

today: consistency & tractable methods

consistency & calibration functions

need to relate  $S(w)$  to  $L(w)$ : look "calibration function" [Steinwart]

relationship is usually very complicated

$\Rightarrow$  current results look mainly at non-parametric setting ( $\infty$  # of parameters)

all functions  $h: X \rightarrow Y$  are considered;  $\rightarrow$  this evacuates the dependence on  $x$  in the analysis

ie. we suppose that  $S(x, y; w)$  can be arbitrary for any  $x$ ;  $\rightarrow$  could use a universal kernel  
(ie.  $w$  is  $\infty$ -dim)  $S(\cdot, \cdot; w) \in \mathcal{H}_{X \times Y}$

kernels:

$$\langle w, \phi(x, y) \rangle = S(x, y; w) \quad \text{if } w = \sum_{i, \tilde{y}} \alpha_i(y) \phi(x_i, \tilde{y})$$

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

often for simplicity,  $k(x, x'; y, y') = K_x(x, x') K_y(y, y')$

"product kernel"

[is equivalent to having  $\phi(x, y) = \phi_x(x) \otimes \phi_y(y) \Leftrightarrow$ ]

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

$\nearrow$   
kronecker product

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

eg:  $\phi: S \rightarrow \mathbb{R}^d$   $d \leq k$   
 $= |S|$

$$K_y(y, y') = \langle \phi(y), \phi(y') \rangle$$

$$K_x(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

RBF kernel  
"universal"

$$= \langle w, w' \rangle \langle v, v' \rangle$$

$$\mathcal{H} = \text{RKHS} = \text{closure}(\text{span}\{k(x, \cdot) : x \in X\})$$

$$h \in \mathcal{H} \text{ is s.t. } \sum_{i=0}^{\infty} \alpha_i v_i$$

↙ basis for  $\mathcal{H}$

"reproducing property"

$$\langle k(x, \cdot), h \rangle_{\mathcal{H}} = h(x) \text{ for } h \in \mathcal{H}$$

$\mathcal{H}$  is  $L_2$ -dense in  $L_2 = \{ \text{space of square integrable functions} \}$

i.e. for any  $f \in L_2$ ,  $\exists$  sequence  $h_n \in \mathcal{H}$  s.t.  $\|h_n - f\|_{L_2} \xrightarrow{n \rightarrow \infty} 0$

$$\|f\|_{L_2}^2 = \int_x f(x)^2 dx$$

$$\langle f, g \rangle = \int_x f(x)g(x) dx$$

$$\begin{aligned} \langle h, h \rangle_{\mathcal{H}} &= \left\langle \sum_i \alpha_i k(x_i, \cdot), \sum_j \alpha_j k(x_j, \cdot) \right\rangle \\ &= \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) \end{aligned}$$

$$L(h) = \mathbb{E}_{x,y} \ell(y, h(x)) = \mathbb{E}_x \left[ \underbrace{\mathbb{E}_{y|x} \ell(y, h(x))}_{q_x(h(x))} \right]$$

↙  $q_x \triangleq p(\cdot|x)$

define  $h^*(x) = \arg \min_{y \in \mathcal{Y}} q_x(y)$

then  $L(h^*) = \min_{\text{all functions } h: X \rightarrow \mathcal{Y}} L(h)$

$$\mathbb{E}_{y|x} \ell(y, y') \geq \mathbb{E}_{y|x} \ell(y, h^*(x)) \quad \forall y'$$

$$\begin{aligned} \mathbb{E}_x \mathbb{E}_{y|x} \ell(y, y'(x)) &\geq \mathbb{E}_x \mathbb{E}_{y|x} \ell(y, h^*(x)) \\ \text{i.e. } L(y'(\cdot)) &\geq L(h^*) \quad \forall y'(\cdot) \end{aligned}$$

[for SBD setup;  $L(w) = \mathbb{E}_{\mathcal{Z}} \tilde{L}(w, \mathcal{Z})$ ]

$$\nabla L(w) = \mathbb{E}_{\mathcal{Z}} \nabla_w \tilde{L}(w, \mathcal{Z})$$

in SBD, you have  $\nabla_w \tilde{L}(w, \mathcal{Z}) \neq \nabla_w L(w)$

$$\mathbb{E}[\nabla_w \tilde{L}(w_t, \mathcal{Z}) | w_t] = \nabla_w L(w_t)$$

$$\mathcal{L}(w) = \mathbb{E}_{(x,y) \sim p} [\mathcal{L}(x, y; w)]$$

$$\nabla \mathcal{L}(w) = \mathbb{E}_{(x,y) \sim p} \nabla_w \mathcal{L}(x,y;w) \rightarrow \text{SbD } (x_i, y_i) \sim p$$

$$\hat{w}_n = \arg \min_w \hat{\mathcal{L}}(w) + \frac{\lambda_n \|w\|^2}{2} \quad w_{t+1} = w_t - \delta_t \nabla_w \mathcal{L}(x_i, y_i; w_t)$$

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

\* binary classification [Bentley & al. 2004] characterized a whole family of consistent surrogate losses

→ binary SVM consistent logistic regression

for multiclass classification [Lee & al. 2004] showed that multiclass hinge loss  $\mathcal{L}_{\text{hinge}}(x,y;w) = \max_{\tilde{y} \neq y} S(x,\tilde{y};w) + \ell(y,\tilde{y})$

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

→ they propose a different surrogate loss that has  $\sum_y$  instead of  $\max_y$  which is consistent 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)  $p(y|x)$  "noise model" is much richer

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

⊗ [Osokin & al. 2017] → we looked at effect of  $\ell(y,y')$  for a convex consistent surrogate loss in the simplest possible setting and were careful about exponential constants

calibration function for structured cost  $\ell$ , surrogate  $\mathcal{L}_q$  and set  $\mathcal{W}$

$$H_{\ell, \ell, \mathcal{W}}(\epsilon) \triangleq \inf_{w \in \mathcal{W}} \left[ \mathcal{L}_q(w) - \min_{w' \in \mathcal{W}} \mathcal{L}_q(w') \right]$$

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

$\mathcal{J}_q(w) \triangleq \mathbb{E}_{q(y)} [\mathcal{J}(x, y; w)]$   
 true conditional risk:  
 $L_q(w) \triangleq \mathbb{E}_{q(y)} [L(y, h(w; x))]$

$q \in \Delta(\mathcal{Y})$  s.t.  $[\mathcal{J}_q(w) - \min_{w' \in \mathcal{W}} \mathcal{J}_q(w')] \geq \epsilon$   
 $\downarrow$   
 $L_q^*$

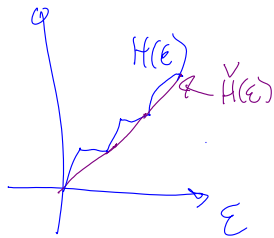
$\hookrightarrow$  smallest optimization surrogate regret (over all dist.  $q$ )  
 s.t. true regret  $\geq \epsilon$

Consequence

(conditional on  $x$  version)

$\forall q: \mathcal{J}_q(w) < \mathcal{J}_q^* + H(\epsilon) \Rightarrow L_q(w) \leq L_q^* + \epsilon$

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



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

$\check{H}(\epsilon) \triangleq H^{**}(\epsilon) \quad f^*(z) \triangleq \sup_x x^T z - f(x) \quad \leftarrow$  Fenchel-Legendre conjugate

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

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

e.g. here  
 $0 \leq L(w) - L^* \leq \sqrt{C(\mathcal{J}(w) - \mathcal{J}^*)}$

Next time: we complete argument  
 for sample complexity using SGD convergence result