## SL : BOOSTING

Boosting is a powerful technique, which can produce very accurate learnears. Boosting combines the prediction of many (weak) learnears to produce a powerful 'committee'/ aka a strong learner.

→ A weak learner $f$ does slightly better than a random guess. The formal definition of a weak learner was introduced by Vaillant (1984) in the context of PAC learning ( $\underline{P}$robably $\underline{A}$pproximately $\underline{C}$orrect ).

→ A strong learner $f$ does arbitrailly well.

For us, it is enough to understand that boosting converts a set of weak learners into a strong one by taking a weighted combination of the weak learner's decisions. Formally, consider

$$\mathcal{F} = \left\{ \sum_{j=1}^{M} \beta_j f_j(\cdot) \quad \middle| \quad \begin{array}{l} \|\beta\|_1 \leq 1 \\ f_j : X \to [-1, 1] = \text{soft classifiers} \end{array} \right\}$$

strong learner
"good prediction accuracy"

weak learner
"prediction accuracy close to a random guess"

### I. RISK BOUND FOR BOOSTING

In this section, we derive an oracle inequality for boosting. In the context of convex relaxation for binary classification, the risk of an element $f \in \mathcal{F}$ is $R_\varphi(f) = \mathbb{E}\,\varphi(-Yf(X))$, where $Y \in \{-1, 1\}$, and $\varphi$ = convex surrogate.

---

The empirical $\varphi$-risk, based on a random sample $\mathcal{L}_n$, is

$$\hat{R}_{n,\varphi} = \frac{1}{n} \sum_{i=1}^{n} \varphi(-Y_i f(X_i)).$$

Put • $\hat{f}_n \in \underset{f \in \mathcal{F}}{\text{argmin}}\ \hat{R}_{n,\varphi}(f)$ = empirical risk minimizer

• $\bar{f} \in \underset{f \in \mathcal{F}}{\text{argmin}}\ R(f)$ = best element in the class

Step I. Bound for $R_\varphi(\hat{f}_n) - R_\varphi(\bar{f})$

↖ the risk of $\hat{f}_n$ is computed conditionally on $\mathcal{F}$.

$$R_\varphi(\hat{f}_n) = R_\varphi(\hat{f}_n) + \underbrace{\hat{R}_{n,\varphi}(\hat{f}_n) - \hat{R}_{n,\varphi}(\hat{f}_n)}_{\leq \hat{R}_{n,\varphi}(\bar{f})} + R_\varphi(\bar{f}) - R_\varphi(\bar{f})$$

$$\leq R_\varphi(\hat{f}_n) + \hat{R}_{n,\varphi}(\bar{f}) - \hat{R}_{n,\varphi}(\hat{f}_n) + R_\varphi(\bar{f}) - R_\varphi(\bar{f})$$

$$\leq R_\varphi(\bar{f}) + 2 \underbrace{\underset{f \in \mathcal{F}}{\sup} \left| \hat{R}_{n,\varphi}(f) - R_\varphi(f) \right|}_{\overset{\|}{\underset{f \in \mathcal{F}}{\sup} \left| \frac{1}{n} \sum_{i=1}^{n} \left\{ \varphi(-Y_i f(X_i)) - \mathbb{E}\varphi(-Yf(X)) \right\} \right|}}$$

We assume now that $\varphi$ is an $L$-Lipschitz convex surrogate:

$$\varphi : [-1, 1] \to \mathbb{R}_+ \quad ; \quad \forall u, v \in [-1, 1], \quad |\varphi(u) - \varphi(v)| \leq L\,|u - v|.$$

EX: • Hinge $\varphi(x) = \max(x+1, 0)$ ; $L = 1$

• Exp $\varphi(x) = e^x$ ; $L = e$

• Logistic $\varphi(x) = \log_2(1 + e^x)$ ; $L = \frac{e}{1+e} \log_2 e \simeq 2.43$

Recall that if $\forall x \in [-1, 1]$ $\varphi'(x) \leq L$, then $\forall u, v \in [-1, 1]$

$$\varphi(u) - \varphi(v) = \int_v^u \varphi'(x)\,dx \leq \int_v^u L\,dx = L(u - v)$$
$$(u > v)$$

Consider the function

$$(x_1, y_1) \cdots (x_n, y_n) \xmapsto{g} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varphi(-y_i f(x_i)) - R_\varphi(f) \right|$$

Then

$$g\Big((x_1, y_1), \ldots, (x_i, y_i), \ldots, (x_n, y_n)\Big)$$

Change only the i-th entry

$$- g\Big((x_1, y_1), \ldots, (x_i', y_i'), \ldots, (x_n, y_n)\Big)$$

$$\leq \frac{1}{n}\Big(\varphi(1) - \varphi(-1)\Big) \leq \frac{2L}{n}.$$

$\varphi$ is non-decreasing     $\varphi$ is L-Lipschitz

$\Rightarrow$ We can apply the bounded difference inequality (p.16 in SL: VAPNIK CHERVONENKIS THEORY)

$$\mathbb{P}\left( \left| \sup_{f \in \mathcal{F}} |\hat{R}_{n,\varphi}(f) - R(f)| - \mathbb{E}(\ldots) \right| > \varepsilon \right) \leq 2 \exp\left( -\frac{2\varepsilon^2}{\Sigma\left(\frac{2L}{n}\right)^2} \right)$$

$\Longleftrightarrow$

$$\sup_{f \in \mathcal{F}} |\hat{R}_{n,\varphi}(f) - R(f)| \leq \mathbb{E}(\ldots) + 2L\sqrt{\frac{\log(2/\delta)}{2n}}$$

w.p. $\geq 1-\delta$.

$\Rightarrow$ We can now focus on bounding $\mathbb{E}\left( \sup_{f \in \mathcal{F}} |\hat{R}_{n,\varphi}(f) - R(f)| \right)$.

Using the symmetrization trick, $\mathbb{E}(\ldots)$ is upper-bounded by twice the Rademacher complexity

$$\mathcal{R}_n(\varphi \circ \mathcal{F}) := \sup_{\substack{(x_1, y_1) \\ \cdots (x_n, y_n)}} \mathbb{E}\left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \sigma_i \varphi(-y_i f(x_i)) \right| \right\}$$

compare with the expression on page 21 in SL: VC THEORY.

$$\sigma_i = \pm 1 \quad \text{w.p. } \frac{1}{2}.$$

---

The convex surrogate $\varphi$ is such that $\varphi(0)=1$. To apply a contraction inequality on $\mathcal{R}_n(\varphi \circ \mathcal{F})$, we need $\varphi(0)=0$ $\Rightarrow$ put $\Psi(x) = \varphi(x) - 1$. Then

$$\mathbb{E}(\ldots) \leq 2 \mathcal{R}_n(\Psi \circ \mathcal{F})$$

and

$$\mathcal{R}_n(\Psi \circ \mathcal{F}) \leq 2L \, \mathcal{R}_n(\mathcal{F}),$$

contraction inequality (not trivial) [details omitted]

where $\mathcal{R}_n(\mathcal{F}) = \sup_{\substack{(x_1, y_1) \cdots \\ (x_n, y_n)}} \mathbb{E}\left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \sigma_i y_i f(x_i) \right| \right\}.$

Summary: So far, we have established that

$$R_\varphi(\hat{f}_n) - R_\varphi(\bar{f}) \leq 8L \, \mathcal{R}_n(\mathcal{F}) + 4L\sqrt{\frac{\log(2/\delta)}{2n}} \quad \text{w.p.} \geq 1-\delta$$

Step II. Bound for $\mathcal{R}_n(\mathcal{F})$.

Recall that $\mathcal{F} = \left\{ \sum_{j=1}^{M} \beta_j f_j(\cdot) \;\middle|\; \|\beta\|_1 \leq 1 \; ; \; f_j : X \to [-1, 1] \right\}$

fixed in advance (and so is M)

Then

$$\mathcal{R}_n(\mathcal{F}) = \sup \mathbb{E}\left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^{M} \beta_j \sum_{i=1}^{n} \sigma_i y_i f_j(x_i) \right| \right\}$$

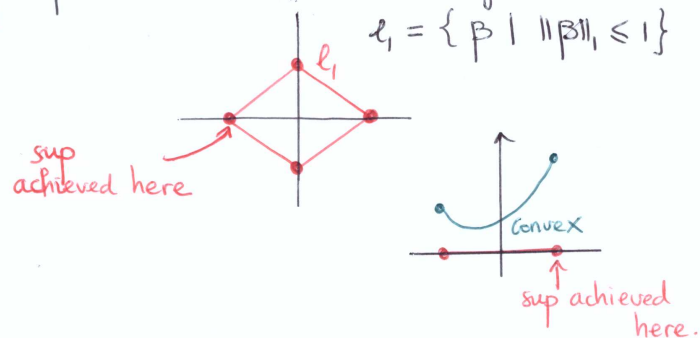Put $|Z_\beta| := \left| \sum_{j=1}^{M} \beta_j \sum_{i=1}^{n} \sigma_i y_i f_j(x_i) \right|.$

$\pm 1 \quad \pm 1$

Note that $\beta \mapsto |Z_\beta|$ is a convex function of $\beta$

$$\left| \sum_{j=1}^{M} (\lambda \beta_{1,j} + (1-\lambda)\beta_{2,j}) \sum_{i=1}^{n} \tau_i y_i f_j(x_i) \right|$$

$$\leq \lambda \left| \sum_j \beta_{1,j} \sum_i \tau_i y_i f_j(x_i) \right|$$

$$+ (1-\lambda)\left| \sum_j \beta_{2,j} \sum_i \tau_i y_i f_j(x_i) \right|$$

$$\Rightarrow \sup_{\|\beta\|_1 \leq 1} |Z_\beta| \text{ is achieved at a vertex of the unit ball}$$

$$\ell_1 = \{ \beta \mid \|\beta\|_1 \leq 1 \}$$



sup achieved here

convex

sup achieved here.

