance of estimation such as trees [14].
 - In the above, we analyzed boosting for <u>binary classification</u>. A similar discussion can be done for <u>multi-class</u> classification in boosting and find an upper bound on the generalization error (see the appendix in [9] for more <u>details</u>).

 $(0.99)(0.99)(0.49)\cdots(0.48)$ 0.0001

Boosting as Maximum Margin Classifier

Boosting as Maximum Margin Classifier

- In another perspective, the found upper bound for boosting shows that boosting can be seen as a method to increase (maximize) the margins of training error which results in a good generalization error [15]. This phenomenon is the base for the theory of Support Vector Machines (SVM) [16, 17].
- In the following, we analyze the <u>analogy between maximum margin classifier (i.e., SVM)</u> and boosting [9]. In addition to [9], some more discussions exist for <u>upper bound and</u> margin of boosting [18, 19] to which we refer the interested readers.
- Assume we have training instances {(x_i, y_i)}ⁿ_{i=1} where y_i ∈ {−1, +1} for binary classification. The two classes may not be <u>linearly separable</u>. In order to handle this case, we map the data to higher dimensional feature space using kernels [20, 21], hoping that they become linearly separable in the feature space. Assume *h*(*x*) is a vector which non-linearly maps data to the feature space. Considering α as the vector of optimization variables, the optimization problem [22] in SVM is [17, 9].

Note that
$$y_i = \pm 1$$
 and $\alpha^{\top} h(x_i) \ge 0$; therefore, the sign of $y_i(\alpha^{\top} h(x_i))$ determines the class of the *i*-th instance. (31)

Boosting as Maximum Margin Classifier

• Eq. (31) was: • Con the other hand, the Eq. (16), $\hat{f}(\mathbf{x}_i) = \frac{\sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i)}{\sum_{j=1}^k \alpha_j}$, can be written in a vector form: $\hat{f}(\mathbf{x}_i) = \frac{\sum_{j=1}^k \alpha_j h_j(\mathbf{x}_i)}{\sum_{j=1}^k \alpha_j}$ (32)

where $h(\mathbf{x}_i) = [h_1(\mathbf{x}_i), \dots, h_k(\mathbf{x}_i)]^\top$ and $\alpha = [\alpha_1, \dots, \alpha_k]^\top$. Note that here, $h(\mathbf{x}_i) = \pm 1$ and α_j is obtained from Eq. (26), $\alpha_j = \frac{1}{2} \log(\frac{1-L_j}{L_j})$, or the formula of α_j in the Algorithm. • The similarity between the Eq. (32) and the cost function in Eq. (31) shows that

- The similarity between the Eq. (32) and the cost function in Eq. (31) shows that boosting can be seen as maximizing the margin of classification resulting in a good generalization error [9].
- In other words, finding a linear combination in the high dimensional feature space having a large margin between the training instances of the classes is performed in the two methods.
- Note that a slight difference is the type of norm which is interpretable because the mapping to feature space in boosting is only to $h(x_i) = \pm 1$ while in SVM, it can be any number where the sign is important. Therefore, ℓ_1 and ℓ_2 norms are suitable in boosting and SVM, respectively [9].

Boosting as Maximum Margin Classifier

Another connection between SVM (maximum margin classifier) and boosting is that some of the training instances are found to be most important instances, called support vectors [17]. In boosting, also, weighting the training instances can be seen as selecting some informative models [23] which can be analogous to support vectors.

Acknowledgment

- For more information on boosting in machine learning, see the book: Trevor <u>Hastie</u>, Robert Tibshirani, Jerome H. Friedman, Jerome H. Friedman. "<u>The elements of statistical</u> learning: <u>data mining</u>, inference, and prediction". Vol. 2. New York: springer, 2009 [24].
- Some slides are based on our tutorial paper: "The Theory Behind Overfitting, Cross Validation, Regularization, Bagging, and Boosting: Tutorial" [25]
- Some slides of this slide deck are inspired by teachings of Prof. Ali Ghodsi at University of Waterloo, Department of Statistics and Rrof. Hoda Mohammadzade at Sharif University of Technology, Department of Electrical Engineering.

