

today: consistency for convex surrogate losses

non-parametric viewpoint on scores

$$\begin{aligned} \hat{\ell}(x, y; w) &= \langle w, \varphi(x, y) \rangle & \text{if } w &= \sum_{i, \tilde{y}} \alpha_i(\tilde{y}) \varphi(x_i, \tilde{y}) \\ &\Rightarrow \langle w, \varphi(x, y) \rangle &= \sum_{i, \tilde{y}} \alpha_i(\tilde{y}) \underbrace{\langle \varphi(x_i, \tilde{y}), \varphi(x, y) \rangle}_{k(x_i, x; \tilde{y}, y)} \end{aligned}$$

often for simplicity: $k(x, x'; y, y') = k_x(x, x') k_y(y, y')$
 [is equivalent to have $\varphi(x, y) \triangleq \varphi_x(x) \otimes \varphi_y(y)$] "product kernel"
↑
 Kronecker product

$$\begin{aligned} V \otimes w \quad v w^T \\ \langle v \otimes w, v' \otimes w' \rangle &= \text{tr} \left((v w^T)^T (v' w'^T) \right) \\ &= \text{tr} \left(\underbrace{w w'^T}_{\langle w, w' \rangle} \underbrace{v^T v'}_{\langle v, v' \rangle} \right) \\ &= \langle v, v' \rangle \text{tr} \left(\underbrace{w w'^T}_{\langle w, w' \rangle} \right) = \langle v, v' \rangle \langle w, w' \rangle // \end{aligned}$$

e.g. $k_x(x, x') = \exp(-\frac{\|x-x'\|}{2\sigma^2})$ RBF kernel (universal)

$\varphi_y: \mathcal{Y} \rightarrow \mathbb{R}^d$ d ≪ |Y| ≜ k ~ 10²³ $k_y(y, y') = \langle \varphi_y(y), \varphi_y(y') \rangle$

$$S(\tilde{y}) = \langle w_y, \varphi(\tilde{y}) \rangle$$

back to consistency & surrogate losses

$$\hat{w}_n = \underset{w}{\text{argmin}} \hat{\Delta}_n(w) + \lambda_n \frac{\|w\|^2}{2}$$

consistency: $L(\hat{w}_n) \xrightarrow{n \rightarrow \infty} \min_w L(w)$

⊛ binary classification [Bartlett & al. 2004] characterized a whole family of consistent (convex) surrogate losses

↳ binary SVM
 logistic regression

for multiclass classification [Lee & d. 2004] showed that multiclass hinge loss

$$\text{Hinge}(x, y; w) = \max_{\tilde{y}} s(\tilde{y}) + \ell(y, \tilde{y}) - s(y)$$

is not consistent for 0-1 loss when have no "majority" class

$$\text{(i.e. } p(\tilde{y} | x) < \frac{1}{2} \forall \tilde{y})$$

↳ true p

they propose a different surrogate loss that we $\sum_{\tilde{y}}$ instead of $\max_{\tilde{y}}$ which is consistent for 0-1 loss

exponential sum
→ could be intractable for structured prediction

2 aspects of structured prediction which gives a richer theory than binary class. for consistency

- 1) "noise model" $p(y|x)$ is much richer
- 2) $\ell(y, y')$ much richer

⊛ [Ostkin & d. 2017] → we looked at effect of $\ell(y, y')$

for a easy to analyze convex surrogate loss & consistent in the simplest possible setting

and we were careful about exponential constants (eg. $|Y| \leq k$)

Calibration function for a structured loss ℓ , surrogate loss L and a set W

$$H_{L, \ell, W}(\epsilon) \triangleq \inf_{w \in W} [Lq(w) - \min_{w' \in W} Lq(w')] \quad q \in \Delta_{|Y|}$$

$$\text{st. } Lq(w) - \min_{w' \in W} Lq(w') \geq \epsilon$$

"for ϵ -bad w 's"

