

today: consistency for convex surrogate losses

non-parametric viewpoint on scores

$$s(x, y; \omega) = \langle \omega, \ell(x, y) \rangle$$

$$\text{if } \omega = \sum_{i, \tilde{y}} \alpha_i(\tilde{y}) \ell(x_i, \tilde{y})$$

$$\Rightarrow \langle \omega, \ell(x, y) \rangle = \sum_{i, \tilde{y}} \alpha_i(\tilde{y}) \underbrace{\langle \ell(x_i, \tilde{y}), \ell(x, y) \rangle}_{k(x_i, x; \tilde{y}, y)}$$

$$\text{often for simplicity: } k(x, x'; y, y') = k_x(x, x') k_y(y, y')$$

$$\left[\text{is equivalent to have } \ell(x, y) = \varphi_x(x) \otimes \varphi_y(y) \right]$$

↑
Kronecker product

"product kernel"

$$V \otimes \omega \rightarrow V\omega^T$$

$$\begin{aligned} \underbrace{\langle V \otimes \omega, V \otimes \omega' \rangle}_{V\omega^T} &= \text{tr}((V\omega^T)^T (V\omega'^T)) \\ &= \text{tr}(V \overbrace{\omega^T V^T}^{\langle V, V' \rangle} \omega'^T) \\ &= \langle V, V' \rangle \text{tr}(\omega^T \overbrace{\omega'}^{\langle \omega, \omega' \rangle}) = \langle V, V' \rangle \langle \omega, \omega' \rangle // \end{aligned}$$

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

$$\varphi_y: \mathcal{T} \rightarrow \mathbb{R}^d \quad d \ll |\mathcal{T}| \triangleq k \quad k_y(y, y') = \langle \varphi_y(y), \varphi_y(y') \rangle$$

back to consistency of surrogate losses

$$\hat{\omega}_n \stackrel{\Delta}{=} \underset{\omega}{\text{arg min}} \quad \hat{f}_n(\omega) + \lambda_n \frac{\|\omega\|^2}{2}$$

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

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

↳ binary SVM

Logistic regression

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

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

is not consistent for 0-1 loss
when have no "majority" loss (i.e. $p(\tilde{y}|x) < \frac{1}{2} \forall y$)

they propose a different surrogate loss that uses

$\sum_{\tilde{y}} s(\tilde{y})$ instead of $\max_{\tilde{y}} s(\tilde{y})$

which is consistent for 0-1 loss

exponential sum

→ could be intractable 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) $l(y, \tilde{y})$ much richer

⊗ [Osokin & al. 2017] → we looked at effect of $l(y, \tilde{y})$

for a easy to analyze convex surrogate loss $\{$ consistent in the simplest possible setting

and we were careful about exponential constants (e.g. $|w|=k$)

calibration function for a structured loss l , surrogate q and set W

$$H_{q,l,W}(\varepsilon) \triangleq \inf_{\substack{w \in W \\ q \in \Delta_{|W|}}} [L_q(w) - \min_{w' \in W} L_q(w')] \quad \text{s.t. } L_q(w) - \min_{w' \in W} L_q(w') \geq \varepsilon$$

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

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

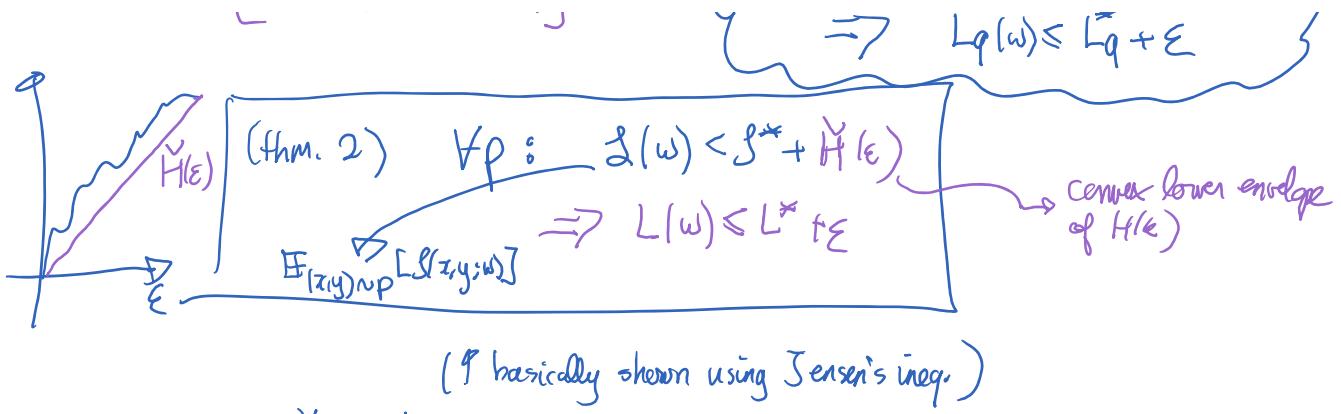
$$L_q(w) \triangleq \mathbb{E}_{q(y)} [l(y, h_w(x))]$$

