# Binary Classification, Part 3
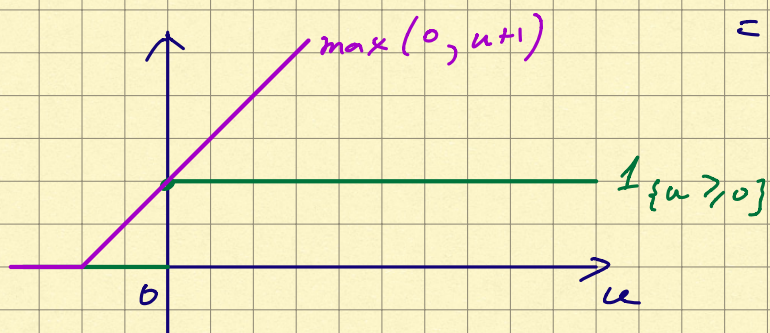
**Review:**

$$f: \mathcal{X} \to \mathbb{R}$$

$$g_f(x) = \text{sgn} \, f(x) = \begin{cases} +1, & f(x) \geq 0 \\ -1, & f(x) < 0 \end{cases}$$

$$(X, Y) \text{ in } \mathcal{X} \times \{-1, +1\}$$

$$L(g_f) = \mathbb{P}[Y \neq g_f(X)] \leq \mathbb{P}[Y f(x) \leq 0]$$

$$= \mathbb{E}[1_{\{-Y f(x) \geq 0\}}]$$



$\max(0, u+1)$

$1_{\{u \geq 0\}}$

penalty fcn $\varphi: \mathbb{R} \to \mathbb{R}_+$
continuous
nondecreasing
$\varphi(u) \geq 1_{\{u \geq 0\}}$

$\varphi(\cdot) \longrightarrow$ surrogate loss

$$\ell_\varphi(y, u) := \varphi(-yu)$$

$$L(g_f) = L(\text{sgn} \, f) \leq \mathbb{E}[1_{\{-Y f(x) \geq 0\}}]$$

$$\leq \mathbb{E}[\varphi(-Y f(x))]$$

$$= \mathbb{E}[\ell_\varphi(Y, f(X))]$$

$$=: A_\varphi(f) \qquad : \text{surrogate loss of } f$$

$$L(g_f) \leq A_\varphi(f)$$

$$L_n(g_f) \leq A_{\varphi,n}(f) \qquad \text{where, e.g., } A_{\varphi,n}(f) = \frac{1}{n} \sum_{i=1}^{n} \varphi(-Y_i f(X_i))$$

**Thm** (Koltchinskii - Panchenko) Let $\varphi: \mathbb{R} \to \mathbb{R}$ be a penalty fcn s.t.:

- $\varphi(-y f(x)) \in [0, 1]$ for all $(x, y)$, all $f \in \mathcal{F}$
- $\varphi$ is $M_\varphi$-Lipschitz: $|\varphi(u) - \varphi(v)| \leq M_\varphi |u - v|$

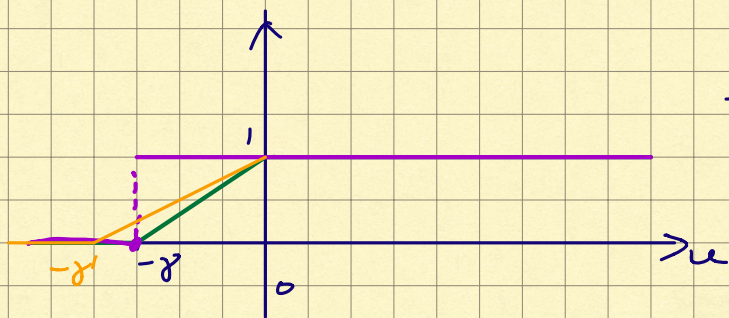Let $\hat{f}_n$ be any element of $\mathcal{F}$, based on data.

