

today: consistency for convex surrogate losses

non-parametric viewpoint on scores

$$s(x, y; w) = \langle w, \varphi(x, y) \rangle$$

$$\text{if } w = \sum_{i, y} \alpha_i(y) \varphi(x_i, y)$$

$$\Rightarrow \langle w, \varphi(x, y) \rangle = \sum_{i, y} \alpha_i(y) \underbrace{\langle \varphi(x_i, y), \varphi(x, y) \rangle}_{K(x_i, x; y, y)}$$

$$\text{often for simplicity: } K(x, x'; y, y') = K_x(x, x') K_y(y, y')$$

$$[\text{is equivalent to have } \varphi(x, y) \triangleq \varphi_x(x) \otimes \varphi_y(y) \quad \text{"product kernel"}]$$

↑
Kronecker product

$$V \otimes w \quad V w^T$$

$$\langle V \otimes w, V \otimes w \rangle = \text{tr}((V w^T)^T (V w^T))$$

$$\begin{aligned} & \text{tr}(w \underbrace{V^T V}_{\langle V, V \rangle} w^T) \\ & = \langle V, V \rangle \underbrace{\text{tr}(w^T w)}_{\langle w, w \rangle} = \langle V, V \rangle \langle w, w \rangle \end{aligned}$$

$$\text{e.g. } K_x(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right) \quad \text{RBF kernel (universal)}$$

$$\varphi_y: \mathcal{Y} \rightarrow \mathbb{R}^d \quad d \ll |\mathcal{Y}| \triangleq k \quad K_y(y, y') = \langle \varphi_y(y), \varphi_y(y') \rangle$$

↑
10²³

back to consistency & surrogate losses

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

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

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

↳ binary SVM
logistic regression

for multiclass classification [Lee & al. 2004] showed that multiclass hinge loss
 McAllister 2007

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

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

(i.e. $p(y|x) < \frac{1}{2} \forall y$)
 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 tractable for structured prediction

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

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

* [Ozokun & al. 2017] → we looked at effect of $l(y, y')$

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

and we were careful about exponential constants (e.g. $(1/5)^k$)

15h18

calibration function for a structured loss l , surrogate loss \mathcal{L} and set \mathcal{W}

$$H_{\mathcal{L}, l, \mathcal{W}}(\epsilon) \triangleq \inf_{\substack{w \in \mathcal{W} \\ q \in \Delta_{|Y|}}} [\mathcal{L}_q(w) - \min_{w' \in \mathcal{W}} \mathcal{L}_q(w')]$$

x is fixed outside
 q is a potential $p(y|x)$

s.t. $\mathcal{L}_q(w) - \min_{w' \in \mathcal{W}} \mathcal{L}_q(w') \geq \epsilon$

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

$$\mathcal{L}_q(w) \triangleq \mathbb{E}_{q(\tilde{y})} [l(\tilde{y}, h_w(x))]$$

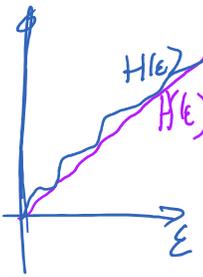
"conditional (logit) risk"
 [conditional on x version]

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

i.e. $\forall q: \mathcal{L}_q(w) < \mathcal{L}_q^* + H(\epsilon)$
 $\Rightarrow \mathcal{L}_q(w) \leq \mathcal{L}_q^* + \epsilon$

L convex on x reason

$$i.e. \forall \theta \in \partial q(w) \cdot \forall q \in \Pi(\epsilon) \Rightarrow L_q(w) \leq L^* + \epsilon$$



(thm. 2) $\forall \rho : \beta(w) < \beta^* + H(\epsilon)$

$$\mathbb{E}_{(x,y) \sim \rho} L(\beta(x,y;w))$$

$$\Rightarrow L(w) \leq L^* + \epsilon$$

convex lower envelope of H(epsilon)

(↑ basically shown using Jensen's ineq.)