⊛ x is fixed outside
 q is a potential $p(\tilde{y}|x)$

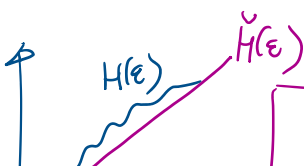
$$Lq(w) \triangleq \mathbb{E}_{q(y)} [L(x, \tilde{y}; w)]$$

$$Lq(w) \triangleq \mathbb{E}_{q(y)} [\ell(\tilde{y}, hw|x)]$$

"conditional (Vapnik) risk"
[conditional on x version]

↳ smallest (over all dist. q) "surrogate optimization regret"
st. true regret $\geq \epsilon$

ie. $\forall q: Lq(w) < Lq^* + H(\epsilon)$
 $\Rightarrow Lq(w) \leq Lq^* + \epsilon$



(thm. 2) $\forall p: L(w) < L^* + \tilde{H}(\epsilon)$



$$\text{(thm. 2)} \quad \forall p: \underbrace{\mathcal{J}(w)}_{\mathbb{E}_{(x,y) \sim p} [\mathcal{L}(x,y;w)]} < \mathcal{J}^* + \check{H}(\epsilon) \Rightarrow L(w) \leq L^* + \epsilon$$

$\check{H}(\epsilon) \triangleq$ convex lower envelope of $H(\epsilon)$

(basically shown using Jensen's inequality)

$$\check{H}(\epsilon) \triangleq H^{**}(\epsilon)$$

$$\mathcal{J}^*(z) \triangleq \sup_x x^T z - f(x) \quad \text{'Fenchel-Legendre conjugate'}$$

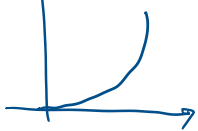
If \check{H} is invertible

$$L(w) - L^* \leq \check{H}^{-1}(\mathcal{J}(w) - \mathcal{J}^*)$$

\mathcal{J} is consistent

iff $H(\epsilon) > 0 \quad \forall \epsilon > 0$
(and $H(\epsilon)$ is finite for some $\epsilon > 0$)

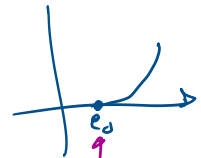
long standard H :



$$H(\epsilon) = \frac{\epsilon^2}{C} \quad H^{-1}(z) = \sqrt{Cz}$$

$$\Rightarrow L(w) - L^* \leq \sqrt{C(\mathcal{J}(w) - \mathcal{J}^*)}$$

You want small C ; for structured pred.
 $C \approx |S|$ often?
(bad)



note: scale of H is arbitrary (scale of \mathcal{L} is arbitrary)

normalize it using stochastic optimization perspective (e.g. SGD)
[next class]

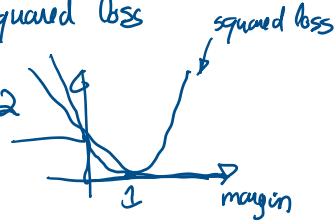
concrete example: simplest surrogate loss & square loss?

$$s(\cdot) \in \mathbb{R}^k \quad (\text{fix } x)$$

$$\mathcal{L}(x, y; s) \triangleq \frac{1}{2k} \|s - (-\log(y; \cdot))\|_2^2 = \frac{1}{2k} \sum_{\tilde{y}} (s(x, \tilde{y}) + \ell(y, \tilde{y}))^2$$

[can be seen as a generalization of squared loss for binary class. to multiclass]

$$[1 - y_i + s_i + \ell(x; \cdot)]^2$$



$$y_i^2 (1 - y_i + s_i)^2$$

$$(y_i - y_i^2 + s_i)^2$$

does not depend on s

$$\mathcal{L}_q(s) \triangleq \mathbb{E}_{q(y)} \mathcal{L}(x, y; s) = \frac{1}{2k} \sum_{\tilde{y}} \mathbb{E}_{q(y)} [s(\tilde{y})^2 + 2s(\tilde{y}) \ell(y, \tilde{y})] + \text{cst.}$$

