

Lecture 8 - scribbles - Calibrated convex surrogate losses

Thursday, January 31, 2019 13:06

today: consistency for convex surrogate losses

Non-parametric viewpoint on scores

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

$$k(x_i, \cdot; \tilde{y}, \cdot)$$

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

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

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

"product kernel"

[is equivalent to having $\varphi(x, y) = \varphi_x(x) \otimes \varphi_y(y)$]

\uparrow
Kronecker product

$$v \otimes w \quad v w^\top$$

$$\begin{aligned} \langle v \otimes w, v' \otimes w' \rangle &= \text{tr}((v w^\top)(v' w'^\top)) \\ &= \text{tr}(w v^\top v' w'^\top) \\ &= \langle w, w' \rangle \langle v, v' \rangle \end{aligned}$$

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

RBF kernel (universal)

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

Back to consistency { surrogate losses}

$$\hat{w}_n \triangleq \operatorname{arg\,min}_{w} \frac{\hat{J}_n(w) + \lambda_n \|w\|^2}{2}$$

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

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

↳ binary SVM
logistic regression

for multiclass classification [Lee et al. 2004]
McAllester 2002 showed that multiclass hinge loss

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

is not consistent for 0-1 loss
when have no majority class (i.e. $\Pr[y|x] < \frac{1}{2}$)

They propose a different surrogate loss that $\sum_{\tilde{y}}$ instead of $\max_{\tilde{y}}$
which is consistent for 0-1 loss
exponential sum
→ could be intractable ↴

for 0-1 loss

exponential sum
→ could be intractable

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

1) "noise model" $p(y|x)$ is much richer

2) $\ell(y, y')$ much richer

④ [Osokin et al. 2017] → we looked at effect of $\ell(y, y')$

for a easy to analyze convex surrogate loss

in the simplest possible setting

and we were careful about exponential constants

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

$$H_{L,\ell,W}(\varepsilon) \triangleq \inf_{\substack{w \in W \\ q \in \Delta^{|\mathcal{Y}|}}} [L_q(w) - \min_{w' \in W} L_q(w')] \quad \text{s.t. } L_q(w) - \min_{w' \in W} L_q(w') \geq \varepsilon$$

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

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

$$L_q(w) \triangleq \mathbb{E}_{q(\tilde{y})} [\ell(\tilde{y}, h_w(x))] \quad \text{"conditional risk"}$$

(conditional on x
version)

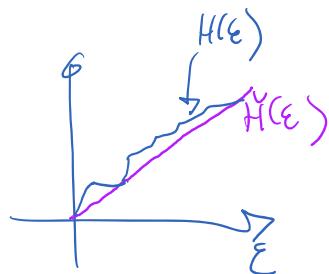
→ smallest optimization surrogate regret (over all dist. q)
s.t. true regret $\geq \varepsilon$

$$\forall q: \quad \mathcal{S}_q(w) < \mathcal{S}_q^* + H(\varepsilon) \Rightarrow L_q(w) \leq L^* + \varepsilon$$

(thm. 2) $\forall p: \quad \mathcal{S}(w) < \mathcal{S}^* + \tilde{H}(\varepsilon) \Rightarrow L(w) \leq L^* + \varepsilon$

$\downarrow_{\mathbb{F}_{x,y} \mathcal{S}(x,y; w)}$

(convex lower envelope of $H(\varepsilon)$)



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

L is consistent iff $H(\varepsilon) > 0 \quad \forall \varepsilon > 0$
(and $H(\varepsilon)$ is finite for some $\varepsilon > 0$)

if H is invertible

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

$$H(\varepsilon) = \frac{\varepsilon^2}{C} \Rightarrow L(w) - L^* \leq \sqrt{C(\mathcal{S}(w) - \mathcal{S}^*)}$$

you want small C ; for structured prediction $C = |\mathcal{Y}|$ often?

Mhw24

note: scale of H is arbitrary;
normalize it using optimization perspective (e.g. SGD)

- [see soon] " "

Simplest surrogate loss: square loss ▷

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

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

[can be seen as generalization of squared loss
for binary classification to multi-class
 $(1 - g_i(w_i^\top \phi(x_i)))^2$]

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

does not depend on s

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

Suppose s is unconstrained

$$\min_s \mathcal{L}_q(s) \quad s(\tilde{y}) + l_{q_x}(\tilde{y}) = 0 \quad \forall \tilde{y}$$

$$\Rightarrow s^*(\tilde{y}) = -l_{q_x}(\tilde{y})$$

$$\arg \max_{\tilde{y}} s^*(\tilde{y}) = \arg \min_{\tilde{y}} l_{q_x}(\tilde{y}) \quad \text{ie. You predict optimally (pt. wise)}$$

so here \mathcal{L} is consistent

$$\text{i.e. } s^* \in \arg \min_{\text{all } s} \mathcal{L}(s) \Rightarrow L(h^*) = \min_{\text{all } h} L(h)$$