$\Rightarrow$ Consider the set

$$B_M := \left\{ \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} \cdots \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} -1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \\ 0 \\ \vdots \end{pmatrix} \cdots \begin{pmatrix} 0 \\ \vdots \\ 0 \\ -1 \end{pmatrix} \right\}$$

$|B_M| = 2M$,

so that $\quad \mathbb{E} \sup_{f \in \mathcal{F}} |Z_\beta| = \mathbb{E} \max_{\beta \in B_M} |Z_\beta|$

Proceed as usual:

$$\mathbb{E} \max_{\beta \in B_M} |Z_\beta| = \frac{1}{s} \log \exp s \, \mathbb{E} \max |Z_\beta|$$

$$\leq \frac{1}{s} \log \mathbb{E} \, e^{s \max |Z_\beta|}$$

$|\tau_i y_i \sum_j \beta_j f_j(x_i)|$
$< \sum_j |\beta_j| \|f_j\| = 1$

$$\leq \frac{1}{s} \log \sum_{\beta \in B_M} e^{s Z_\beta}$$ Hoeffding lemma

$$\leq \frac{1}{s} \log \left\{ |B_M| \left( e^{\frac{s^2 2^2}{8}} \right)^n \right\}$$

$$\mathbb{E} \max_{\beta \in B_M} |Z_\beta| \leq \frac{\log |B_M|}{s} + \frac{sn}{2}$$

Take $s = \sqrt{(2 \log |B_M|)/n}$

$$\leq 2 \sqrt{\frac{n}{2} \log |B_M|}, \quad \text{with } |B_M| = 2M.$$

$$\Rightarrow \mathcal{R}_n(\mathcal{F}) \leq \sqrt{\frac{2 \log 2M}{n}}$$

Summary:

$$R_\varphi(\hat{f}_n) - R_\varphi(\bar{f}) \leq 8L \sqrt{\frac{2 \log 2M}{n}} + 4L \sqrt{\frac{\log(2/\delta)}{2n}} \quad \text{w.p} \geq 1-\delta$$

Step III. Bound for the excess risk $R(\text{sign} \, \hat{f}_n) - R^*$

↳ risk under 0/1 loss

$$R(\text{sign} \, \hat{f}_n) - R^* \leq 2c \left( R_\varphi(\hat{f}_n) - R_\varphi^* \right)^\gamma$$

Zhang lemma (p.7 in SL: CONVEX RELAXATION)

$$= 2c \left( R_\varphi(\hat{f}_n) - R_\varphi(\bar{f}) + R_\varphi(\bar{f}) - R_\varphi^* \right)^\gamma$$

$$\leq 2c \left( 8L \sqrt{\frac{2 \log 2M}{n}} + 4L \sqrt{\frac{\log(2/\delta)}{2n}} \right.$$

$$\left. + R_\varphi(\bar{f}) - R_\varphi^* \right)^\gamma$$

since for $\alpha_i \geq 0$
and $\gamma \in [0,1]$,

$(\alpha_1 + \alpha_2 + \alpha_3)^\gamma \leq \alpha_1^\gamma + \alpha_2^\gamma + \alpha_3^\gamma$

$$\leq 2c \left( R_\varphi(\bar{f}) - R_\varphi^* \right)^\gamma + 2c \left( 8L \sqrt{\frac{2 \log 2M}{n}} \right)^\gamma$$

$$+ 2c \left( 4L \sqrt{\frac{\log(2/\delta)}{2n}} \right)^\gamma$$

## Theorem: Risk Bound for Boosting

Consider. $\mathcal{F} = \left\{ \sum_{j=1}^{M} \beta_j f_j(\cdot) \;\middle|\; \|\beta\|_1 \leq 1, \; f_j : X \to [-1,1] \right\}$

- $\varphi$ = an $L$ - lipschitz convex surrogate.
- $\hat{f}_n = \underset{f \in \mathcal{F}}{\arg\min} \; \hat{R}_{n,\varphi}(f)$.

Then

$$R(\text{sign}\,\hat{f}_n) - R^* \leq 2c\left(R_\varphi(\bar{f}) - R_\varphi^*\right)^\gamma + 2c\left(8L\sqrt{\frac{2\log 2M}{n}}\right)^\gamma$$
$$+ 2c\left(4L\sqrt{\frac{\log(2/\delta)}{2n}}\right)^\gamma$$

$\underbrace{\phantom{R(\text{sign}\,\hat{f}_n) - R^*}}$
excess risk of
the hard
classifier

$\qquad\qquad\qquad\qquad\qquad\qquad w.p \geq 1-\delta$

---

## II - Ada Boost

### II.1. Derivation of the algorithm

AdaBoost (introduced by Freund & Shapire (1995) ) uses an
exponential loss:

- $\hat{f}_n \in \underset{f \in \mathcal{F}}{\arg\min} \; \hat{R}_{n,\varphi}(f)$

- $\hat{R}_{n,\varphi}(f) = \frac{1}{n}\sum_{i=1}^{n} \varphi(-Y_i f(X_i)) = \frac{1}{n}\sum_{i=1}^{n} e^{-Y_i f(X_i)}$

- $\mathcal{F} = \left\{ \sum_{j=1}^{M} \beta_j f_j(\cdot) \;\middle|\; \|\beta\|_1 \leq 1 \quad f_j : X \to [-1,1] \right\}$

AdaBoost is an algorithm that approximately computes $\hat{f}_n$.
Several assumptions are relaxed =

→ M is not fixed in advance // coefficients are not constrained in $\ell_1$.

→ $f_j(\cdot)$ are not fixed in advance : the dictionary adapts
to the data . We only require that $f_j$ take a
particular form

---

→ AdaBoost does not minimize $\frac{1}{n}\sum_{i=1}^{n} e^{-Y_i f(X_i)}$ directly,
but performs this recursively.

Suppose that at iteration $(m-1)$, the current estimate is

$$f^{(m-1)}(x) = \sum_{k=1}^{m-1} \beta_k f_k(x) \quad \text{taking values in } \{-1,1\}$$

Iteration $m$ is searching for the 'best' value
of $\beta_m$ and $f_m$ such that:

$$\boxed{(\beta_m, f_m) \in \underset{\beta, f}{\arg\min} \; \sum_{i=1}^{n} \varphi\left(-y_i\left[f^{(m-1)}(x_i) + \beta f(x_i)\right]\right)}$$

OPTIMIZATION PROBLEM SOLVED BY AdaBoost.

Put $w_i^{(m)} := \exp\left(-y_i f^{(m-1)}(x_i)\right)$. Then

$$(\beta_m, f_m) \in \underset{\beta, f}{\arg\min} \left\{ \sum_{i=1}^{n} w_i^{(m)} e^{-\beta f(x_i) y_i} \right\}$$
$\qquad\qquad\qquad\qquad \hookleftarrow \text{call this } W^{(m)}$

$$W^{(m)} = e^{-\beta} \sum_{\substack{i \\ y_i = f(x_i)}} w_i^{(m)} + e^{\beta} \sum_{\substack{i \\ y_i \neq f(x_i)}} w_i^{(m)}$$

$$= e^{-\beta} \sum_{i=1}^{n} w_i^{(m)} \mathbf{1}(y_i = f(x_i)) + e^{\beta} \sum_{i=1}^{n} w_i^{(m)} \mathbf{1}(y_i \neq f(x_i))$$

$$= e^{-\beta} \sum_{i=1}^{n} w_i^{(m)} \left(1 - \mathbf{1}(y_i \neq f(x_i))\right)$$
$$\qquad\qquad + e^{\beta} \sum_{i=1}^{n} w_i^{(m)} \mathbf{1}(y_i \neq f(x_i))$$

$$= (e^{\beta} - e^{-\beta}) \sum_{i=1}^{n} w_i^{(m)} \mathbf{1}(y_i \neq f(x_i))$$
$$\qquad\qquad + \underbrace{e^{-\beta} \sum_{i=1}^{n} w_i^{(m)}}_{\text{indpt of } f}$$

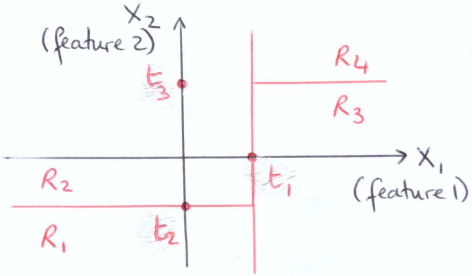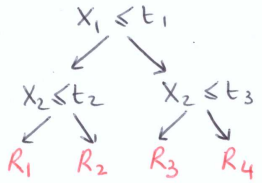⇒ At step $m$, the weak classifier $f_m$ is selected according to:

$$f_m \in \underset{f}{\operatorname{argmin}} \sum_{i=1}^{n} w_i^{(m)} \mathbf{1}(y_i \neq f(x_i))$$

↑ Weighted version of the empirical 0/1 loss

Example: Fit a decision tree to the weighted version of the learning sample

__CART procedure__.



$X_1 \leq t_1$

$X_2 \leq t_2$   $X_2 \leq t_3$

$R_1$   $R_2$   $R_3$   $R_4$

4 terminal nodes = 4 regions.

