

Regularization and Convex Surrogate Loss

§1

Motivation

①

$$\mathcal{H} = \{h_k\}_{k \in \mathbb{N}} \text{ s.t. } h_k \subseteq \mathcal{H}_{\text{all}}$$

$$\hat{R}_s(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq y_i\}}$$

→ Regularization

②

• Extension to \mathcal{H}_r , where $r \in \mathbb{R}_{>0}$:

$\{\mathcal{H}_r\}_{r \in \mathbb{R}_{>0}}$ s.t. \mathcal{H}_r varies "continuously" w.r.t. r .

• Alternative of $\hat{R}_s(h)$ that is easier to optimize.

↳ Convex surrogate loss.

§2

Regulazierung

①

$$\{\mathcal{H}_r\}_{r > 0}$$

• $\mathcal{H}_r = \{x \mapsto \text{sgn}(\langle w, \Phi(x) \rangle) \mid w \in \mathbb{R}^k, \|w\|_p \leq r\}$.

$\Phi: x \rightarrow \mathbb{R}^k$... feature map.

$$\|w\|_p := \left(\sum_{i=1}^k |w_i|^p \right)^{\frac{1}{p}}$$

• $\lambda > 0$: Hyperparameter.

$$w \in \underset{w}{\operatorname{argmin}} \hat{R}_s(h_w) + \lambda \|w\|_p$$

§3

Convex Surrogate Loss

①

Assumption

$\forall h \in \mathcal{H}, \exists f: \mathcal{X} \rightarrow \mathbb{R}$ s.t. $h(x) = \text{sgn}(f(x))$.

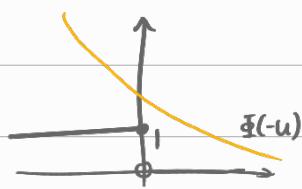
$$\text{sgn}(s) = \begin{cases} 1 & \text{if } s \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

Denote $h_f := \text{sgn} \circ f$.

(2)

Convex Surrogate Loss: $\Phi: \mathbb{R} \rightarrow \mathbb{R}$ s.t.

- Φ is convex
- Φ is nondecreasing
- $\forall u \in \mathbb{R}, \mathbb{1}_{\{u \leq 0\}} \leq \Phi(-u)$.



$$\Phi(u) = \max(0, 1+u) \quad \cdots \text{hinge loss}$$

$$\Phi(u) = \exp(u) \quad \cdots \text{exponential loss}.$$

$$\Phi(u) = \log_2(1+e^u) \quad \cdots \text{logistic loss}.$$

(3)

Why $\mathbb{1}_{\{u \leq 0\}} \leq \Phi(-u)$?

$$\begin{aligned} \hat{R}_s(h) &= \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) + y_i \leq 0\}} \\ &\circ \mathbb{1}_{\{h_f(x_i) + y_i \leq 0\}} = \mathbb{1}_{\{\text{sgn}(f(x_i)) + y_i \leq 0\}} \\ &\leq \mathbb{1}_{\{f(x_i)y_i \leq 0\}} \\ &\leq \Phi(-f(x_i)y_i) \end{aligned}$$

$$\text{Thus, } \hat{R}_s(h) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h_f(x_i) + y_i \leq 0\}} \leq \frac{1}{m} \sum_{i=1}^m \Phi(-f(x_i)y_i) =: \hat{R}_s^\Phi(f).$$

Lemma

 \hat{R}_s^Φ is convex.

