

Lecture 18 - scribbles - consistency

Wednesday, March 22, 2017

10:40

today: • consistency
• CRF objective
and variance reduced SGD

Consistency & calibration function

need to relate $J(w)$ to $L(w)$ \rightarrow "calibration function" [Steinwart]

relationship is usually very complicated;

\Rightarrow current results look mainly at: non-parametric setting
where dependence on x is evacuated (?)

\rightarrow we suppose that $s(x, y; w)$ can be arbitrary for any x
(i.e. w is so -dim)

[can use a universal kernel to formalize this and
re-include x]

* binary classification, Bartlett & al. characterized a whole family
of consistent surrogate losses

for multiclass: [Lee & al. 2004] showed that multiclass hinge loss is not consistent
McAllester 2002

for 0-1 loss
when no majority class (i.e. $p(y|x) \leq \frac{1}{2} \forall y$)

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

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

2) $l(y, y')$ much richer

\rightarrow proposed a fix with a surrogate loss that has $\sum_{\tilde{y}} \dots$ instead of $\max_{\tilde{y}}$

exponential $\# \Rightarrow$ might cause problem?

[Osokin & al. 2017] \rightarrow looked at this exponential aspect more carefully, in the simplest setup

calibration function for structured cost ℓ , surrogate loss \mathcal{J} , and set \mathcal{W}

$$H_{\mathcal{J}, \ell, \mathcal{W}}(\varepsilon) \triangleq \inf_{w \in \mathcal{W}, q \in \Delta(\mathcal{Y}) \text{ s.t. } \mathcal{J}_q(w) - \min_{w' \in \mathcal{W}} \mathcal{J}_q(w') \geq \varepsilon} \mathcal{J}_q(w) - \min_{w' \in \mathcal{W}} \mathcal{J}_q(w')$$

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

$$\mathcal{J}_q(w) \triangleq \mathbb{E}_{q(y)} [\mathcal{J}(x, y, w)]$$

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

Smallest optimization surrogate regret possible (over all dist. q)
 \Rightarrow true regret is $\geq \varepsilon$

consequence (Thm. 2)

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

convex envelope of $H(\varepsilon)$ i.e. H^{**}

if \check{H} is invertible
(note: it is increasing as Gauthier pointed out)

$$L_q(w) - L_q^* \leq \check{H}^{-1}(\mathcal{J}_q(w) - \mathcal{J}_q^*)$$

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

Sample complexity:

can link learning with optimization using SGD : while running SGD with $y^{(i)} \sim q$
we are optimization population \mathcal{J}_q
(instead of \mathcal{J})

thus from SGD convergence on \mathcal{P}_1 can see how many iterations needed

thus from SGD convergence on l_q , can see how many iterations needed
 (= # of samples)
 to get surrogate-regret of $H(\epsilon)$

\Rightarrow translates to ϵ -true regret

e.g. constant step-size SGD projected on a ball of diameter D
 gives $\mathbb{E}[l_q(w_t)] - l_q^* \leq \frac{DM}{\sqrt{t}}$ after t -iterations

$$M^2 \geq \mathbb{E}_{q,y} [\|\nabla_w l(x,y,w)\|^2]$$

(thm 4): to get $\mathbb{E}[l_q(w_t)] - l_q^* \leq \epsilon$
 need at most $t \geq \left(\frac{DM}{H(\epsilon)} \right)^2$

gives
 (sample complexity)

in our paper, consider the consistent convex surrogate

$$l(y,w) \triangleq \frac{1}{|S|} \sum_{\tilde{y}} \frac{1}{2} (S_w(\tilde{y}) + l(y,\tilde{y}))^2$$

$$l_q(w) = \mathbb{E}_q[l(y,w)] = \frac{1}{|S|} \sum_{\tilde{y}} \frac{1}{2} (S_w(\tilde{y}) + \underbrace{\mathbb{E}_q[l(y,\tilde{y})]}_{l_q(\tilde{y})})^2 + \text{const.}$$

$$\text{so } w^* \text{ is s.t. } S_{w^*}(\tilde{y}) = -l_q(\tilde{y})$$

l is consistent for any l and q
 as long as $S_w \supseteq \text{span}\{l(y, \cdot) : y \in |S|\}$
 $S_w \triangleq \{S_w(\cdot) : w \in W\}$

max score \Rightarrow min conditional risk
 and thus consistent

\Rightarrow in paper, we show that if no constraint on $S_w(\cdot)$

then $H(\epsilon) \leq \epsilon^2$ for any l

$2T\epsilon$ } need exponential accuracy \Rightarrow exponential sample complexity?

also, for 0-1 loss, to be consistent, we basically need no constraint on $S(\cdot)$

but for flaming loss, if add constraint that $S(\tilde{y}) = \sum_p \mathbb{I}(\tilde{y}_p)$

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

⊗ Moral here:!

- some losses are harder to learn than others (0-1 difficult in general)
- have linked computation to statistical performance in consistency framework
 \hookrightarrow (convex surrogate loss)

but lost dependence on α (need to use SGD with RKHS ^(kernel stuff) \rightarrow fall in defaults)

\rightarrow different approach than gen. error bound; but gives insights

- still need more theory!

- main pointer covered today:
 - Anton Osokin, Francis Bach, Simon Lacoste-Julien
On Structured Prediction Theory with Calibrated Convex Surrogate Losses
<https://arxiv.org/abs/1703.02403>
- other pointers:
 - canonical paper which presented consistency analysis for binary classification:
 Bartlett, Peter L., Jordan, Michael I., and McAuliffe, Jon D
Convexity, classification, and risk bound
[Journal of the American Statistical Association](#), 101(473):138–156, 2006.
 - paper which showed that multiclass SVM is not consistent (for the 0-1 loss) and proposed a consistent alternative:
 Lee, Yoonkyung, Lin, Yi, and Wahba, Grace.
Multicategory support vector machines: Theory and application to the classification of microarray data and satellite radiance data.
[Journal of the American Statistical Association](#), 99(465):67–81, 2004.
 - see also the McAllester 2007 paper:

Generalization Bounds and Consistency for Structured Labeling in Predicting Structured Data, edited by G. Bakir, T. Hofmann, B. Scholkopf, A. Smola, B. Taskar, and S. V. N. Vishwanathan. MIT Press, 2007

<http://nagoya.uchicago.edu/~dmcalister/colbounds.pdf>

- and interestingly, this recent paper shows that the multiclass SVM *is* consistent for a loss on 3 classes with an "abstain" notion: Ramaswamy, Harish G. and Agarwal, Shivani.
Convex calibration dimension for multiclass loss matrices.
[JMLR](#), 17(14):1–45, 2016.
- see also the extensive related work section of the arxiv 2017 paper Osokin et al.