Q: How to choose the split variable + split point ?

Q̄: How to classify point in each region ?

↳ Answer to this question is relatively easy:
For each class $\in \{-1, 1\}$, compute the weighted proportion of points falling into region $R_\ell$:

$$\hat{P}_{\ell k}^{(m)} = \frac{\sum_{x_i \in R_\ell} w_i^{(m)} \mathbf{1}(y_i = k)}{\sum_{x_i \in R_\ell} w_i^{(m)}}$$

Then classify according to $\hat{k} = \underset{k \in \{-1, 1\}}{\operatorname{argmax}} \hat{P}_{\ell k}^{(m)}$

↳ growing a tree; ie finding the terminal regions, is more tricky = use a __GREEDY ALGORITHM__ + the prune it (using e.g. __cost-complexity pruning__).
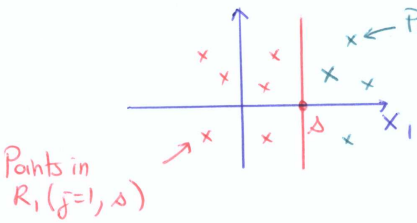
---

More precisely, for a splitting variable $j$ and a splitting point $s$, consider

$$R_1(j, s) = \{x_1, \dots, x_n \mid x_{kj} \leq s\}$$
$$x_k = (x_{k1}, \dots, x_{kd}) \in \mathbb{R}^d$$
$$R_2(j, s) = \{x_1, \dots, x_n \mid x_{kj} > s\}$$



← Points in $R_2(j=1, s)$

Points in $R_1(j=1, s)$

Perform an exhaustive search: for each $j$ and $s$ (there are finitely many — why ?), compute

$$\sum_{x_i \in R_1(j, s)} w_i^{(m)} \mathbf{1}(y_i \neq \hat{y}_1^{(m)}) + \sum_{x_i \in R_2(j, s)} w_i^{(m)} \mathbf{1}(y_i \neq \hat{y}_2^{(m)})$$

where $\hat{y}_k^{(m)} = \underset{y \in \{-1, 1\}}{\operatorname{argmax}} \left( \frac{\sum_{x_i \in R_k(j, s)} w_i^{(m)} \mathbf{1}(y_i = y)}{\sum_{x_i \in R_k(j, s)} w_i^{(m)}} \right)$

And return the value of $(j, s)$ that minimizes this.

⇒ Grow a large tree $T$; and then prune it to avoid overfitting.

• __Complexity criterion used is__

$$C_\lambda(T) = \sum_{\ell=1}^{|T|} |R_\ell| Q_\ell(T) + \lambda |T|$$

good ness of fit Measure ↗   Tuning Parameter ↑   # of terminal leaves

↳ $Q_\ell(T) = \frac{1}{|R_\ell|} \sum_{x_i \in R_\ell} w_i^{(m)} \mathbf{1}(y_i \neq y_\ell^{(m)})$  [weighted misclass]

where $y_\ell^{(m)} = \underset{j \in \{-1, 1\}}{\operatorname{argmax}} \hat{P}_{\ell j}^{(m)}$

↳ $Q_\ell(T) = \sum_{j \in \{-1, 1\}} \hat{P}_{\ell j}^{(m)} (1 - \hat{P}_{\ell j}^{(m)})$  [Gini index]
↳ CART

↳ $Q_\ell(T) = -\sum_{j \in \{-1, 1\}} \hat{P}_{\ell j}^{(m)} \log \hat{P}_{\ell j}^{(m)}$

For each $\lambda$, find a subtree $T_\lambda \subset T$ minimizing ⑪
$C_\lambda(T)$.

How? Using **weakest link pruning** = successively collapse internal nodes that result in the smallest increase of $\sum_\ell |R_\ell| Q_\ell(T)$, until you reach the root. This gives a finite sequence of trees, and it can be shown that this sequence must contain $T_\lambda$. The value of $\lambda$ is usually selected using cross-validation.

**Remarks:** (i) In boosting, growing large trees + pruning them is not needed. Typically, a rule of thumb is that up to 5/6 terminal nodes is OK. You can also consider a tree with a single split / two terminal nodes (aka STUMPS).

(ii) In the approach above, the weak (tree) classifiers $f_m$ are restricted to taking values in $\{1, ..., K\}$, and in the context of binary classification, in $\{-1, 1\}$.
AdaBoost then combines scaled classification trees $\beta_m f_m$.

<span style="color:blue">"Discrete AdaBoost"</span>

Alternatively, you may boost classification trees without this restriction :

$$\{R_{jm}, \gamma_{jm}\} = \operatorname{argmin} \sum_{i=1}^n \ell\left(y_i, f^{(m-1)}(x_i) + T(x_i)\right)$$

terminal regions

prediction in each terminal region

min taken over the parameters of the tree $T$; ie the split variables, split points, and predicted value.

For example, with a **square loss**, the problem reduces to : ⑪ₐ

$$\operatorname{argmin} \sum_{i=1}^n \left( \underbrace{y_i - f^{(m-1)}(x_i)}_{\text{current residuals}} - T(x_i) \right)^2$$

$\equiv$ fit a regression tree that best fits the current residuals.

<span style="color:green">make the problem (p.9)
$\sum_{i=1}^n w_i^{(m)} \mathbb{1}(...)$
convex</span>

With an **exponential loss**, we need to solve

$$\operatorname{argmin} \sum_{i=1}^n w_i^{(m)} \exp\left(- y_i T(x_i)\right).$$

Once the terminal regions are computed, the predicted value in $R_{jm}$ is

$$\gamma_{jm} = \frac{1}{2} \log\left( \frac{\sum_{x_i \in R_{jm}} w_i^{(m)} \mathbb{1}(y_i = 1)}{\sum_{x_i \in R_{jm}} w_i^{(m)} \mathbb{1}(y_i = -1)} \right)$$

<span style="color:blue">"Real AdaBoost"</span>

<span style="color:teal">Not restricted to $\{-1, 1\}$</span>

<u>Back on track</u>: once $f_m$ is constructed, the next step is to compute $\beta_m$ :

$$\beta_m = \operatorname*{argmin}_\beta \sum_{i=1}^n w_i^{(m)} e^{-\beta y_i f_m(x_i)}$$

$$W^{(m)} = e^{-\beta} \left\{ \sum_{y_i = f_m(x_i)} w_i^{(m)} + e^{2\beta} \sum_{y_i \neq f_m(x_i)} w_i^{(m)} \right\}$$

$$= e^{-\beta} \left\{ \sum_{i=1}^n w_i^{(m)} + (e^{2\beta} - 1) \sum_{i=1}^n w_i^{(m)} \mathbb{1}(y_i \neq f_m(x_i)) \right\}$$

$$= e^{-\beta} \left( \sum_{i=1}^n w_i^{(m)} \right) \left\{ 1 + (e^{2\beta} - 1) \underbrace{\frac{\sum_{i=1}^n w_i^{(m)} \mathbb{1}(y_i \neq f_m(x_i))}{\sum_{i=1}^n w_i^{(m)}}}_{\substack{\| \\ \text{err}^{(m)} \in (0,1)}} \right\}$$
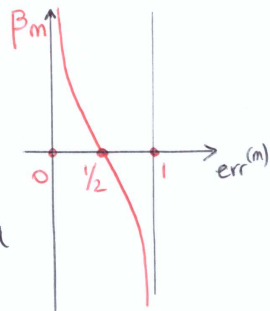
$$= e^{-\beta} \left( \sum_{i=1}^n w_i^{(m)} \right) \left\{ 1 + (e^{2\beta} - 1) \, \text{err}^{(m)} \right\}.$$

Thus
$$\beta_m = \underset{\beta}{\arg\min} \; e^{-\beta} \left\{ 1 + (e^{2\beta} - 1) \, \text{err}^{(m)} \right\}$$

$$\boxed{\beta_m = \frac{1}{2} \log \left( \frac{1 - \text{err}^{(m)}}{\text{err}^{(m)}} \right)}$$



<u>Meaning</u> = if classifier $f_m$ (weak) performs well; so that $\text{err}^{(m)}$ is close to 0, the weight associated to $f_m$ is large positively $\Rightarrow f_m$ has a strong contribution to the final classification $\sum_{k=1}^{M} \beta_k f_k$