(pf)

$$\begin{aligned} f, g \in [X \rightarrow \mathbb{R}], \alpha \in (0, 1), \text{ then } \hat{R}_s^\Phi(\alpha f + (1-\alpha)g) &\leq \alpha \hat{R}_s^\Phi(f) + (1-\alpha) \hat{R}_s^\Phi(g). \\ \hat{R}_s^\Phi(\alpha f + (1-\alpha)g) &= \frac{1}{m} \sum_{i=1}^m \Phi(-(\alpha f(x_i) + (1-\alpha)g(x_i))y_i) \\ &= \frac{1}{m} \sum_{i=1}^m \Phi(\alpha(-f(x_i)y_i) + (1-\alpha)(-g(x_i)y_i)) \\ &\leq \frac{1}{m} \sum_{i=1}^m [\alpha \Phi(-f(x_i)y_i) + (1-\alpha) \Phi(-g(x_i)y_i)] \\ &= \alpha \hat{R}_s^\Phi(f) + (1-\alpha) \hat{R}_s^\Phi(g). \end{aligned}$$

(4)

Bayes' Scoring Function $f^* \cdots R$

$$\text{Let } \mathbb{E}_{S \sim D^m} \hat{R}_s^\Phi(f) = R^\Phi(f).$$

$$R(h_{f^*}) = R(f^*) \text{ minimal}$$

$$R^\Phi(f^*) \text{ minimal}$$

$$R(f^*)$$

||

$$R(h_f) - R(h_{\text{Bayes}}) \quad (\text{Excess Error})$$

$$R^{\Phi}(h_f) - R^{\Phi}(f^*)$$

] Relation?

Proposition

$$f^*(x) = \eta(x) - \frac{1}{2}, \text{ where } \eta(x) = P(y=1|x).$$

$$f^*(x) = \underset{u \in [-\infty, \infty]}{\operatorname{argmin}} \eta(x)\Phi(-u) + (1-\eta(x))\Phi(u)$$

(pf)

$$\textcircled{1} \quad R(h_f) \geq R(h_{f^*}) \quad \forall f.$$

$$\begin{aligned} \therefore R(h_f) &= E[\mathbb{1}_{\{h_f(x) \neq y\}}] \rightarrow E[\mathbb{1}_{\{h_f(x)=1 \wedge y=1\}}] \\ &= E[E[\mathbb{1}_{\{h_f(x) \neq y\}} | x]] \quad + E[\mathbb{1}_{\{h_f(x)=-1 \wedge y=1\}}] \\ &= E[P(\operatorname{sgn}(f(x)) \neq y) | x] \quad (\because \eta(x) > \frac{1}{2} \rightarrow \operatorname{sgn}(f^*(x)) = 1) \\ &\geq E[P(\operatorname{sgn}(f^*(x)) \neq y) | x]. \quad (P(y=1|x) > \frac{1}{2}) \\ &= E[E[\mathbb{1}_{\{h_{f^*}(x) \neq y\}} | x]] \quad \Rightarrow P(y=1|x) > P(y=-1|x) \\ &= E[\mathbb{1}_{\{h_{f^*}(x) \neq y\}}] \quad \Rightarrow P(\operatorname{sgn}(f^*(x)) = y | x) \\ &= R(h_{f^*}) \quad \geq P(\operatorname{sgn}(f(x)) = y | x), \\ & \quad \text{given } \operatorname{sgn}(f(x)) \in \{-1, 1\}. \end{aligned}$$

$$\textcircled{2} \quad R^{\Phi}(f) \geq R^{\Phi}(f^*) \quad \forall f.$$

$$\begin{aligned} \therefore R^{\Phi}(f) &= E[\Phi(-yf(x))] \\ &= E[E[\Phi(-yf(x)) | x]] \end{aligned}$$

It suffices to prove $E[\Phi(-yf(x)) | x] \geq E[\Phi(-yf^*(x)) | x]$
for $\forall x$.