Then, w.p. $\geq 1 - e^{-2t^2}$ $(t>0)$,

$$L(\text{sgn} \bar{f_n}) \leq \underbrace{A_{\varphi,n}(\bar{f_n})}_{\substack{\text{computable} \\ \text{from data}}} + 4M_\varphi \underbrace{\mathbb{E}R_n(\mathcal{F}(x^n))}_{\substack{\text{typically easy} \\ \text{to bound}}} + \frac{t}{\sqrt{n}}.$$

<u>Example</u> (ramp penalty + margin)

$\gamma > 0$

$$\varphi(u) := \begin{cases} 0, & u < -\gamma \\ 1 + u/\gamma, & -\gamma \leq u < 0 \\ 1, & u \geq 0 \end{cases}$$



$$1_{\{u > -\gamma\}} \geq \varphi(u) \geq 1_{\{u \geq 0\}}$$

— for any $f: \mathcal{X} \to \mathbb{R}$,

$$L(\text{sgn} f) \leq A_\varphi(f) \leq L^\gamma(f),$$

where $L^\gamma(f) := \mathbb{P}[Y f(X) < \gamma]$

$Y f(X)$ : ==margin of $f$== on $(X, Y)$

$$L^\gamma(f) = \underbrace{\mathbb{P}[Y f(X) < 0]}_{\mathbb{P}[Y \neq \text{sgn} f(X)]} + \underbrace{\mathbb{P}[0 \leq Y f(X) < \gamma]}_{\text{prob. of margin} < \gamma}$$

$\varphi(\cdot)$ bdd between $[0,1]$, $\frac{1}{\gamma}$-Lipschitz

<u>Corollary</u>　　For any $\hat{f}_n$, w.p. $\geq 1 - e^{-2t^2}$,

$$L(\mathrm{sgn}\,\hat{f}_n) \leq L_n^{\gamma}(\hat{f}_n) + \frac{4}{\gamma} \mathbb{E}R_n(\mathcal{F}(x^n)) + \frac{t}{\sqrt{n}}$$

　　　　　　　　　increases w. $\gamma$　　　　decreases w. $\gamma$

— would like to make $\gamma$ data-dependent !

$$L(\mathrm{sgn}\,\hat{f}_n) \lesssim \inf_{\gamma \in (0,1]} \left\{ L_n^{\gamma}(\hat{f}_n) + \frac{c}{\gamma} R_n + \cdots \right\}$$
$$\text{w.p.} \geq 1 - e^{-O(t^2)} \qquad + \frac{t}{\sqrt{n}}$$

<u>Thm</u> (K.–P.)　Let $\varphi$ be a pen fcn, which is:

- bdd between $0$ and $1$
- $1$–Lipschitz
- monotone : $1 \geq \gamma \geq \gamma' > 0 \implies \varphi(u/\gamma) \geq \varphi(u/\gamma')$ for all $u$.

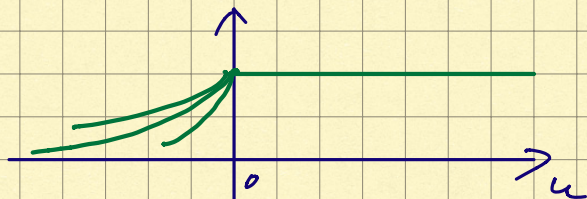Then, for any $\hat{f}_n \in \mathcal{F}$, w.p. $\geq 1 - 2e^{-2t^2}$,

$$L(\mathrm{sgn}\,\hat{f}_n) \leq \inf_{0 < \gamma \leq 1} \left\{ A_{\varphi(\cdot/\gamma),n}(\hat{f}_n) + \frac{c}{\gamma} \mathbb{E}R_n(\mathcal{F}(x^n)) \right.$$
$$\left. + C\sqrt{\frac{\log\log(2/\gamma)}{n}} \right\} + \frac{t}{\sqrt{n}}.$$

(cf. lecture notes for exact constants + proof)

<u>Examples of $\varphi$</u> :

1) ramp,　$\varphi(u) = \min\{1, \max\{0, u+1\}\}$

2) truncated exp,　$\varphi(u) = \min\{1, e^u\}$



$$\varphi(u/\gamma) = \min\{1, e^{u/\gamma}\}$$

# 1) Generalized majority vote

$\mathcal{G}$ : fixed collection of classifiers, $g: \mathcal{X} \to \{\pm 1\}$

$V(\mathcal{G}) < \infty$        (base classifiers)