→ Once $\beta_m$ and $f_m$ are computed, update

• $f^{(m)}(x) = f^{(m-1)}(x) + \beta_m f_m(x)$

• $w_i^{(m+1)} = \exp\left(- y_i f^{(m)}(x_i)\right)$

$\qquad = \exp\left(- y_i f^{(m-1)}(x_i)\right) \exp\left(- y_i \beta_m f_m(x_i)\right)$

$\qquad = w_i^{(m)} \exp\left(- y_i \beta_m f_m(x_i)\right)$

<u>Remark:</u>

↗ If $y_i = f_m(x_i)$, then $- y_i f_m(x_i) = -1$
$\qquad\qquad\qquad\qquad\qquad = 2 \, \mathbb{1}(y_i \neq f_m(x_i)) - 1$

↘ If $y_i \neq f_m(x_i)$, then $- y_i f_m(x_i) = 1$
$\qquad\qquad\qquad\qquad\qquad = 2 \, \mathbb{1}(y_i \neq f_m(x_i)) - 1$

Thus
$$w_i^{(m+1)} = w_i^{(m)} \exp\left\{ \beta_m \left( 2 \, \mathbb{1}(y_i \neq f_m(x_i)) - 1 \right) \right\}$$

$$= w_i^{(m)} \exp\left\{ \underset{!!}{\underset{2\beta_m}{\alpha_m}} \, \mathbb{1}(y_i \neq f_m(x_i)) \right\} e^{-\beta_m}$$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\uparrow$ independent of $i$
$\qquad\qquad\qquad\qquad\qquad\Rightarrow$ can be dropped (cancels out in the def of $\beta_k^{(m)}$)

---

<u>AdaBoost Algorithm.</u>

(i) Initialize observation weights to $w_i^{(1)} = 1/n$.

(ii) For $m = 1$ to $M$

→ Fit a classifier $f_m$ to the learning sample using weights $w_i^{(m)}$

→ Compute $\text{err}^{(m)} = \dfrac{\sum_{i=1}^{n} w_i^{(m)} \mathbb{1}(y_i \neq f_m(x_i))}{\sum_{i=1}^{n} w_i^{(m)}}$

→ Compute $\alpha_m = \log\left( \dfrac{1 - \text{err}^{(m)}}{\text{err}^{(m)}} \right)$

provided $\alpha_m > 0$, a misclassified observation has its weight increased at the next iteration

→ Update $w_i^{(m+1)} \leftarrow w_i^{(m)} \left(\exp\left( \alpha_m \, \mathbb{1}(y_i \neq f_m(x_i)) \right)\right)$

(iii) Output $f^{(M)}(x) = \sum_{m=1}^{M} \beta_m f_m(x)$ and classify according to sign $f^{(M)}$

<u>Remark:</u> The weights $w_i^{(m+1)}$ are **sometimes** renormalized so that they sum up to 1 ($\equiv$ probability distribution on observations).

Update is $w_i^{(m+1)} = \dfrac{w_i^{(m)}}{Z_m} e^{-\beta_m y_i f_m(x_i)}$

so that $\left[ \sum_{i=1}^{n} w_i^{(m+1)} = 1 \right]$

## II.2. Performance bound

• We show that provided $\text{err}^{(\ell)} = \frac{1}{2} - \gamma_\ell$; for $\gamma_\ell \geq \gamma > 0 \quad \forall \ell$, then
$$\frac{1}{n} \sum_{i=1}^{n} \mathbb{1}(y_i \neq f^{(M)}(x_i)) \leq e^{-2\gamma^2 M}$$

$\Rightarrow$ Training error can be made arbitrarily small, provided $M$ is sufficiently large.

→ **Step I** Cascading the weights,

$$w_i^{(M+1)} = \frac{w_i^{(M)}}{z_M} \exp\left(-y_i \beta_M f_M(x_i)\right)$$

$$w_i^{(M)} = \frac{w_i^{(M-1)}}{z_{M-1}} \exp\left(-y_i \beta_{M-1} f_{M-1}(x_i)\right)$$

$$\vdots$$

$$w_i^{(2)} = \frac{w_i^{(1)}}{z_1} \exp\left(-y_i \beta_1 f_1(x_i)\right)$$

$$w_i^{(1)} = \frac{1}{n}$$

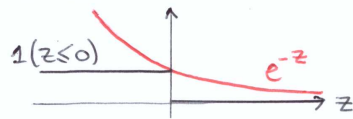$$w_i^{(M+1)} = \frac{1}{n} \frac{1}{\prod_{l=1}^{M} z_l} \exp\left(-y_i \sum_{l=1}^{M} \beta_l f_l(x_i)\right)$$

$$w_i^{(M+1)} = \frac{1}{n \prod_{l=1}^{M} z_l} \exp\left(-y_i f^{(M)}(x_i)\right)$$

→ **Step II**

$$\frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\left(y_i \neq f^{(M)}(x_i)\right) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\left(y_i f^{(M)}(x_i) \leq 0\right)$$

*training error*

$$\leq \frac{1}{n} \sum_{i=1}^{n} \exp\left(-y_i f^{(M)}(x_i)\right)$$

$$\mathbb{1}(z \leq 0) \qquad e^{-z}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left(n\left[\prod_{l=1}^{M} z_l\right] w_i^{(M+1)}\right)$$

$$= \prod_{l=1}^{M} z_l \left(\sum_{i=1}^{n} w_i^{(M+1)}\right) \longmapsto = 1$$

$$= \prod_{l=1}^{M} z_l$$

---

Next, $z_l = \sum_{i=1}^{n} w_i^{(l)} \exp\left(-\beta_l y_i f_l(x_i)\right)$

$$= \sum_{y_i = f_l(x_i)} w_i^{(l)} e^{-\beta_l} + \sum_{y_i \neq f_l(x_i)} w_i^{(l)} e^{\beta_l}$$

$$= e^{-\beta_l} + err^{(l)}\left(e^{\beta_l} - e^{-\beta_l}\right)$$

$$err^{(l)} = \sum_{i=1}^{n} w_i^{(l)} \mathbb{1}\left(y_i \neq f_l(x_i)\right)$$

[compare with page 11 / weights sum to 1]

$$= 2\sqrt{err^{(l)}(1 - err^{(l)})}$$

since $\beta_l = \frac{1}{2} \log\left(\frac{1 - err^{(l)}}{err^{(l)}}\right)$

Thus $\boxed{\dfrac{1}{n} \sum_{i=1}^{n} \mathbb{1}\left(y_i \neq f^{(M)}(x_i)\right) \leq 2^M \prod_{l=1}^{M} \sqrt{err^{(l)}(1 - err^{(l)})}}$

→ **Step III. Conclusion.**

Under the assumption that $err^{(l)} = \frac{1}{2} - \gamma_l$,

$$2\sqrt{err^{(l)}(1 - err^{(l)})} = 2\sqrt{\frac{1}{4} - \gamma_l^2}$$

$$= \sqrt{1 - 4\gamma_l^2}$$

$$\leq \sqrt{e^{-4\gamma_l^2}}$$

$$= e^{-2\gamma_l^2}$$

$e^x \qquad 1+x$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\left(y_i \neq f^{(M)}(x_i)\right) \leq \prod_{l=1}^{M} \exp\left(-2\gamma_l^2\right)$$

$$= \exp\left(-2 \sum_{l=1}^{M} \gamma_l^2\right)$$

$$\leq \exp\left(-2 M \gamma^2\right)$$

provided $\gamma_l \geq \gamma > 0$

⇒ Provided weak classifiers perform slightly better than a random guess, AdaBoost returns a proportion of the training samples misclassified arbitrarily small. Note that if $err^{(\ell)} > \frac{1}{2}$, the corresponding coefficient $\beta_\ell$ is negative, and AdaBoost does not update the weights in the right direction.

⇒ The bound suggests that we can achieve any degree of accuracy, provided $M$ is large enough. However, AdaBoost is prone to overfitting, and $M$ should be limited in size. In practice, one approach is to hold out a part of the training sample, and to test the performance of $f^{(M)}$ as a function of $M$. In fact, we can show that $VC(\mathcal{F}_M) = O\left(VC(\mathcal{F})\, M \log M\right)$, for $\mathcal{F}_M = \left\{ sign\left( \sum_{m=1}^{M} \beta_m f_m \right),\ \beta_m \in \mathbb{R},\ f_m \in \mathcal{F} \right\}$.

## II. 3. AdaBoost for K-class classification

* Coding - scheme adopted here for K-class classification is $Y = (Y_1, \cdots, Y_K)^t$, where $Y_k = \begin{cases} 1 & \text{if } X \text{ belongs to class } k \\ -\frac{1}{K-1} & \text{otherwise} \end{cases}$ $(K \geqslant 3)$

* Goal: Construct $f(x) = (f_1(x), \cdots, f_K(x))^t$, under the constraint $\sum f_k(x) = 0$, and classify $x$ according to the largest value of $f_k(x)$, $k = 1, \cdots, K$.