$$\text{For any } u, \eta(x)\Phi(-f^*(x)) + (1-\eta(x))\Phi(f^*(x)) \leq \eta(x)\Phi(-u) + (1-\eta(x))\Phi(u).$$

$$\begin{aligned} E[\Phi(-yf(x)) | x] &= E[\Phi(-yf(x)) \mathbb{1}_{\{y=1\}} | x] + E[\Phi(-yf(x)) \mathbb{1}_{\{y=-1\}} | x] \\ &= E[\Phi(-f(x)) \mathbb{1}_{\{y=1\}} | x] + E[\Phi(f(x)) \mathbb{1}_{\{y=-1\}} | x] \\ &= \Phi(-f(x))\eta(x) + \Phi(f(x))(1-\eta(x)) \\ &\geq \eta(x)\Phi(-f^*(x)) + (1-\eta(x))\Phi(f^*(x)). \\ &= E[\Phi(-yf^*(x)) | x] \end{aligned}$$

$$\textcircled{3} \quad \Phi(x) = \max(0, 1+x), \quad \Psi(x) = \exp(x).$$

$$\min_u \eta(x) \max(0, 1-u) + (1-\eta(x)) \max(0, 1+u).$$

$$u=1 \Rightarrow 2(1-\eta(x))$$

$$u=-1 \Rightarrow 2\eta(x).$$

$$\Rightarrow f_{\Psi}^*(x) = \begin{cases} 1 & \text{if } \eta(x) \geq \frac{1}{2} \\ -1 & \text{if } \eta(x) < \frac{1}{2}. \end{cases}$$

$$\Phi(u) = \exp(u) \cdot \eta(u) e^{-u} + (1-\eta(u)) e^u.$$

$$\frac{\partial}{\partial u} (\Phi(u)) = 0 \Rightarrow f_{\Psi}^*(x) = \frac{1}{2} \log \frac{\eta(x)}{1-\eta(x)}$$

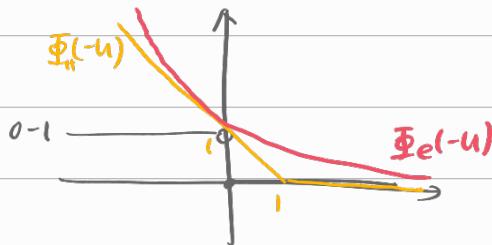
Examples

$$\textcircled{1} \quad \Phi_H(u) = \max\{0, 1+u\}$$

(Hinge loss)

$$\textcircled{2} \quad \Phi_e(u) = \exp(u)$$

(Exponential loss)



(Relation Between SUM and Hinge Loss)

$$H = \{x \mapsto \operatorname{sgn}(\langle w, x \rangle + b) \mid w \in \mathbb{R}^n, b \in \mathbb{R}\}$$

$$\operatorname{argmin}_{w,b} \hat{R}_S^{\Phi_H}(x \mapsto \langle w, x \rangle + b) + \frac{1}{2} \|w\|^2.$$

(3)

Risk

$$R(f) = R(h_f) = \mathbb{E}_{(x,y) \sim D} [\hat{R}_S(h_f)] = \mathbb{E}_{(x,y) \sim D} [1_{\{\operatorname{sgn}(f(x)) \neq y\}}]$$

Bayes' scoring function f^* s.t. $h_{f^*} = h^*$.

$$f^*(x) = \eta(x) - \frac{1}{2}, \text{ where } \eta(x) = P[y=1|x]$$

Bayes' error $R^* = R(h^*) = R(h_{f^*})$ **Φ -Risk**

$$R^{\Phi}(f) = \mathbb{E}_{(x,y) \sim D} [\Phi(-f(x)y)]$$

$$f_{\Phi}^*(x) = \operatorname{argmin}_{u \in [-\alpha, \alpha]} \eta(x) \Phi(-u) + (1-\eta(x)) \Phi(u).$$

$$\text{Then, } R^{\Phi}(f_{\Phi}^*) \leq R^{\Phi}(g) \quad \forall g.$$

$\frac{!!}{R^{\Phi,*}}$

{

Excess Error

$$R(f) - R^*$$

v.s.

$$R^{\Phi}(f) - R^{\Phi,*}$$

Theorem 4.7

$$\exists s \geq 1, c > 0 \text{ s.t. } \forall x \in \mathcal{X} \quad |f^*(x)|^s \leq c^s [\Phi(0) - \mathbb{E}_{\mathbb{P}[y|x]} [\Phi(-y f^{\Phi,*}(x))]]$$

Then for $\forall f \in \mathcal{K} \rightarrow \mathbb{R}$,

$$(R(f) - R^*)^s \leq (2c)^s [R^{\Phi}(f) - R^{\Phi,*}]$$

example) Hinge loss $s=1, c = \frac{1}{2}$

Exponential loss $s=2, c = \frac{1}{\sqrt{2}}$.