$$= \frac{1}{2k} \|s + \ell_{q_x}\|_2^2 + \text{cst.}$$

$$\ell_{q_x}(\tilde{y}) \triangleq \mathbb{E}_{q(y)} \ell(y, \tilde{y})$$

Suppose s is unconstrained $\min_s \mathcal{L}_Q(s) \Rightarrow s^*(\tilde{y}) = -l_{Q_x}(\tilde{y})$

$$\arg\max_{\tilde{y}} s^*(\tilde{y}) = \arg\min_{\tilde{y}} l_{Q_x}(\tilde{y})$$

ie. you predict optimally pairwise on x

so here \mathcal{L} is consistent ie. $s^* \in \arg\min_{s \in \mathbb{R}^k} \mathcal{L}(s)$

$$\Rightarrow L(h_{s^*}) = \min_{\text{all } h} L(h)$$

$$\mathcal{L}_Q(s) - \underbrace{\min_{s' \in \mathbb{R}^k} \mathcal{L}_Q(s')}_{\mathcal{L}_Q^*} = \frac{1}{2k} \|s - (-l_{Q_x})\|_2^2$$

Let L be a $k \times k$ matrix where $L_{\tilde{y}, \tilde{y}} \triangleq \ell(\tilde{y}, \tilde{y})$ $l_{Q_x}(\tilde{y}) = \sum_y q(y|x) \ell(y, \tilde{y})$

$$l_{Q_x} = L^{\leftrightarrow} q_x$$

recall: $s^* = -l_{Q_x} = -L^{\leftrightarrow} q_x \in \text{span}(L^{\leftrightarrow})$ ie. $\sum \alpha_y L^{\leftrightarrow}(y, \cdot)$

⊛ to get consistency for Q , it is sufficient to consider $s \in \text{span}(L^{\leftrightarrow})$

or that $s \in \text{span}(F) \supseteq \text{span}(L^{\leftrightarrow})$

restriction on scores

$F \in \mathbb{R}^{k \times r}$ matrix
can be chosen arbitrarily depending on L^{\leftrightarrow}
(row of F as $\ell(y^i) \in \mathbb{R}^r$)

$$s = F\theta \quad \theta \in \mathbb{R}^r \quad (\text{ie. } s(\tilde{y}) = \langle \ell(\tilde{y}), \theta \rangle)$$

$$\mathcal{L}_Q(\theta) - \min_{\theta \in \mathbb{R}^r} \mathcal{L}_Q(\theta) = \frac{1}{2k} \|F\theta - (-L^{\leftrightarrow} q)\|_2^2$$

thm. 7 if $\text{span}(F) \supseteq \text{span}(L^{\leftrightarrow})$

$$H_{\text{Squad}, Q, F}(\epsilon) \geq \frac{\epsilon^2}{2k \max_{i \neq j} \|\Delta_{ij}\|_2^2} \geq \frac{\epsilon^2}{4k}$$

lower bound \Rightarrow crisp result

$$\Delta_{ij} \triangleq e_i - e_j \in \mathbb{R}^k$$

P_F is orthogonal projection on $\text{span}(F)$ $P_F = F(F^T F)^{\dagger} F^T$

• in paper, we show that for 0-1 loss, $H(\epsilon) = \frac{\epsilon^2}{4k}$

thm. 8: if $\text{span}(F) = \mathbb{R}^k$ (ie. no constraints) hardness result

thm 8: if $\text{span}(F) = \mathbb{R}^k$ (ie. no constraints)
 then $H(\epsilon) \leq \frac{\epsilon^2}{4k}$ for any loss? hardness result

⊗ but for Hamming loss, if add constraints that $s(\tilde{y}) = \sum_{p \in \text{points}} S_p(\tilde{y}_p)$

over T binary variables

$$H(\epsilon) = \frac{\epsilon^2}{8T}$$

} not too big

→ we can learn!

note: computation how to compute $\sum_{\tilde{y}} \ell(y, \tilde{y}) s(\tilde{y})$

→ efficient to compute for Hamming loss ϵ (additive)
separable here
function