* Is this a good idea? It is, under a generalized version of the exponential loss

$$\boxed{\ell(Y, f(X)) = \exp\left(-\frac{1}{K} Y^t f(X)\right)}$$

Exp loss for K-class classification.

⇒ Population minimizer $f^*(x) = \underset{f}{\text{argmin}}\ \mathbb{E}\left( \ell(Y, f(X)) \mid X = x \right)$

$$f^*(x) = (f_1^*(x), \cdots, f_K^*(x))$$

is such that

$$f_k^*(x) = (K-1)\left\{ \log \mathbb{P}(Y_k = 1 \mid X = x) - \frac{1}{K} \sum_{j=1}^{K} \log \mathbb{P}(Y_j = 1 \mid X = x) \right\}$$

Proba dos ∈ class k

And thus, classifying $x$ according to the largest value of $f_k(x)$ makes sense.

proof = $f^*(x) = \underset{f}{\text{argmin}}\ \mathbb{E}\left\{ \exp\left(-\frac{1}{K} Y^t f(X)\right) \mid X = x \right\}$

Constrained optimization problem

The Lagrangian is

$$\mathcal{L}(f, \lambda) = \underbrace{\mathbb{E}\left\{ \exp\left(-\frac{1}{K} Y^t f(X)\right) \mid X = x \right\}}_{} - \lambda \sum_{k=1}^{K} f_k(x)$$

$$\| $$

$$\mathbb{E}\left\{ \exp\left(-\frac{1}{K}\left[ Y_1 f_1 + \cdots + Y_K f_K \right]\right) \mid X = x \right\}$$

$$\| $$

$$\exp\left\{ -\frac{1}{K}\left[ f_1 - \frac{1}{K-1} f_2 - \cdots - \frac{1}{K-1} f_K \right] \right\} \mathbb{P}(Y_1 = 1 \mid X = x)$$

$$+ \exp\left\{ -\frac{1}{K}\left[ -\frac{1}{K-1} f_1 + f_2 - \cdots - \frac{1}{K-1} f_K \right] \right\} \mathbb{P}(Y_2 = 1 \mid X = x)$$

$$\cdots$$

$$+ \exp\left\{ -\frac{1}{K}\left[ -\frac{1}{K-1} f_1 - \frac{1}{K-1} f_2 - \cdots + f_K \right] \right\} \mathbb{P}(Y_K = 1 \mid X = x)$$

Making use of $-\frac{1}{K-1} = 1 - \frac{K}{K-1}$,

$$
\overset{\shortparallel}{} \exp\left\{ -\frac{1}{K}\Big( \underbrace{f_1 + \cdots + f_K}_{=0} - \frac{K}{K-1}\underbrace{(f_2 + \cdots + f_K)}_{=-f_1} \Big) \right\} \mathbb{P}(Y_1 = 1 \mid X = x)
$$

$$
+ \exp\left\{ -\frac{1}{K}\Big[ \underbrace{f_1 + \cdots + f_K}_{=0} - \frac{K}{K-1}\underbrace{(f_1 + f_3 + \cdots + f_K)}_{=-f_2} \Big] \right\} \mathbb{P}(Y_2 = 1 \mid X = x)
$$

$$
\vdots
$$

$$
+ \exp\left\{ -\frac{1}{K}\Big[ \underbrace{f_1 + \cdots + f_K}_{=0} - \frac{K}{K-1}\underbrace{(f_1 + \cdots + f_{K-1})}_{=-f_K} \Big] \right\} \mathbb{P}(Y_K = 1 \mid X = x)
$$

$$
\overset{\shortparallel}{}
$$

$$
\exp\left\{ \frac{-f_1}{K-1} \right\} \mathbb{P}(Y_1 = 1 \mid X = x) + \cdots + \exp\left\{ -\frac{f_K}{K-1} \right\} \mathbb{P}(Y_K = 1 \mid X = x)
$$

$\Rightarrow$ Lagrangian becomes

$$
\mathscr{L}(f, \lambda) = \sum_{k=1}^{K} \exp\left\{ -\frac{f_k}{K-1} \right\} \mathbb{P}(Y_k = 1 \mid X = x) - \lambda \sum_{k=1}^{K} f_k
$$

Taking derivatives with respect to $f_k$ and $\lambda$ yields

$$
\begin{cases}
\dfrac{\partial \mathscr{L}(f, \lambda)}{\partial f_k} = -\dfrac{1}{K-1} \exp\left\{ -\dfrac{f_k^*}{K-1} \right\} \mathbb{P}(Y_k = 1 \mid X = x) - \lambda = 0 \\[4mm]
\dfrac{\partial \mathscr{L}(f, \lambda)}{\partial \lambda} = -\sum_{k=1}^{K} f_k^* = 0
\end{cases}
$$

We get from the first relation that $\exp\left\{ -\dfrac{f_k^*}{K-1} \right\} \times \mathbb{P}(\cdots) = -\lambda(K-1)$

$$
\Updownarrow
$$

$$
f_k^* = (K-1)\log \mathbb{P}(Y_k = 1 \mid X = x)
$$
$$
- (K-1)\log\big[ -\lambda(K-1) \big]
$$

---

Enforcing that the $f_k^*$ sum to zero, we get

$$
\frac{1}{K-1} \sum_k f_k^* = \sum_k \log \mathbb{P}(Y_k = 1 \mid X = x) - \sum_k \log\big[ -\lambda(K-1) \big] = 0
$$

$$
\Longleftrightarrow
$$

$$
\lambda = -\frac{1}{K-1} \exp\left\{ \frac{1}{K} \sum_k \log \mathbb{P}(Y_k = 1 \mid X = x) \right\}
$$

Plugging this value of $\lambda$ back into the expression of $f_k^*$ yields the result.

Remark. Once the $f_k^*$ are (somehow) estimated, the posterior class probabilities can be computed:

$$
\mathbb{P}(Y_k = 1 \mid X = x) = \underset{\underset{\text{sum to one}}{\uparrow}}{\exp\left\{ \frac{f_k^*}{K-1} \right\}} \left( \prod_{j=1}^{K} \mathbb{P}(Y_j = 1 \mid X = x) \right)^{\frac{1}{K}}
$$

$$
1 = \left( \prod_{j=1}^{K} \mathbb{P}(Y_j = 1 \mid X = x) \right)^{\frac{1}{K}} \left( \sum_{k=1}^{K} \exp\left\{ \frac{f_k^*}{K-1} \right\} \right)
$$

Plugging this expression of $\left( \prod_{j=1}^{K} \mathbb{P}(\cdots) \right)^{\frac{1}{K}}$ back into the posterior proba gives:

$$
\boxed{ \mathbb{P}(Y_k = 1 \mid X = x) = \frac{\exp\left\{ \dfrac{f_k^*}{K-1} \right\}}{\sum_{j=1}^{K} \exp\left\{ \dfrac{f_j^*}{K-1} \right\}} }
$$

- AdaBoost algorithm.

  ↘ Classical AdaBoost $(K=2)$ : $f^{(n)}(x) = \sum_{\ell=1}^{M} \beta_\ell \underset{\underset{\in \{-1, 1\}}{\uparrow}}{f_\ell(x)}$

  ↑ soft classifier

  (or, more generally, in $[-1, 1]$)

↘ K-class AdaBoost: $f^{(M)}(x) = \sum_{\ell=1}^{M} \beta_\ell f_\ell(x)$

↑ Each weak classifier take values in

$$\left\{ \begin{array}{c} \left( 1, -\frac{1}{K-1}, \cdots, -\frac{1}{K-1} \right) \\ \left( -\frac{1}{K-1}, 1, \cdots, -\frac{1}{K-1} \right) \\ \vdots \\ \left( -\frac{1}{K-1}, -\frac{1}{K-1}, \cdots, 1 \right) \end{array} \right\}$$

The $(\beta_\ell, f_\ell)$ are fit in a stagewise manner:

$$(\beta_\ell, f_\ell) \in \underset{(\beta, f)}{\operatorname{argmin}} \sum_{i=1}^{n} \ell\left( y_i, f^{(\ell-1)}(x_i) + \beta f(x_i) \right)$$

↗ vectors in $\mathbb{R}^K$

$$= \underset{(\beta, f)}{\operatorname{argmin}} \sum_{i=1}^{n} \exp\left\{ -\frac{1}{K} y_i^t \left( f^{(\ell-1)}(x_i) + \beta f(x_i) \right) \right\}$$

$$\boxed{(\beta_\ell, f_\ell) = \underset{(\beta, f)}{\operatorname{argmin}} \sum_{i=1}^{n} w_i^{(\ell)} \exp\left\{ -\frac{1}{K} \beta y_i^t f(x_i) \right\}}$$

where $w_i^{(\ell)} = \exp\left\{ -\frac{1}{K} y_i^t f^{(\ell-1)}(x_i) \right\}$

We proceed as for binary classification:

$$\sum_{i=1}^{n} w_i^{(\ell)} \exp\left\{ -\frac{1}{K} \beta y_i^t f(x_i) \right\}$$

$$= \underset{\substack{\text{correctly} \\ \text{classified} \\ \text{observations}}}{\sum} w_i^{(\ell)} \exp\left\{ -\frac{1}{K} \beta y_i^t f(x_i) \right\} + \underset{\substack{\text{incorrectly} \\ \text{classified} \\ \text{obs}}}{\sum} w_i^{(\ell)} \exp\left\{ \cdots \right\}$$

