

Binary Classification, Part 2

(Surrogate losses)

• (X, Y) in $\mathcal{X} \times \{-1, +1\}$ (instead of $\{0, 1\}$)

Sign-based classifiers:

\mathcal{F} - class of fcn's $f: \mathcal{X} \rightarrow \mathbb{R}$

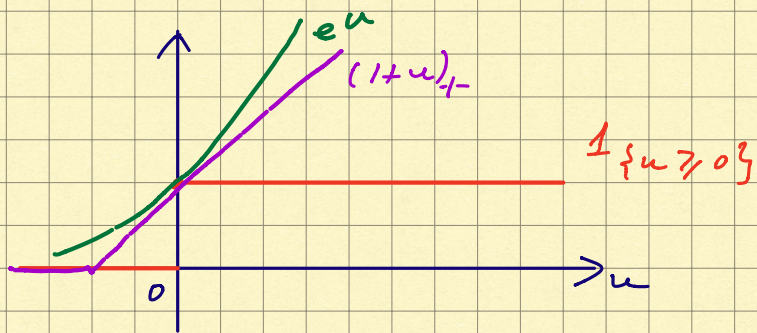
$$f \in \mathcal{F} \rightarrow g_f(x) := \text{sgn } f(x) = \begin{cases} 1, & f(x) \geq 0 \\ -1, & f(x) < 0 \end{cases}$$

$$\begin{aligned} L(g_f) &= P[Y \neq g_f(x)] \\ &= P[Y g_f(x) \leq 0] \\ &\leq P[Y f(x) \leq 0] \\ &= \mathbb{E}[1_{\{-Y f(x) \geq 0\}}] \end{aligned}$$

$$\begin{aligned} Y \text{sgn } f(x) \leq 0 \\ \Rightarrow Y f(x) \leq 0 \end{aligned}$$

Replace 0-1 loss $l(y, \hat{y}) = 1_{\{y \neq \hat{y}\}}$ with a tractable bound.

• Penalty fcn's / Surrogate losses



$$\varphi(u) = e^u$$

$$\varphi(u) = \log_2(1 + e^u)$$

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

$$\varphi(u) = \min\{1, (1 + u)_+\}$$

penalty fcn $\varphi: \mathbb{R} \rightarrow \mathbb{R}$

a) nondecreasing

b) continuous

c) $\varphi(u) \geq 1_{\{u \geq 0\}}$

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

$$L(g_f) = \mathbb{P}[Y \neq \text{sgn } f(x)]$$

$$\leq \mathbb{P}[Y f(x) \leq 0]$$

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

$$\leq \mathbb{E}[\varphi(-Y f(x))] \quad \text{for any pen fcn } \varphi$$

$$=: A_\varphi(f) \quad \text{— surrogate loss of } f \text{ (for pen fcn } \varphi)$$

$$l(y, \hat{y}) = \mathbb{1}_{\{y \neq \hat{y}\}} \longrightarrow l_\varphi(y, u) := \varphi(-yu)$$

$$A_\varphi(f) = \mathbb{E}[l_\varphi(Y, f(x))]$$

Empirical surrogate loss:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

$$A_{\varphi, n}(f) := \frac{1}{n} \sum_{i=1}^n \varphi(-y_i f(x_i))$$

$$\geq \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{-y_i f(x_i) \geq 0\}}$$

$$\geq \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i \neq \text{sgn } f(x_i)\}} = L_n(g_f)$$

easier to
min. over f ,
easier to
analyze

Analysis plan:

1) ERM for surrogate loss
— excess risk bounds, symmetrization, Rademacher

2) for any $\hat{f}_n \in \mathcal{F}$,

$$L(\text{sgn } \hat{f}_n) \leq A_\varphi(\hat{f}_n) \leq \dots$$

e.g. $L(\text{sgn } \hat{f}_n) \leq \inf_{f \in \mathcal{F}} A_\varphi(f) + \text{Rad}(\dots)$

4) Under additional mild assumptions on φ ,
 $\text{Rad}(\dots) = \text{Rad}(\mathcal{F})$