"conditional (Vapnik) risk"

[conditional on x version]

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

e.g. $H_q: L_q(w) \leq d_q + H(\varepsilon)$
 $\Rightarrow L_q(w) \leq d_q + \varepsilon$

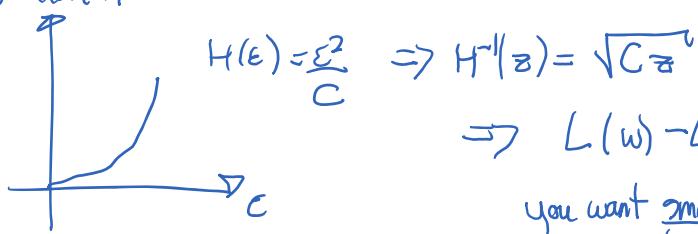


$$H(\epsilon) \triangleq H^{**}(\epsilon) \quad f^*(z) \triangleq \sup_z z^T z - f(z) \Leftrightarrow \text{"Fenchel-Legendre conjugate"}$$

If H is invertible

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

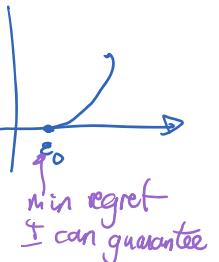
standard H :



you want small C; for structural prediction
 $C = |\omega|$ often?
 (bad)

L is consistent

(if $H(\epsilon) > 0 \forall \epsilon > 0$
 (and $H(\epsilon)$ is finite for some $\epsilon \geq 0$)



note: scale of H is arbitrary

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

14h34

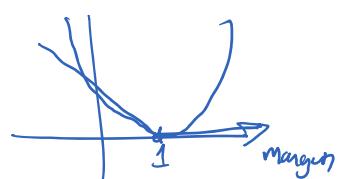
Concrete example: simplest surrogate loss: square loss?

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

$$\delta(x_i, y_i; s) \triangleq \frac{1}{2k} \|s - (-l(y_i; \cdot))\|_2^2 = \frac{1}{2k} \sum_{j=1}^k (s(x_i) + l(y_i, j))^2$$

L can be seen as generalization of squared loss
 for binary class. to multiclass

$$\begin{aligned} y_i = 1 & [1 - y_i \langle w, \ell(x_i) \rangle]^2 \\ & = [y_i - \langle w, \ell(x_i) \rangle]^2 \end{aligned}$$



$$\delta_q(s) \triangleq \mathbb{E}_{y \sim p} \delta(x_i, y; s)$$

$$= \frac{1}{n} \sum_{i=1}^n [s(\hat{y})^2 + 2s(\hat{y})l(y_i, \hat{y})] + \text{const.}$$

does not depend on s

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

$$= \frac{1}{2k} \|s + l_{q_x}\|_2^2 + \text{cst.} \quad l_{q_x}(\tilde{y}) \triangleq \mathbb{E}_{q(y)} [l(y, \tilde{y})]$$

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

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

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

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

$$\text{Let } \overset{\leftrightarrow}{L} \text{ be a } k \times k \text{ matrix where } L_{\tilde{y}, y} = l(y, \tilde{y}) \quad l_{q_x}(\cdot) = \sum_{\tilde{y}} q(\tilde{y}|x) l(\tilde{y}, \cdot)$$

$$l_{q_x} = \overset{\leftrightarrow}{L} q_x$$

$$\text{recall: } s^* = -l_{q_x} = -\overset{\leftrightarrow}{L} q_x \in \text{span}(\overset{\leftrightarrow}{L}) \text{ i.e. } \sum_y q(y|x) \overset{\leftrightarrow}{L}(y, \cdot)$$

④ to get consistency for \mathcal{L} , it's sufficient to consider $s \in \text{span}(\overset{\leftrightarrow}{L})$
 or that $s \in \text{span}(F) \supseteq \text{span}(\overset{\leftrightarrow}{L})$
restriction on scores

$F \in \mathbb{R}^{K \times r}$ matrix
 can be chosen cleverly depending on $\overset{\leftrightarrow}{L}$

$$s = F\theta \quad \theta \in \mathbb{R}^r$$

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

thm. 7 if $\text{span}(F) \supseteq \text{span}(\overset{\leftrightarrow}{L})$ $H_{\text{Squared}, \ell, F}(\epsilon) \geq \frac{\epsilon^2}{2k \max_{i \in \mathcal{Y}} \|F \Delta_{ij}\|_2^2} \geq \frac{\epsilon^2}{4K}$ this is bad

lower bound \Rightarrow excess result

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

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

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

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

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

i.e. for any loss, we need an exp. # of samples (in worst case)
to learn well!

[correct \rightarrow all these bounds are
worst case.]

\rightarrow without structure on F , you're doomed in worst case over $p(x,y)$

② but for Hamming Loss, if add constraints that $s(y') = \sum_{y \neq y'} s_p(y)$

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

note : computation how to compute $\sum_y s(y, \hat{y}) s(\hat{y})$

(polynomial sample complexity)

\rightarrow efficient to compute for

Hamming loss $\{$ separable score fct. e.g.