More formally, with $y_i = (y_{i,1}, \cdots, y_{i,K})$
correctly classified obs $= \left\{ (x_i, y_i) \mid y_i = f(x_i) \right\}$

---

↪ For a correctly classified observation,
$$y_i^t f(x_i) = 1 + \underbrace{\left(\frac{1}{K-1}\right)^2 + \cdots + \left(\frac{1}{K-1}\right)^2}_{K-1} = \frac{K}{K-1}$$

↪ For an incorrectly classified observation,
$$y_i^t f(x_i) = -\frac{2}{K-1} + \underbrace{\left(\frac{1}{K-1}\right)^2 + \cdots + \left(\frac{1}{K-1}\right)^2}_{K-2} = -\frac{K}{(K-1)^2}$$

Thus,
$$\sum_{i=1}^{n} w_i^{(\ell)} \exp\left\{ -\frac{1}{K} \beta y_i^t f(x_i) \right\}$$

$$= \underset{\text{correct}}{\sum} w_i^{(\ell)} \exp\left( -\frac{\beta}{K-1} \right) + \underset{\text{incorrect}}{\sum} w_i^{(\ell)} \exp\left( \frac{\beta}{(K-1)^2} \right).$$

$$= e^{-\frac{\beta}{K-1}} \sum_{i=1}^{n} w_i^{(\ell)} + \left( e^{\frac{\beta}{(K-1)^2}} - e^{-\frac{\beta}{K-1}} \right) \underline{\sum w_i^{(\ell)} \mathbb{1}(y_i \neq f(x_i))}$$

Call this term $W^{(\ell)}$

To find $f_\ell$, construct a classifier minimizing a weighted 0/1 loss. Ex: grow a tree.

• It remains to find $\beta_\ell$.

$$W^{(\ell)} = e^{-\frac{\beta}{K-1}} \left( \sum_{i=1}^{n} w_i^{(\ell)} \right) \left[ 1 + \left( \frac{e^{\frac{\beta}{(K-1)^2}} - e^{-\frac{\beta}{K-1}}}{e^{-\frac{\beta}{K-1}}} \right) \underline{\frac{\sum w_i^{(\ell)} \mathbb{1}(y_i \neq f(x_i))}{\sum w_i^{(\ell)}}} \right]$$

call this term $err^{(\ell)}$.

$$= e^{-\frac{\beta}{K-1}} \left( \sum_{i=1}^{n} w_i^{(\ell)} \right) \left[ 1 + \left\{ \exp\left( \frac{\beta}{(K-1)^2} + \frac{\beta}{K-1} \right) - 1 \right\} err^{(\ell)} \right]$$

Thus
$$\beta_\ell = \underset{\beta}{\operatorname{argmin}} \, e^{-\frac{\beta}{K-1}} \left[ 1 + \left\{ \exp\left( \frac{\beta K}{(K-1)^2} \right) - 1 \right\} err^{(\ell)} \right]$$

Put $u = \exp\left(\frac{\beta}{K-1}\right)$, and consider

$$\Psi(u) = u^{-1}\left(1 + err^{(\ell)}\left[u^{K/(K-1)} - 1\right]\right)$$

Solve $\Psi'(u^*) = 0 \implies \frac{K}{K-1}\log u^* = \log(K-1) + \log\left(\frac{1-err^{(\ell)}}{err^{(\ell)}}\right)$

Thus

$$\boxed{\beta_\ell = \frac{(K-1)^2}{K}\left\{\log(K-1) + \log\left(\frac{1-err^{(\ell)}}{err^{(\ell)}}\right)\right\}}$$

$$= \frac{(K-1)^2}{K}\alpha_\ell$$

- **Update**: * $f^{(\ell)}(x) = f^{(\ell-1)}(x) + \beta_\ell f_\ell(x)$

  * $w_i^{(\ell+1)} = \exp\left\{-\frac{1}{K}y_i^t f^{(\ell)}(x_i)\right\}$

  $$= w_i^{(\ell)}\exp\left\{-\frac{1}{K}\beta_\ell\, y_i^t f_\ell(x_i)\right\}$$

  $$= \begin{cases} w_i^{(\ell)}\exp\left\{-\frac{K-1}{K}\alpha_\ell\right\} & \text{if } y_i = f_\ell(x_i) \\ w_i^{(\ell)}\exp\left\{\frac{\alpha_\ell}{K}\right\} & \text{if } y_i \neq f_\ell(x_i) \end{cases}$$

Multiply weights by $\exp\left(\frac{\alpha_\ell}{K}\right) = $ indpt of $i$

$$\overset{\Leftrightarrow}{update} \begin{cases} w_i^{(\ell)}\exp(-\alpha_\ell) & \text{if } y_i = f_\ell(x_i) \\ w_i^{(\ell)} & \text{if } y_i \neq f_\ell(x_i) \end{cases}$$

$$\overset{\Leftrightarrow}{update} \begin{cases} w_i^{(\ell)} & \text{if } y_i = f_\ell(x_i) \\ w_i^{(\ell)}\exp(\alpha_\ell) & \text{if } y_i \neq f_\ell(x_i) \end{cases}$$

$$\implies w_i^{(\ell+1)} = w_i^{(\ell)}\exp\left[\alpha_\ell \mathbb{1}(y_i \neq f_\ell(x_i))\right]$$

For the weights to be updated in the right direction, we need $\alpha_\ell > 0$; that is $\log(K-1) + \log\left(\frac{1-err^{(\ell)}}{err^{(\ell)}}\right) > 0$: in K-class classif, we do not longer require that the base classifier perform better than a random guess.

---

## AdaBoost (K-class ; $K \geqslant 3$)

① Initialize weights $w_i^{(1)} = \frac{1}{n}$

② For $\ell = 1$ to $M$,

→ Fit a classifier $f_\ell$ to the training sample with weights $w_i^{(\ell)}$.

→ Compute its error $err^{(\ell)} = \dfrac{\sum\limits_{i=1}^n w_i^{(\ell)}\mathbb{1}(y_i \neq f_\ell(x_i))}{\sum\limits_{i=1}^n w_i^{(\ell)}}$

→ $\alpha_\ell = \log(K-1) + \log\left(\frac{1-err^{(\ell)}}{err^{(\ell)}}\right)$

→ Update $w_i^{(\ell+1)} = w_i^{(\ell)}\exp\left[\alpha_\ell \mathbb{1}(y_i \neq f_\ell(x_i))\right]$

③ Output $f^{(M+1)}(x) = \sum\limits_{\ell=1}^M \alpha_\ell f_\ell(x)$

SAMME algorithm (Zhu et al 2009) Plugging in $\beta_\ell$ makes no difference.

## III - GRADIENT BOOSTING

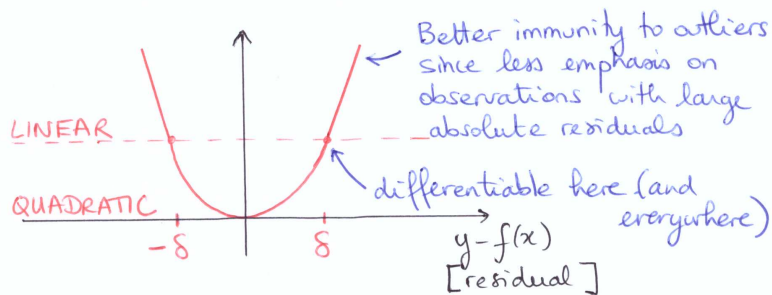### III.1. Gradient Boosting for regression

Gradient boosting algorithms are procedures minimizing general differentiable loss functions (other than the exponential loss in the context of classification).

Gradient → the gradient of the loss function is computed & implemented using a gradient descent-type algorithm

Boosting → the final (soft) classifier is expressed using a stagewise additive expansion.

**Remark:** In regression problems, the <u>HUBER</u> loss is a good alternative to the square loss (sensitive to outliers) and to the absolute loss (non-differentiable).
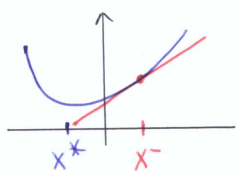
$$\ell(y, f(x)) = \begin{cases} (y - f(x))^2 & \text{if } |y - f(x)| < \delta \\ 2\delta |y - f(x)| - \delta^2 & \text{otherwise} \end{cases}$$



Better immunity to outliers since less emphasis on observations with large absolute residuals

LINEAR

QUADRATIC

$-\delta$ $\quad$ $\delta$

differentiable here (and everywhere)

$y - f(x)$ [residual]

- <u>Gradient descent</u>: relatively cheap & it works (even if better descent algorithms exist).
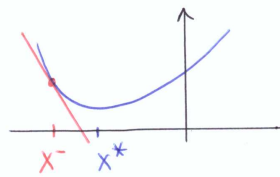
  <u>Goal</u>: find the minimizer $x^*$ of a function $\phi(x)$ ; $x \in \mathbb{R}^n$.
  <u>Idea</u>: move along the opposite direction of the gradient.
  $\Rightarrow$ Recursive procedure.