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

$$f^*(z) \triangleq \sup_x x^T z - f(x) \Leftrightarrow \text{Fenchel-legendre conjugate}$$

If \check{H} is invertible

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

\mathcal{L} is consistent iff $H(\epsilon) > 0 \forall \epsilon > 0$

(and $H(\epsilon)$ is simple for some $\epsilon > 0$)

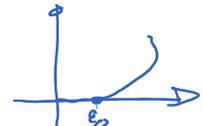
standard H is



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

$$L(w) - L^* \leq \sqrt{C(L(w) - \beta^*)}$$

you want small C; for structured prediction $C = |\mathcal{S}|$ often (bad)



min regret I can guarantee

note: scale of H is arbitrary

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

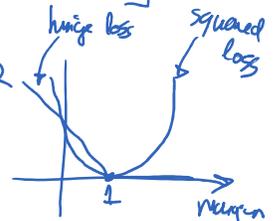
concrete example: simplest surrogate loss: square loss?

$$S(\cdot) \in \mathbb{R}^k \quad (S(x, z))$$

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

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

$$[1 - y_i \langle w, \phi(x; \tilde{y}) \rangle]^2$$



$$\mathcal{L}_q(s) \triangleq \mathbb{E}_{q(y)} \mathcal{L}(x, y; s)$$

$$= \frac{1}{2k} \sum_{\tilde{y}} \mathbb{E}_{q(y)} [s(\tilde{y})^2 + 2s(\tilde{y}) \ell(y, \tilde{y}) + \ell(y, \tilde{y})^2]$$

does not depend on s

$$\begin{aligned} E_{q|y} l(y, \tilde{y}) &\triangleq l_{q_x}(\tilde{y}) \\ &= \frac{1}{2k} \|s + l_{q_x}\|_2^2 + \text{cst.} \end{aligned}$$

Suppose s is unconstrained $\min_s l_{q_x}(s) \Rightarrow s^*(\tilde{y}) = -l_{q_x}(\tilde{y})$
 $\arg \max_{\tilde{y}} s^*(\tilde{y}) = \arg \min_{\tilde{y}} l_{q_x}(\tilde{y})$
 i.e. you predict optimally pointwise on x

so here l is consistent i.e. $s^* \in \arg \min_{s: X \rightarrow \mathbb{R}^k} l(s)$

$$l_q(s) - \min_{s' \in \mathbb{R}^k} l_q(s') = \frac{1}{2k} \|s - (-l_{q_x})\|_2^2 \Rightarrow L(h_{s^*}) = \min_{\text{all } h} L(h)$$

let \tilde{L} be a $k \times k$ matrix where $\tilde{L}_{\tilde{y}, y} = l(y, \tilde{y})$ $l_{q_x} = \sum_y q(y|x) l(y, \cdot)$

$$l_{q_x} = \tilde{L} q_x$$

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

⊛ to get consistency for l , it is sufficient to consider $s \in \text{span}(\tilde{L})$
 or that $s \in \text{span}(F) \supseteq \text{span}(\tilde{L})$
restriction on errors

$F \in \mathbb{R}^{k \times r}$ matrix
 can be chosen clearly depending on \tilde{L}

$$S = F \Theta \quad \Theta \in \mathbb{R}^r$$

$$l_q(\Theta) - \min_{\Theta \in \mathbb{R}^r} l_q(\Theta) = \frac{1}{2k} \|F\Theta - (-\tilde{L}q)\|_2^2$$

Thm. 7

if $\text{span}(F) \supseteq \text{span}(\tilde{L})$

$H_{l_{\text{square}}, l, F}(\tilde{\epsilon}) \geq$

$$\frac{\epsilon^2}{2k \max_{i \neq j} \|P_F \Delta_{ij}\|_2^2} \geq \frac{\epsilon^2}{4k}$$

lower bound \Rightarrow robustness result

this is bad

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

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}
 then $H(\epsilon) \leq \frac{\epsilon^2}{4k}$ for any loss?

ie. for any loss, we need an exp # of samples (in worst case)
 to learn well [concat \rightarrow all these bounds] one worst case

⊛ but for Hamming loss, if add constraints that $s(\vec{y}) = \sum_{p \text{ parts}} s_p(\vec{y}_p)$

over T binary variables, $H(\epsilon) = \frac{\epsilon^2}{8T}$ } not too big
 \rightarrow we can learn?

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

\rightarrow efficient to compute for
 Hamming loss & separable score fct. eg