

Lecture 8 - scribbles - calibration

Friday, February 16, 2018

14:30

today: finish calibrated surrogate loss story

recall: $J(x, y; s) \triangleq \frac{1}{2k} \sum_{\tilde{y}} (s_x(\tilde{y}) - (-l(y, \tilde{y})))^2$ $J(s) = \|s - (-l_q)\|^2 + \text{csl}$

$$J_q(s) \triangleq \mathbb{E}_q J(x, y; s)$$

$$s \in \mathbb{R}^k$$

$$J_q(s) = \min_{\tilde{s} \in \mathbb{R}^k} J_q(\tilde{s})$$

$$= \frac{1}{2k} \|s_x(\cdot) - (-l_q(\cdot))\|_2^2$$

$$\hookrightarrow \in \mathbb{R}^k \quad \vec{l}_{q,x} = \sum_y q(y|x) l(y, \cdot)$$

let \vec{L} be a $k \times k$ matrix

$$\text{where } L_{\tilde{y}, y} = l(y, \tilde{y})$$

$$\vec{l}_{q,x} = \vec{L} \vec{q}_x$$

$$s^* = \vec{L} \vec{q}_x \in \text{span}(\vec{L}) \text{ i.e. } \sum_y \alpha_y l(\cdot, y)$$

to get consistency, it is sufficient to only consider $s \in \text{span}(\vec{L})$ column space

$$\text{or that } s \in \text{span}(F) \supseteq \text{span}(\vec{L})$$

$$\text{restriction on scores } F \in \mathbb{R}^{k \times r}$$

F can be chosen cleverly depending on \vec{L}

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

$$\text{then have } J_q(s) = \min_{s \in \mathbb{R}^k} J_q(s) = \frac{1}{2k} \left\| \vec{F} \vec{s} - (-\vec{l}_q) \right\|_2^2$$

then have $S_q(\theta) - \min_{\theta \in \mathbb{R}^k} S_q(\theta) = \frac{1}{2k} \|\tilde{F}\tilde{\theta} - (-\tilde{L}q)\|_2^2$

$S = F\theta$

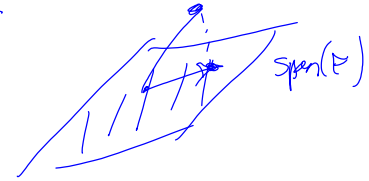
lower bound \rightarrow carries results this side is hard

thm. 7: $H_{\text{Squared}, L, F}(\epsilon) \geq \frac{\epsilon^2}{2k \max_{i \neq j} \|P_F \Delta_{ij}\|_2^2} \geq \frac{\epsilon^2}{4k}$ $|S| \text{ eg } 10^{23}$

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

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

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



in the paper, show that for OI loss, $H(\epsilon) = \frac{\epsilon^2}{4k}$

thm. 8: if $\text{span}(F) \not\subseteq \mathbb{R}^k$ (ie. no constraints) "hardness result"

then $H(\epsilon) \leq \frac{\epsilon^2}{2k}$

ie. for any loss, we need exponential # of samples in the worst case

(current \rightarrow all these are bounds and worst case)

* but for Hamming loss, if add constraint that $S(\tilde{y}) = \sum_{p \leftarrow \text{points}} S_p(\tilde{y}_p)$

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

missing link: going from calibration function to sample complexity

setup: let $S(x, y)$ be of the form $F \circ \Phi(x)$

$$\Theta(x) \in \mathbb{R}^r$$

$$\Theta_j(\cdot) \in \mathcal{H} \quad \text{RKHS}$$

optimization variable

run projected SGD on $J(\Theta)$ i.e. $J(\Theta) = \mathbb{E}_{(x,y) \sim p} J(x, y, \Theta)$

$$\Theta^{(t+1)} = \text{Proj}_{\mathcal{D}} \left[\Theta^{(t)} - \gamma \nabla_{\Theta} J(x^{(t)}, y^{(t)}, \Theta) \right]$$

small ball of radius D

$(x^{(t)}, y^{(t)}) \stackrel{\text{i.i.d.}}{\sim} p$

feature map for \mathcal{H}

$\Phi(x^{(t)})^T$

$K(x^{(t)}, \cdot)$

convergence result: (thm. 5)

if $\|\Theta^*\|_{\mathcal{H}_S} \leq D$

if $\mathbb{E}_{(x,y) \sim p} \|\nabla_{\Theta} J(x, y, \Theta)\|_{\mathcal{H}_S} \leq M^2$

then that averaged SGD with step size $\gamma = \frac{2D}{M\sqrt{n}}$,

$$\mathbb{E} \left[J \left(\frac{1}{n} \sum_{t=1}^n \Theta^{(t)} \right) \right] - J(\Theta^*) \leq \frac{2DM}{\sqrt{n}} \quad (\text{convergence result})$$

thm. 6: Learning complexity

if Θ^* which minimizes $L(\Theta)$, has $\|\Theta^*\|_{\mathcal{H}_S} \leq D$

choosing $n \geq \frac{4D^2M^2}{\epsilon^2}$ implies $\mathbb{E}[L(\bar{\Theta}^n)] < L(\Theta^*) + \epsilon$

define
meaningful
scale

$H^q(\mathcal{E})$

sample complexity
(backgrounds)

in the paper: we compute $D \& M \& H^q(\mathcal{E})$ for specific losses ℓ
and the quadratic \mathcal{J}

to get sample complexity

→ Moral here:

* some losses are harder to learn than others (worse sample complexity)
[0-1 is difficult in general]

→ have linked computation to statistical performance in consistency framework
↳ (convex surrogate loss)

* could handle dependence on x using RKHS

caveats: • distribution-free result (i.e. worst case over all distributions)

→ still need more theory? (e.g. role of $p(y|x)$?)
or other surrogates?