  <u>Indeed</u>,



$x^*$ $\quad$ $x^-$ $\qquad$ $x^-$ $\quad$ $x^*$

gradient $\geqslant 0$
$x^+ = x^- - \varrho \, (\text{positive number})$

$\Rightarrow$ decrease the value of your current approximation towards $x^*$

gradient $\leqslant 0$
$x^+ = x^- - \varrho \, (\text{neg } \#)$

$\Rightarrow$ increase the value of $x^-$ towards $x^*$

$\boxed{\varrho > 0}$

---

$\Rightarrow$ Update is $x^+ = x^- - \varrho \nabla_x \phi(x)\big|_{x = x^-}$

learning rate

$\multimap$ Convergence towards a local minimum
$\multimap$ $\varrho$ is updated at each iteration via a <u>line search</u>:

$$\varrho^- = \arg\min_\varrho \phi\left(x^- - \varrho \nabla_x \phi(x)\big|_{x = x^-}\right)$$

Search in a particular direction $\Rightarrow$ it should be OK in many cases to perform the minimization.

- <u>Gradient boosting for regression.</u>

  Consider a differentiable loss function $\ell$. Empirical Risk Minimization is the task of selecting a function $f$ in a class of candidates $\mathcal{F}$ minimizing the empirical risk

$$\boxed{\sum_{i=1}^n \ell(y_i, f(x_i))}$$

Forget for now that $f$ belongs to some $\mathcal{F}$. We incorporate this constraint later

Consider the vector $\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \in \mathbb{R}^n$. Gradient boosting algorithms compute a vector in $\mathbb{R}^n$ minimizing the function $\phi(z) := \sum_{i=1}^n \ell(y_i, z_i)$ , $z = (z_1, -, z_n)^t \in \mathbb{R}^n$

$\hookrightarrow$ Use Gradient Descent to minimize $\phi$:
Given the current estimate $z^{(\ell-1)}$, compute the gradient of $\phi$ evaluated at $z^{(\ell-1)}$:

$z^{(\ell-1)} = \begin{pmatrix} z_1^{(\ell-1)} \\ \vdots \\ z_n^{(\ell-1)} \end{pmatrix}$

$$\nabla \phi(z) = \left( \frac{\partial \ell(y_1, z_1)}{\partial z_1}, \cdots, \frac{\partial \ell(y_n, z_n)}{\partial z_n} \right)$$

$\Rightarrow z^{(\ell)} = z^{(\ell-1)} - \varrho_{\ell-1} \nabla \phi(z)\big|_{z = z^{(\ell-1)}}$

where $\quad \rho_{\ell+1} = \underset{\rho}{\text{argmin}} \; \phi \left( z^{(\ell-1)} - \rho \nabla \phi(z) \big|_{z = z^{(\ell-1)}} \right)$

$z^{(\ell-1)}$ represents our current estimate of $(y_1, \cdots, y_n)^t \in \mathbb{R}^n$

Denote it $\quad z^{(\ell-1)} = \left( f^{(\ell-1)}(x_1), \cdots, f^{(\ell-1)}(x_n) \right)^t$

$\quad (= (z_1^{(\ell-1)}, \cdots, z_n^{(\ell-1)})^t)$

This approach is however not satisfying; since we cannot generalize the prediction to unseen points $x$. The main idea of gradient boosting is that before updating the current estimate, we generalize the value of the gradient for any $x \in \mathbb{R}^d$ = and prevents us from overfitting.

Ex: in $\mathbb{R}^2$



your learning sample

At iteration $\ell$, the gradient of $\phi$ is computed at the training point $x_1, \cdots, x_5$

generalize the value of
$-\dfrac{\partial \ell(y_i, z_i)}{\partial z_i}\bigg|_{z_i = z_i^{(\ell-1)}}$ $\longrightarrow$
to all $x$ by viewing the task as a regression problem on
$\left( x_i, -\dfrac{\partial \ell(y_i, z_i)}{\partial z_i} \right)_{z_i = z_i^{(\ell-1)}}$

$-\dfrac{\partial \ell(y_5, z_5)}{\partial z_5}\bigg|_{z_5 \, = \, z_5^{(\ell-1)}}$

$-\dfrac{\partial \ell(y_1, z_1)}{\partial z_1}\bigg|_{z_n \, = \, z_1^{(\ell-1)}}$

---

Let $\quad h(x; \Theta)$ = regression function. The parameter $\Theta$ is estimated using the training sample
$$\left\{ \left( x_i, \; -\dfrac{\partial \ell(y_i, z_i)}{\partial z_i}\bigg|_{z_i = z_i^{(\ell-1)}} \right) \right\}_{i=1, -, n}$$

Ex: for a regression tree, $\Theta$ includes the split variables, the split points and the predicted value in each terminal region.

Typically, this regression problem is considered with a square loss function:
$$\Theta_\ell = \underset{\Theta}{\text{argmin}} \quad \left( -\dfrac{\partial \ell(y_i, z_i)}{\partial z_i}\bigg|_{z_i = z_i^{(\ell-1)}} - h(x_i, \Theta) \right)^2$$

Back to our example in $\mathbb{R}^2$, you may want to fit a decision tree with, say, three terminal nodes:



You have generalized the value of the gradient at the points $x_1, \cdots, x_n$, to all points in $\mathbb{R}^2$!

Once the gradient is estimated everywhere, it remains to perform a line search:
$$\rho_\ell = \underset{\rho}{\text{argmin}} \; \sum_{i=1}^{n} \ell \left( y_i, \; f^{(\ell-1)}(x_i) + \rho \, h(x_i; \Theta_\ell) \right)$$

and update
$$f^{(\ell)}(x) = f^{(\ell-1)}(x) + \rho_\ell \, h(x; \Theta_\ell).$$
Continue until convergence.

$\Rightarrow$ After M steps: $f^{(M)}(x) = \sum_{\ell=1}^{M} \rho_\ell \, h(x; \Theta_\ell)$ — additive expansion — BOOSTING! — weak learners

—— Gradient-Boost algorithm (Friedman 99) ——

1. Put $f^{(0)}(x) = \underset{\varrho}{\text{argmin}} \sum_{i=1}^{n} l(y_i, \varrho)$

2. For $l = 1$ to $M$, do
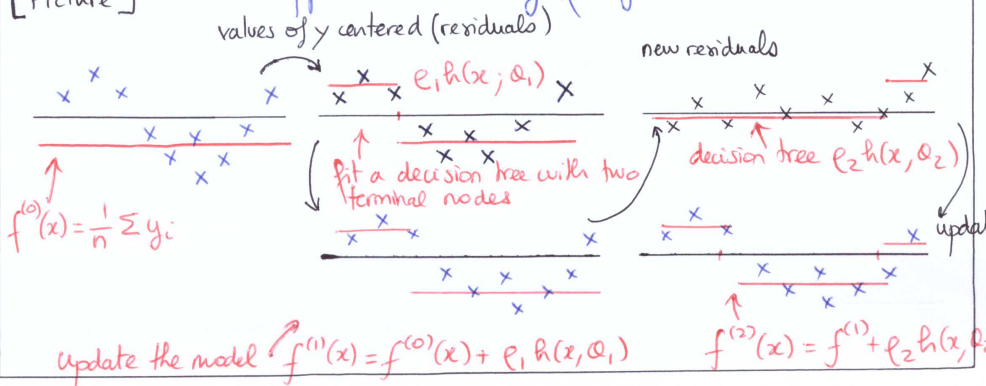
3. • Compute the pseudo-response $\tilde{y}_{il}$ :

$$\tilde{y}_{il} = -\left[ \frac{\partial l(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i) = f^{(l-1)}(x_i)}$$

4. • $\Theta_l = \underset{\Theta}{\text{argmin}} \sum_{i=1}^{n} \left( \tilde{y}_{il} - h(x_i; \Theta) \right)^2$

5. • $\varrho_l = \underset{\varrho}{\text{argmin}} \sum_{i=1}^{n} l\left( y_i, f^{(l-1)}(x_i) + \varrho\, h(x_i; \Theta_l) \right)$

6. • $f^{(l)}(x) = f^{(l-1)}(x) + \varrho_l\, h(x; \Theta_l)$

7. End for

8. End Algorithm

Examples: (i) <u>Square loss</u> $l(y, f(x)) = \frac{1}{2}(y - f(x))^2$

The pseudo response is $\tilde{y}_{il} = y_i - f^{(l-1)}(x_i)$

$\quad = $ residuals
$\quad = $ negative gradient

Gradient boosting with a square loss $\equiv$ stagewise approach iteratively fitting the current residuals.

