

today: consistency for convex surrogate losses

non-parametric viewpoint on scores

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

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

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

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

[is equivalent to having $\phi(x, y) \triangleq \phi_x(x) \otimes \phi_y(y)$]
 "product kernel"
 ↑
 Kronecker product

$$v \otimes w \quad v w^T$$

$$\begin{aligned} \langle v \otimes w, v' \otimes w' \rangle &= \text{tr}((v w^T)^T (v' w'^T)) \\ &= \text{tr}(w v'^T v w'^T) \\ &= \langle w, w' \rangle \langle v, v' \rangle // \end{aligned}$$

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

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

Back to consistency & surrogate losses

$$\hat{w}_n \triangleq \underset{w}{\text{argmin}} \hat{L}_n(w) + \frac{\lambda_n \|w\|_2^2}{2}$$

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

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

↳ binary SVM
logistic regression

for multiclass classification [Lee & d. 2004] showed that multiclass hinge loss
McAlister 2007
is not consistent for 0-1 loss when we have no "majority" class
(i.e. $p(y|x) < \frac{1}{2} \forall y$)

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

they propose a different surrogate loss that uses $\sum_{\tilde{y}} \ell(y, \tilde{y})$ instead of $\max_{\tilde{y}}$
which is consistent for 0-1 loss
exponential sum
→ could be intractable for structured prediction

2 aspects of structured pred. which give a much richer theory than binary classifi. for Consistency

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

⊛ [Osokin & 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 (e.g. 1/51)

Calibration function for a structured loss ℓ , surrogate loss L and set W (x is fixed outside and q is a potential $p(y|x)$)

$$H_{L, \ell, W}(\epsilon) \triangleq \inf_{\substack{w \in W \\ q \in \Delta_{\mathcal{Y}}}} [Lq(w) - \min_{w' \in W} Lq(w')] \quad \text{s.t. } Lq(w) - \min_{w' \in W} Lq(w') \geq \epsilon$$

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

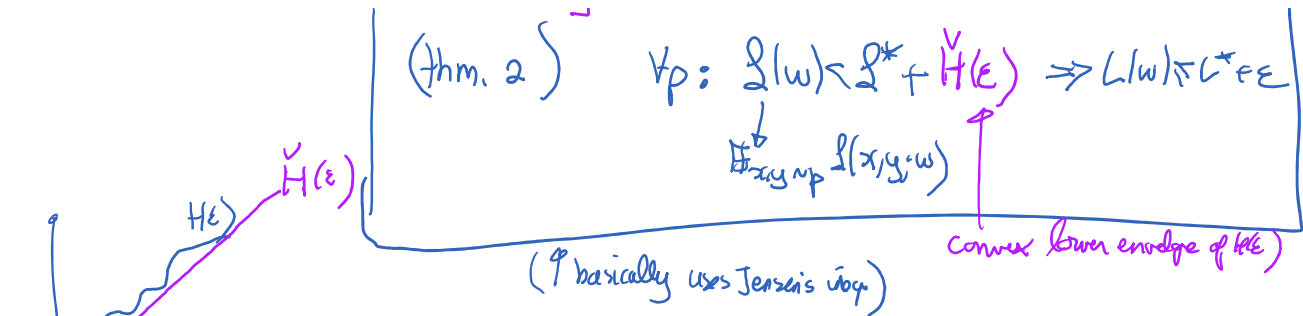
$$Lq(w) \triangleq \mathbb{E}_{q(y)} [L(\tilde{y}, hw(x))] \quad \text{"conditional risk"}$$

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

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

[condition on x version]

(thm. 2) $\forall p: L(w) < L^* + H(\epsilon) \Rightarrow L(w) \leq L^* + \epsilon$