$$h_g^*(z) = \operatorname{argmax}_{\tilde{y}} \hat{L}(z, \tilde{y})$$

$$\mathcal{L}_q(s) = \|s - (-l_q)\|^2 + \underline{\text{const}}$$

$$\mathcal{L}_q(s) - \min_{s' \in \mathbb{R}^k} \mathcal{L}_q(s') = \frac{1}{2k} \|s - (-l_q)\|_2^2$$

$$l_q(x) = \sum_y q(y|x) l(y, \cdot)$$

Let $\overset{\leftrightarrow}{L}$ be a $k \times k$ matrix where $\overset{\leftrightarrow}{L}\tilde{y}, y = l(y, \tilde{y})$

$$l_q(x) = \overset{\leftrightarrow}{L} q_x$$

$$s^* = -l_q(x) = -\overset{\leftrightarrow}{L} q_x \in \operatorname{span}(\overset{\leftrightarrow}{L}) \text{ i.e. } \sum_{\tilde{y}} q_{\tilde{y}} \overset{\leftrightarrow}{L}(\cdot, \tilde{y})$$

To get consistency for l , it is sufficient to consider $s \in \operatorname{span}(\overset{\leftrightarrow}{L})$

or that $s \in \operatorname{span}(F) \supseteq \operatorname{span}(\overset{\leftrightarrow}{L})$
restriction on scores

$F \in \mathbb{R}^{K \times r}$ matrix

can be chosen cleverly depending on $\overset{\leftrightarrow}{L}$

$$s = FG \quad G \in \mathbb{R}^{r \times r}$$

if $\operatorname{span}(F) \supseteq \operatorname{span}(\overset{\leftrightarrow}{L})$

$$\mathcal{L}_q(F) - \min_{G \in \mathbb{R}^{r \times r}} \mathcal{L}_q(G) = \frac{1}{2K} \|FG - (\overset{\leftrightarrow}{L}q)\|_2^2$$

, convex bound \Rightarrow easiness result

.. -

thm. 7

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

$$H_{\ell_2^{\text{square}}, F}(\varepsilon) \geq \frac{\varepsilon^2}{2K \max_{i \neq j} \|P_F \Delta_{ij}\|_2^2} \geq \frac{\varepsilon^2}{4K}$$

Cover bound \Rightarrow easiness result

this bad

$$\Delta_{ij} \triangleq \ell_i - \ell_j \in \mathbb{R}^K$$

P_F is orthogonal projection on $\text{span}(F)$ $P_F = F(F^T F)^{-1} F^T$

- in the paper, we show that for O+I loss, $H(\varepsilon) = \frac{\varepsilon^2}{4K}$

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

$$\text{then } H(\varepsilon) \geq \frac{\varepsilon^2}{2K} \text{ for any loss?}$$

i.e. for any loss, we need an exponential # of samples (in the worst case)
to learn "well" {caveat \rightarrow all those are bounds and worst case}

⊗ but for hamming loss, if add constraint that $s(\tilde{y}) = \sum_{P_k} \text{sp}(\tilde{y})$

$$\text{over } T \text{ binary variables, } H(\varepsilon) = \frac{\varepsilon^2}{8T} \quad \{ \text{not too big} \Rightarrow \text{we can learn!} \}$$

Note: computation now to compute $\sum_{\tilde{y}} l(y, \tilde{y}) S(\tilde{y})$
 → efficient for Hamming Loss & separable score ...

optimization normalization

Setup Set $s(x, \tilde{y})$ be of the form $F\Theta(x) \quad \Theta(x) \in \mathbb{R}^r$

$$\Theta_j(\cdot) \in \mathcal{H} \leftarrow \text{RHS}$$

optimization variables

$$S(G) = \mathbb{E}_{(x,y) \sim p} s(x, y; G)$$

run projected SGD on $S(\Theta)$ i.e. $\Theta^{(t+1)} = P_{\Theta} \{ \Theta^{(t)} - \gamma \nabla_S S(x^{(t)}, y^{(t)}, \Theta^{(t)}) \}$

$$\begin{aligned} & (x^{(t)}, y^{(t)}) \text{ iid } p \\ & \downarrow \\ & \text{a ball of radius } D \\ & \downarrow \\ & F^T \nabla_S S(x^{(t)}, y^{(t)}, \Theta^{(t)}) \Phi^{(t)} \Gamma \\ & \downarrow \text{ gradient map} \\ & K(x^{(t)}, \cdot) \end{aligned}$$

Convergence result: if $\|\Theta^*\|_{HS} \leq D$
 (thm. 5)

$$\text{if } \mathbb{E}_{(x,y) \sim p} \|\nabla_S S(x, y, \Theta)\|_{HS} \leq M^2$$

then averaged SGD with stepsize $\gamma = \frac{2D}{M\sqrt{n}}$

$$\mathbb{E} \left[S(\Theta^{(t)}) \right] - S(\Theta^*) \leq \frac{2DM}{\sqrt{n}} \quad (\text{convergence result})$$

$$\mathbb{E} [L(\hat{\theta}^{(t)})] - J(\theta) \leq \frac{\alpha D M}{\sqrt{n}} \quad (\text{convergence result})$$

$\sum_{n=1}^N \hat{\theta}^{(t)}$

thm. 6 Learning complexity (to be cted) \Rightarrow