[Picture]

values of y centered (residuals)   new residuals



$f^{(0)}(x) = \frac{1}{n}\sum y_i$

fit a decision tree with two terminal nodes

$\varrho_1 h(x; \Theta_1)$

decision tree $\varrho_2 h(x, \Theta_2)$

updal

update the model $f^{(1)}(x) = f^{(0)}(x) + \varrho_1 h(x, \Theta_1)$     $f^{(2)}(x) = f^{(1)} + \varrho_2 h(x, \Theta_2)$

---

(ii) <u>Absolute loss</u> $l(y, f(x)) = |y - f(x)|$

The pseudo response is $\tilde{y}_{il} = \text{sign}(y_i - f^{(l-1)}(x_i))$.

$\Rightarrow$ The regression function is fit to the sign of the residuals in line 4 of the algorithm.

Then, line 5 becomes

$\varrho_l = \underset{\varrho}{\text{argmin}} \sum_{i=1}^{n} \left| y_i - f^{(l-1)}(x_i) - \varrho\, h(x_i, \Theta_l) \right|$

$\quad = \underset{\varrho}{\text{argmin}} \sum_{i=1}^{n} \underbrace{|h(x_i, \Theta_l)|}_{} \left| \underbrace{\frac{y_i - f^{(l-1)}(x_i)}{h(x_i, \Theta_l)}}_{} - \varrho \right|$

$\quad = \underset{\varrho}{\text{argmin}} \sum_{i=1}^{n} \underbrace{w_i}_{} \left| \underbrace{z_i}_{} - \varrho \right|$

$\varrho_l = \underset{w}{\text{median}} \{z_1, .., z_n\}$

$\quad = $ weighted median of $\{z_1, .., z_n\}$ with weights $w_1, .., w_n$.

(iii) <u>Huber loss</u> $l(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & |y-f| \leq \delta \\ \delta(|y-f| - \frac{\delta}{2}) & \underline{\quad\quad} > \delta \end{cases}$

The pseudo response is

$\tilde{y}_{il} = \begin{cases} y_i - f^{(l-1)}(x_i) & \text{if } |y_i - f^{(l-1)}(x_i)| \leq \delta \\ \delta\, \text{sign}(y_i - f^{(l-1)}(x_i)) & \text{if } \underline{\quad\quad\quad} > \delta \end{cases}$

RK = $\delta$ may be allowed to vary with the iteration number:

$\delta_l = \alpha$-quantile of $\{y_i - f^{(l-1)}(x_i)\}_{i=1,..}$

$\quad = $ break-down point

and

$\varrho_l = \underset{\varrho}{\text{argmin}} \sum_{i=1}^{n} l(y_i, f^{(l-1)}(x_i) + \varrho\, h(x; \Theta_l))$

solve numerically; see Huber (1964)

In the case where the (weak) learners $h(x; \theta)$ are decision trees, the Gradient-Boost algorithm takes a simple form for specific loss functions.

Line 4 → parameter $\theta_\ell$ corresponds to the terminal regions

Line 5 → parameter $\rho_\ell$ corresponds to the predicted value in each terminal region.

Specifically,

---

**Gradient Tree-Boosting Algorithm**

1. Put $\quad f^{(0)}(x) = \underset{\rho}{\text{argmin}} \sum\limits_{i=1}^{n} \ell(y_i, \rho)$

2. For $\quad \ell = 1$ to $M$, do

3. $\qquad \bullet \; \tilde{y}_{i\ell} = -\left[\dfrac{\partial \ell(y_i, f(x_i))}{\partial f(x_i)}\right]_{f(x_i) = f^{(\ell-1)}(x_i)}$

4. $\qquad \bullet \; \underline{\text{Compute}}$ the $J$ terminal regions:
$\qquad\qquad \{R_{j\ell}\}_{j=1}^{J} = J\text{-terminal tree fitted to } (x_i, \tilde{y}_{i\ell}).$

5. $\qquad \bullet \; \underline{\text{Compute}}$
$\qquad\qquad \gamma_{j\ell} = \underset{\gamma}{\text{argmin}} \sum\limits_{x_i \in R_{j\ell}} \ell\left(y_i, f^{(\ell-1)}(x_i) + \gamma\right)$

6. $\qquad \bullet \; \underline{\text{Update}}$
$\qquad\qquad f^{(\ell)}(x) = f^{(\ell-1)}(x) + \sum\limits_{j=1}^{J} \gamma_{j\ell} \, \mathbb{1}(x \in R_{j\ell})$

7. End For

8. End Algorithm

---

**Ex**: With an absolute loss, line 5 becomes:
$$\gamma_{j\ell} = \underset{x_i \in R_{j\ell}}{\text{median}} \left\{ y_i - f^{(\ell-1)}(x_i) \right\}.$$

---

• **Context**: $K$-class classification

• **Coding scheme**: $Y_i = (Y_{i1}, \dots, Y_{iK})$, where $Y_{ik} = 1$ if observation $X_i$ belongs to class $k$, and $Y_{ik} = 0$ otherwise.

$\quad$ Put $\quad p_k(x_i) = P(Y_{ik} = 1 \mid X = x_i)$

$\qquad\qquad = \dfrac{\exp f_k(x_i)}{\sum\limits_{j=1}^{K} \exp f_j(x_i)}, \quad k = 1, \dots, K$

*softmax representation*

adding a constant value to all the $f_j$ does not change the value of $p_k$
$\Rightarrow$ impose that $\sum\limits_{j=1}^{K} f_j = 0$

• **Goal**: Estimation of $(f_1, \dots, f_K) =: f$

• **loss**: the multinomial deviance (for $n$ obs) [minus the log lik]

$\ell(y, f(x)) = -\sum\limits_{i=1}^{n} \sum\limits_{k=1}^{K} Y_{ik} \log p_k(X_i)$

$\qquad\qquad = -\sum\limits_{i,k} Y_{ik} f_k(X_i) - \sum\limits_{i=1}^{n} \underbrace{\sum\limits_{k=1}^{K} Y_{ik}}_{=1} \log\left[\sum\limits_{j=1}^{K} f_j(X_i)\right]$

$\qquad\qquad = -\sum\limits_{i,k} Y_{ik} f_k(X_i) - \sum\limits_{i} \log\left(\sum\limits_{j} f_j(X_i)\right)$

• **Pseudo-residuals**

$\tilde{y}_{ik\ell} = \dfrac{\partial}{\partial f_k(x_i)}\Big[-\ell(y_i, f_k(x_i))\Big]\Big|_{f_k(x_i) = f_k^{(\ell-1)}(x_i)}$

obs $i$ / class $k$ / iteration $\ell$

$\qquad\quad = y_{ik} - p_{k, \ell-1}(x_i),$

$\qquad\qquad$ where $p_{k,\ell-1}(x_i) := \dfrac{\exp f_k^{(\ell-1)}(x_i)}{\sum\limits_{j=1}^{K} \exp f_j^{(\ell-1)}(x_i)}.$

—— Gradient Boosting for K-class classification ——

1. Initialise $f_k^{(0)}(x)$ for $k=1, -, K$

2. For $\ell=1$ to $M$, repeat

3.     For $k=1$ to $K$, repeat

4.         • Compute the pseudo response

$$\tilde{y}_{ik\ell} = y_{ik} - p_{k,\ell-1}(x_i), \qquad i=1,-,n$$

$$\text{where } p_{k,\ell-1}(x_i) = \frac{\exp f_k^{(\ell-1)}(x_i)}{\sum_{j=1}^{K} \exp f_j^{(\ell-1)}(x_i)}$$

5.         • Fit a decision tree $T_{k\ell}$ to $\{(x_1, \tilde{y}_{1k\ell}), \cdots, (x_n, \tilde{y}_{nk\ell})\}$, producing regions $\{R_{jk\ell}\}_{j=1}^{J}$

6.         • Compute

$$\{\gamma_{jk\ell}\} = \arg\min_{\gamma} \sum_{x_i \in R_{jk\ell}} \ell(y_i, f_k^{(\ell-1)}(x_i) + \gamma)$$

7.         • Update

$$f_k^{(\ell)}(x) = f_k^{(\ell-1)}(x) + \sum_{j=1}^{J} \gamma_{jk\ell} \, \mathbb{1}(x \in R_{jk\ell})$$

8.     End For

9. End For

10. Compute the posterior probabilities $\hat{p}_k(x) = \dfrac{\exp f_k^{(M)}(x)}{\sum_{j=1}^{K} \exp f_j^{(M)}(x)}$,

and classify a new observation $x$ according to $\arg\max_{1 \leq k \leq K} \hat{p}_k(x)$.

Remark: Computing the values $\gamma_{jk\ell}$ in line 6 is not trivial, and can be approximately done using a one-step Newton-Raphson update using only the diagonal of the Hessian matrix of $\ell$, see (Friedman 99)

References

Freund Y and Schapire R (1996). Experiments with a new boosting algorithm. In: Machine Learning: Proceedings of the Thirteenth International Conference, 148-156.

Zhu J, Zou H, Rosset S, and Hastie T (2009). Multi-class AdaBoost. Statistics and its Interface. Vol 2, p. 349-360.

Friedman J (2001). Greedy Function Approximation: A Gradient Boosting Machine. Annals of Statistics, Vol 29, No 5, 1189-1232.