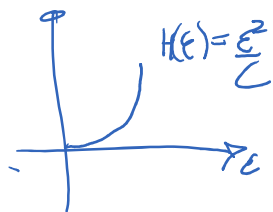
$$\tilde{H}(\epsilon) \triangleq H^{**}(\epsilon) \quad f^*(z) \triangleq \sup_x x^T z - f(x) \Leftrightarrow \text{"Fenchel-Legendre conjugate"}$$

if \tilde{H} is invertible

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

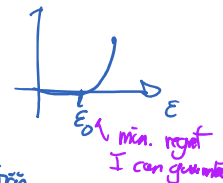
\mathcal{L} is constant iff $H(\epsilon) > 0 \forall \epsilon > 0$ (and $H(\epsilon)$ is finite) for some $\epsilon > 0$

standard H :



$$H^{-1}(z) = \sqrt{Cz} \Rightarrow \mathcal{L}(w) - \mathcal{L}^* \leq \sqrt{C(\mathcal{L}(w) - \mathcal{L}^*)}$$

you want small C ; for structured prediction $C = |S|$ often?



14h48

note: scale of H is arbitrary;

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

simplest surrogate loss: square loss?

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

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

[can be seen as a generalization of squared loss for binary class. to multiclass $(1 - y_i < w_i, \ell(x_i) > 2)$

$$\begin{aligned} \mathcal{L}_q(s) &\triangleq \mathbb{E}_{q(y)} \mathcal{L}(x, y; s) \\ &= \frac{1}{2k} \sum_{\tilde{y}} \mathbb{E}_{q(y)} [s(\tilde{y})^2 + 2s(\tilde{y})\ell(y, \tilde{y}) + \text{const.}] \quad \text{does not depend on } s \\ &\quad \mathbb{E}_{q(y)} \ell(y, \tilde{y}) \triangleq q_x(\tilde{y}) \\ &= \frac{1}{2k} \|s + q_x\|^2 + \text{const.} \end{aligned}$$

suppose s is unconstrained $\min_s \mathcal{L}_q(s) \Rightarrow s^*(\tilde{y}) = -q_x(\tilde{y})$

$$\arg \max_y s^*(y) = \arg \min_y Lq_x(y)$$

i.e. you predict optimally
pointwise

so here L is consistent

$$\text{i.e. } s^* \in \arg \min_{\text{all } s: \mathbb{R}^k \rightarrow \mathbb{R}^r} f(s)$$

$$\Rightarrow L(h_s) = \min_{\text{all } h} L(h)$$

$$Lq(s) - \min_{s' \in \mathbb{R}^k} Lq(s') = \frac{1}{2K} \|s - (-Lq)\|_2^2$$

Let \overleftrightarrow{L} be a $k \times k$ matrix where $\overleftrightarrow{L} \tilde{y} = L(y, \tilde{y})$ $(Lq_x = \sum_y q(y|x) L(y, \cdot))$

$$Lq_x = \overleftrightarrow{L} q_x$$

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

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

or that $\text{span}(F) \supseteq \text{span}(\overleftrightarrow{L})$

restriction on scores

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

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

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

$$Lq(\theta) - \min_{\theta \in \mathbb{R}^r} Lq(\theta) = \frac{1}{2K} \|F\theta - \overleftrightarrow{L}q\|_2^2$$

Thm. 7
if $\text{span}(F) \supseteq \text{span}(\overleftrightarrow{L})$

$$H_{\text{square}, L, F}(\epsilon) \stackrel{\text{lower bound} \Rightarrow \text{easier result}}{\geq} \frac{\epsilon^2}{2K \max_{i \neq j} \|F \Delta_{ij}\|_2^2} \geq \frac{\epsilon^2}{4K} \quad \text{this is bad}$$

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

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

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

Thm. 8: if $\text{rank}(F) = \mathbb{R}^k$ (i.e. no constraints) hardness result

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

i.e. for any loss, we need an exp # of samples (in the worst case) to learn "well" [caution \rightarrow all these are bounds] and worst case

⊗ but for Hamming loss, if add constraints that $f(\vec{y}) = \sum_{j \in J} \mathbb{1}(y_j)$

over J binary variables, $H(\epsilon) = \frac{\epsilon^2}{8|J|}$ not too big \rightarrow we can recurse

note: computation how to compute $\sum_{\vec{y}} \ell(y_j) s(y_j)$
 \rightarrow efficient to compute for Hamming loss & separable score e.g.