ERM: $\hat{f}_n = \underset{f \in \mathcal{F}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \varphi(-Y_i f(x_i))$

$$L(\text{sgn } \hat{f}_n) \leq A_\varphi(\hat{f}_n) \leq \inf_{f \in \mathcal{F}} A_\varphi(f) + 2 \sup_{f \in \mathcal{F}} |A_\varphi(f) - A_{\varphi,n}(f)|$$

$$\begin{aligned} A_\varphi(\hat{f}_n) - \inf_{f \in \mathcal{F}} A_\varphi(f) &= A_\varphi(\hat{f}_n) - A_{\varphi,n}(\hat{f}_n) + \underbrace{A_{\varphi,n}(\hat{f}_n) - A_{\varphi,n}(f^*)}_{\leq 0} + A_{\varphi,n}(f^*) - A_\varphi(f^*) \\ &\leq 2 \sup_{f \in \mathcal{F}} |A_\varphi(f) - A_{\varphi,n}(f)| \end{aligned}$$

assume $\exists f^* \in \mathcal{F}$ s.t.
 $A_\varphi(f^*) \leq A_\varphi(f), \forall f \in \mathcal{F}$

$$\Delta_n(\mathcal{Z}^n) := \sup_{f \in \mathcal{F}} |A_\varphi(f) - A_{\varphi,n}(f)|$$

$z_i = (x_i, y_i) \quad i \in [n]$

Symmetrization:

$$\mathbb{E} \Delta_n(\mathcal{Z}^n) \leq 2 \text{Rad}_n(\dots)$$

$$(x, y) \mapsto \varphi(-y f(x)), \quad f \in \mathcal{F}$$

$\text{Rad}_n(\dots)$:

$$\left\{ \left(\varphi(-y_1 f(x_1)), \dots, \varphi(-y_n f(x_n)) \right) : f \in \mathcal{F} \right\}$$

$$\hookrightarrow \left\{ (f(x_1), \dots, f(x_n)) : f \in \mathcal{F} \right\} - \text{contraction principle}$$

Lemma Suppose that φ is Lipschitz-cont.:
 $\exists M_\varphi > 0$ s.t.

$$|\varphi(u) - \varphi(v)| \leq M |u - v| \quad \text{for all } u, v \in \mathbb{R}.$$

Let \mathcal{H}_φ be the class of fcn's

$$(x, y) \mapsto \varphi(-y f(x)) - \varphi(0), \quad f \in \mathcal{F}.$$

$$\text{then } R_n(\mathcal{H}_\varphi(z^n)) \leq 2 M_\varphi R_n(\mathcal{F}(X^n)).$$

Proof

o) Contraction principle (Prop. 6.2 in lec. notes)

let \mathcal{H} be a class of fcn's $h: \mathcal{Z} \rightarrow \mathbb{R}$, and
 let $F: \mathbb{R} \rightarrow \mathbb{R}$ be a Lip. cont. fcn, s.t. $F(0) = 0$.

let $F \circ \mathcal{H}$ be the class of fcn's $z \mapsto F \circ h(z)$, $h \in \mathcal{H}$

$$\mathcal{Z} \xrightarrow{h \in \mathcal{H}} \mathbb{R} \xrightarrow{F} \mathbb{R}$$

$$\text{Then } R_n(F \circ \mathcal{H}) \leq 2 \|F\|_{\text{Lip}} R_n(\mathcal{H})$$

where $\|F\|_{\text{Lip}}$ is the Lipschitz const. of F .

$$1) A_{\varphi, n}(f) = \frac{1}{n} \sum_{i=1}^n \varphi(-y_i f(x_i))$$

$$A_\varphi(f) = \mathbb{E}[\varphi(-Y f(X))]$$

$$\Delta_n(z^n) = \sup_{f \in \mathcal{F}} |A_{\varphi, n}(f) - A_\varphi(f)|$$

$$= \sup_{f \in \mathcal{F}} |P_n(\ell_{\varphi, f}) - P(\ell_{\varphi, f})|$$

$$\ell_{\varphi, f}(x, y) := \varphi(-y f(x))$$

$$P_n(\ell_{\varphi, f}) - P(\ell_{\varphi, f}) = P_n(\ell_{\varphi, f} - \varphi(0)) - P(\ell_{\varphi, f} - \varphi(0))$$

$$\therefore \Delta_n(z^n) = \sup_{h \in \mathcal{H}_\varphi} |P_n(h) - P(h)|$$

where $\mathcal{H}_\varphi = \{ l_{\varphi, f} - \varphi(0) : f \in \mathcal{F} \}$

if $f(x) = 0$, then $l_{\varphi, f}(x, y) - \varphi(0) = 0$

$\mathcal{H} := \{ (x, y) \mapsto -y f(x) : f \in \mathcal{F} \}$

in fact, $\mathcal{H}_\varphi = F \circ \mathcal{H}$, where

$$F(v) := \varphi(v) - \varphi(0)$$

$$|F(u) - F(v)| = |\varphi(u) - \varphi(v)| \leq M_\varphi |u - v|$$

$$F(0) = \varphi(0) - \varphi(0) = 0$$

$$\Rightarrow \boxed{R_n(\mathcal{H}_\varphi(z^n)) \leq 2M_\varphi R_n(\mathcal{H}(z^n))}$$

by the contraction principle.