Remark

$$f^*(x) = \eta(x) - \frac{1}{2}$$

$$\eta(x) = \mathbb{P}[y=1|x]$$

$$f^{\Phi,*}(x) = \operatorname{argmin}_u \eta(x) \Phi(-u) + (1-\eta(x)) \Phi(u).$$

Lemma 4.5

$$R(f) - R^* \leq 2 \mathbb{E}_{x \sim D_{\mathcal{X}}} [|f^*(x)| \mathbb{1}_{\{f(x) f^*(x) \leq 0\}}]$$

(pf of lemma)

$$\begin{aligned} R(f) - R^* &= \mathbb{E}_{(x,y) \sim D} [\mathbb{1}_{\{f(x) \neq y\}} - \mathbb{1}_{\{f^*(x) \neq y\}}] \\ &\leq \eta(x) \mathbb{E}[\\ &\quad + (1-\eta(x)) \mathbb{E}[\end{aligned}$$

See below

(pf of thm 4.7) ① $\Phi(-2f(x)f^*(x)) \leq \mathbb{E}_{\mathbb{P}[y|x]} [\Phi(-y f(x))]$

$$\begin{aligned} \textcircled{2}) \quad \Phi(-2(\eta(x) - \frac{1}{2})f(x)) &= \Phi((1-2\eta(x))f(x)) \\ &= \Phi(\eta(x)(-f(x)) + (1-\eta(x))(f(x))) \\ &\leq \eta \Phi(-f(x)) + (1-\eta(x)) \Phi(-(-f(x))) \end{aligned}$$

$$= \mathbb{E}_{y|x} [\Phi(-y f(x))]$$

$$\begin{aligned}
② R(f) - R^* &\leq 2 \mathbb{E}_x [|f^*(x)| \mathbb{1}_{\{f(x)f^*(x) \leq 0\}}] \quad (\text{lemma 4.5}) \\
&= 2 \mathbb{E}_x [(|f^*(x)|^s \mathbb{1}_{\{f(x)f^*(x) \leq 0\}})^{1/s}] \quad (\text{if } f \in \{0, 1\}). \\
&\leq 2 (\mathbb{E}_x [|f^*(x)|^s \mathbb{1}_{\{f(x)f^*(x) \leq 0\}}])^{1/s} \quad (\text{Jensen for concave}). \\
&\leq 2 (\mathbb{E}_x [c^s (\Phi(0) - \mathbb{E}_{y|x} [\Phi(-y f^*(x))]) \mathbb{1}_{\{f(x)f^*(x) \leq 0\}}])^{1/s} \\
&\quad (\text{By assumption}) \\
&\leq 2 (\mathbb{E}_x [c^s (\Phi(-f(x)f^*(x))) - \mathbb{E}_{y|x} [\Phi(-y f^{\Phi,*}(x))]] \mathbb{1}_{\{f(x)f^*(x) \leq 0\}})^{1/s} \\
&\quad (\Phi \text{ is nondecreasing}) \\
&\leq 2 (\mathbb{E}_x [c^s (\Phi(-f(x)f^*(x)) - \mathbb{E}_{y|x} [\Phi(-y f^{\Phi,*}(x))])]^{1/s} \\
&\leq 2 (\mathbb{E}_x [c^s \mathbb{E}_{y|x} [\Phi(-y f(x))] - \mathbb{E}_{y|x} [\Phi(-y f^{\Phi,*}(x))]])^{1/s} \quad (\text{By ①}) \\
&= 2c \left(\mathbb{E}_{(x,y) \sim D} [\Phi(-y f(x)) - \mathbb{E}_{y|x} [\Phi(-y f^{\Phi,*}(x))]] \right)^{1/s} \\
&= 2c (R(f) - R^*)^{1/s}
\end{aligned}$$