References

- M. Kearns, "Thoughts on hypothesis boosting," *Technical Report, Machine Learning class project*, pp. 1–9, 1988.
- [2] M. Kearns and L. Valiant, "Cryptographic limitations on learning boolean formulae and finite automata," *Journal of the ACM (JACM)*, vol. 41, no. 1, pp. 67–95, 1994.
- [3] Y. Freund and R. E. Schapire, "Experiments with a new boosting algorithm," in Proceedings of the Thirteenth International Conference on Machine Learning, pp. 148â–156, Morgan Kaufman, San Francisco, 1996.
- [4] T. J. Hastie and R. Tibshirani, "Generalized additive models," *Statistical Science*, vol. 1, no. 3, pp. 297–318, 1986.
- [5] J. Friedman, T. Hastie, and R. Tibshirani, "Additive logistic regression: a statistical view of boosting," *The annals of statistics*, vol. 28, no. 2, pp. 337–407, 2000.
- [6] R. Rojas, "Adaboost and the super bowl of classifiers: a tutorial introduction to adaptive boosting," *Freie University, Berlin, Technical Report*, 2009.
- [7] D. S. Broomhead and D. Lowe, "Multivariable functional interpolation and adaptive networks," *Complex Systems*, vol. 2, pp. 321–355, 1988.
- [8] F. Schwenker, H. A. Kestler, and G. Palm, "Three learning phases for radial-basis-function networks," *Neural networks*, vol. 14, no. 4-5, pp. 439–458, 2001.

References (cont.)

- [9] R. E. Schapire, Y. Freund, P. Bartlett, and W. S. Lee, "Boosting the margin: A new explanation for the effectiveness of voting methods," *The annals of statistics*, vol. 26, no. 5, pp. 1651–1686, 1998.
- [10] R. Polikar, "Ensemble learning," in Ensemble machine learning, pp. 1-34, Springer, 2012.
- [11] J. A. Hoeting, D. Madigan, A. E. Raftery, and C. T. Volinsky, "Bayesian model averaging: a tutorial," *Statistical science*, pp. 382–401, 1999.
- [12] G. Claeskens and N. L. Hjort, Model selection and model averaging. Cambridge Books, Cambridge University Press, 2008.
- [13] L. Breiman, "Arcing classifier (with discussion and a rejoinder by the author)," The annals of statistics, vol. 26, no. 3, pp. 801–849, 1998.
- [14] J. R. Quinlan, "Bagging, boosting, and c4.5," in AAAI/IAAI Conference, vol. 1, pp. 725–730, 1996.
- [15] B. E. Boser, I. M. Guyon, and V. N. Vapnik, "A training algorithm for optimal margin classifiers," in *Proceedings of the fifth annual workshop on Computational learning theory*, pp. 144–152, ACM, 1992.
- [16] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.

References (cont.)

- [17] C. J. Burges, "A tutorial on support vector machines for pattern recognition," Data mining and knowledge discovery, vol. 2, no. 2, pp. 121–167, 1998.
- [18] L. Wang, M. Sugiyama, C. Yang, Z.-H. Zhou, and J. Feng, "On the margin explanation of boosting algorithms," in *COLT*, pp. 479–490, Citeseer, 2008.
- [19] W. Gao and Z.-H. Zhou, "On the doubt about margin explanation of boosting," Artificial Intelligence, vol. 203, pp. 1–18, 2013.
- [20] B. Scholkopf and A. J. Smola, Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT press, 2001.
- [21] T. Hofmann, B. Schölkopf, and A. J. Smola, "Kernel methods in machine learning," The annals of statistics, pp. 1171–1220, 2008.
- [22] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [23] Y. Freund, "Boosting a weak learning algorithm by majority," *Information and computation*, vol. 121, no. 2, pp. 256–285, 1995.
- [24] T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*, vol. 2. Springer, 2009.

References (cont.)

[25] B. Ghojogh and M. Crowley, "The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial," arXiv preprint arXiv:1905.12787, 2019.