Fix $\lambda > 0$

$$\bar{\mathcal{F}}_\lambda := \left\{ \sum_{k=1}^{N} c_k \, g_k(x) : N \in \mathbb{N}, \ |c_1| + \cdots + |c_N| \leq \lambda, \right\}$$
$$g_1, \ldots, g_N \in \mathcal{G}$$

– gen. maj. vote:     $c_1 = \cdots = c_N = \lambda/N$

<u>Note:</u> $\bar{\mathcal{F}}_\lambda$ may not be a VC class    (unlike $\mathcal{G}$)

$$R_n(\bar{\mathcal{F}}_\lambda(X^n)) = R_n(\lambda \cdot \text{absconv}(\mathcal{G}(X^n)))$$
$$= \lambda \cdot R_n(\text{absconv}(\mathcal{G}(X^n)))$$
$$= \lambda \cdot R_n(\mathcal{G}(X^n))$$
$$\leq C\lambda \cdot \sqrt{\frac{V(\mathcal{G})}{n}}$$

# 2) AdaBoost (Y. Freund – R. Schapire, 1997)

$\mathcal{G}$ : collection of base classifiers    $g: \mathcal{X} \to \{-1, +1\}$
    (weak learners)

$\mathcal{F} = \text{conv}(\mathcal{G})$

Data:    $(X_1, Y_1), \ldots, (X_n, Y_n)$ iid, in $\mathcal{X} \times \{-1, +1\}$

Algo:   iterative update of classifiers in $\text{conv}(\mathcal{G})$

$K \geq 1$ iterations

__Init:__ $w^{(1)} = (w_1^{(1)}, ..., w_n^{(1)})$ , $w_i^{(1)} = 1/n \ \forall i$

for $k = 1, ..., K$ :

- $e_k(g) := \sum_{i=1}^{n} w_i^{(k)} 1_{\{Y_i \neq g(X_i)\}}$ , $g \in \mathcal{G}$

  -- weighted class. error

Note: $k=1$    $e_1(g) = \frac{1}{n} \sum_{i=1}^{n} 1_{\{Y_i \neq g(X_i)\}} = L_n(g)$

$g_k := \underset{g \in \mathcal{G}}{argmin} \ e_k(g)$

$e_k := e_k(g_k) = \underset{g \in \mathcal{G}}{min} \sum_{i=1}^{n} w_i^{(k)} 1_{\{Y_i \neq g(X_i)\}}$

__Key Assumption:__    $e_k \leq 1/2$

- update $w^{(k)} \longrightarrow w^{(k+1)}$

$\forall i \in [n]: \quad w_i^{(k+1)} = \dfrac{w_i^{(k)} \exp(- \alpha_k Y_i g_k(X_i))}{Z_k}$

where $\alpha_k := \frac{1}{2} \log \frac{1-e_k}{e_k}$    $( \geq 0 )$    $\begin{array}{l} 0 \leq e_k \leq \frac{1}{2} \\ 1-e_k \geq e_k \end{array}$

$Z_k := \sum_{i=1}^{n} w_i^{(k)} \exp(- \alpha_k Y_i g_k(X_i))$

$\exp(- \alpha_k Y_i g(X_i)) = \begin{cases} e^{\alpha_k} & if \quad Y_i g_k(X_i) = -1 \\ e^{-\alpha_k} & if \quad Y_i g_k(X_i) = +1 \end{cases}$

- After $K$ steps, return

$\hat{f}_n(x) = \dfrac{\sum_{k=1}^{K} \alpha_k g_k(x)}{\sum_{k=1}^{K} \alpha_k}$    $\in conv(\mathcal{G})$

Note: $\quad \text{sgn} \, \hat{f}_n(x) = \text{sgn} \left( \sum_{k=1}^{K} \alpha_k \, g_k(x) \right)$

Preview: $\quad L(\hat{f}_n) \leq \prod_{k=1}^{K} 2 \sqrt{e_k(1 - e_k)}$

— if $2\sqrt{e_k(1-e_k)} < 1 \quad \forall k$, then error will decay with $K$!

— implicit margin minimization!

$$\gamma \sim \frac{1}{\sum_{k=1}^{K} \alpha_k} .$$