(pf of lemma 4.5) ① $g \in [x \rightarrow \mathbb{R}]$, $R(g) = \mathbb{E}_{x \sim D_x} [2f^* \mathbb{1}_{\{g(x) < 0\}} + (1-\eta(x))]$

$$\begin{aligned}
\therefore R(g) &= \mathbb{E}_{x \sim D_x} \left[\mathbb{E}_{y|x} \left[\mathbb{1}_{\{y \neq \text{sgn}(g(x))\}} \right] \right] \\
&= \mathbb{E}_x \left[\mathbb{E}_{y|x} \left[\mathbb{1}_{\{y=-1, f(x) < 0\}} + \mathbb{1}_{\{y=0, f(x) \geq 0\}} \right] \right] \\
&= \mathbb{E}_x \left[-\eta(x) \mathbb{1}_{\{g(x) < 0\}} + (1-\eta(x)) (1 - \mathbb{1}_{\{g(x) < 0\}}) \right] \\
&= \mathbb{E}_x \left[2(\eta(x) - \frac{1}{2}) \mathbb{1}_{\{g(x) < 0\}} + (1-\eta(x)) \right] \\
&= \mathbb{E}_x [2f^*(x) \mathbb{1}_{\{g(x) < 0\}} + (1-\eta(x))].
\end{aligned}$$

$$\begin{aligned}
 ② R(f) - R^* &= \mathbb{E}_{x \sim D_{\text{fit}}} [2f^*(x) \mathbf{1}_{\{f(x) < 0\}} - 2f^*(x) \mathbf{1}_{\{f^*(x) < 0\}}] \quad \text{by ①} \\
 &= 2\mathbb{E}[f^*(x)(\mathbf{1}_{\{f(x) < 0\}} - \mathbf{1}_{\{f^*(x) < 0\}})] \\
 &\leq 2\mathbb{E}\left[\mathbf{1}_{\{f^*(x) \geq 0\}} f^*(x) (\mathbf{1}_{\{f(x) < 0\}} - \mathbf{1}_{\{f^*(x) < 0\}}) + \mathbf{1}_{\{f^*(x) < 0\}} f^*(x) (\quad \quad \quad)\right] \\
 &= 2\mathbb{E}\left[\mathbf{1}_{\{f^*(x) \geq 0\}} |f^*(x)| (\mathbf{1}_{\{f(x) < 0\}} - \mathbf{1}_{\{f^*(x) < 0\}}) + \mathbf{1}_{\{f^*(x) < 0\}} |f^*(x)| (\mathbf{1}_{\{f^*(x) < 0\}} - \mathbf{1}_{\{f(x) < 0\}})\right] \\
 &= 2\mathbb{E}\left[\mathbf{1}_{\{f^*(x) \geq 0\}} |f^*(x)| \mathbf{1}_{\{f(x) < 0\}} + \mathbf{1}_{\{f^*(x) < 0\}} |f^*(x)| \mathbf{1}_{\{f^*(x) < 0\}}\right] \\
 &= 2\mathbb{E}\left[\mathbf{1}_{\{f^*(x) \geq 0\}} |f^*(x)| \mathbf{1}_{\{f(x) < 0\}} + \mathbf{1}_{\{f^*(x) < 0\}} |f^*(x)| \mathbf{1}_{\{f(x) > f^*(x) < 0\}}\right] \\
 &\quad - 2\mathbb{E}[|f^*(x)| \mathbf{1}_{\{f(x) > f^*(x) < 0\}}].
 \end{aligned}$$