2) $R_n(\mathcal{H}(z^n))$

condition on z^n !!!

$$= \frac{1}{n} \mathbb{E}_\varepsilon \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \varepsilon_i \gamma_i f(x_i) \right| \right\}$$

$$= \frac{1}{n} \mathbb{E}_\varepsilon \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right\}$$

$$= R_n(\mathcal{F}(X^n))$$

$\varepsilon_1, \dots, \varepsilon_n \stackrel{i.i.d.}{\sim} \text{Rad}$

$\varepsilon_i \gamma_i, \dots, \varepsilon_n \gamma_n \stackrel{i.i.d.}{\sim} \text{Rad}$
for any fixed $\gamma_1, \dots, \gamma_n \in \{\pm 1\}$



Key takeaway: $\mathbb{E} \left\{ \sup_{f \in \mathcal{F}} |A_{\varphi_n}(f) - A_\varphi(f)| \right\}$

$$\leq 2 \mathbb{E} R_n(\mathcal{H}_\varphi(z^n))$$

$$\leq 4M_\varphi \mathbb{E} R_n(\mathcal{F}(X^n)).$$

Thm (Koltchinskii-Panchenko, 2002)

Let φ be a penalty fcn which is M_φ -Lipschitz, and assume $\exists B < \infty$ s.t.

$$\varphi(-y f(x)) \in [0, B] \quad \text{for all } (x, y) \in \mathcal{X} \times \mathcal{Y} \text{ and all } f \in \mathcal{F}.$$

Then, $\forall t > 0$, the surrogate ERM classifier \hat{f}_n satisfies

$$\begin{aligned} L(\text{sgn } \hat{f}_n) &\leq A_\varphi(\hat{f}_n) \\ &\leq \inf_{f \in \mathcal{F}} A_\varphi(f) + 8M_\varphi \mathbb{E} R_n(\mathcal{F}(z^n)) + \frac{2Bt}{\sqrt{n}} \end{aligned}$$

$$\text{w. p. } \geq 1 - e^{-2t^2}.$$

$$\text{NB: } \begin{aligned} e^{-2t^2} &\leq \delta \\ e^{2t^2} &\geq 1/\delta \end{aligned}$$

$$2t^2 \geq \log(1/\delta)$$

$$t = \sqrt{\frac{1}{2} \log(1/\delta)}$$

$$\Rightarrow \frac{2Bt}{\sqrt{n}} = B \sqrt{\frac{2 \log(1/\delta)}{n}}$$

Proof

$$\begin{aligned} L(\text{sgn } \hat{f}_n) &\leq A_\varphi(\hat{f}_n) \\ &\leq \inf_{f \in \mathcal{F}} A_\varphi(f) + 2\Delta_n(z^n) \end{aligned}$$

$$\mathbb{E} \Delta_n(z^n) \leq 4M_\varphi \mathbb{E} R_n(\mathcal{F}(X^n))$$

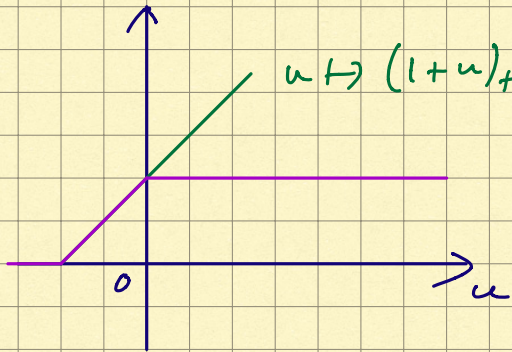
$$\mathbb{P} \left\{ \Delta_n(z^n) \geq \mathbb{E} \Delta_n(z^n) + \frac{Bt}{\sqrt{n}} \right\} \leq e^{-2t^2}$$

$$\text{— McDiarmid: } 0 \leq \varphi(-y f(x)) \leq B$$

$$\Delta_n(z^n) \text{ has bdd dif. } c_1 = \dots = c_n = B/n$$

$$P \left\{ \Delta_n(z^n) \geq \mathbb{E} \Delta_n(z^n) + \frac{Bt}{\sqrt{n}} \right\}$$

$$\leq \exp \left(- \frac{2B^2 t^2}{\cancel{\pi} \cdot \frac{B^2}{\cancel{\pi}}} \right); \text{ by McDiarmid.}$$



$$u \mapsto (1+u)_+ = \max(0, 1+u)$$

$$u \mapsto \min(1, (1+u)_+)$$

