

Ch 13 part 2 Analysis and Stability

Thursday, March 19, 2020

4:10 PM

(following Shalev-Shwartz and Ben-David)

Analysis of \textcircled{II} term

Recall $\textcircled{II} = \mathbb{E}_S \hat{L}_S(A(S))$
 $L = \arg \min_{w \in \mathcal{H}} \hat{L}_S(w) + \lambda \|w\|^2$

The larger λ is,

the larger \textcircled{II} becomes: this represents more bias
or underfitting

So unlike \textcircled{I} term,
now want λ small

To analyze, note that since $\lambda \|w\|^2 \geq 0$,

$$\begin{aligned} (\forall w) \quad \hat{L}_S(A(S)) &\leq \hat{L}_S(A(S)) + \lambda \|A(S)\|^2 \\ &\leq \hat{L}_S(w) + \lambda \|w\|^2 \end{aligned}$$

Hence

$$\begin{aligned} \textcircled{II} := \mathbb{E}_S \hat{L}_S(A(S)) &\leq \mathbb{E}_S \hat{L}_S(w) + \lambda \|w\|^2 (\forall w \in \mathcal{H}) \\ &= L_D(w) + \lambda \|w\|^2 \end{aligned}$$



If w fixed, then draw S iid \mathcal{D}^m ,

$$\mathbb{E}_S \hat{L}_S(w) = L_D(w) \quad \text{but} \quad \mathbb{E}_S \hat{L}_S(A(S)) \neq L_D(A(S))$$

So we can choose the best (or at least good) w to minimize this bound (such a w is an "oracle")

(and can also choose λ as we wish ...)

though in practice choose λ via cross-validation or similar)

For example, and using a bound on \textcircled{I} from last lecture,

Corollary 13.9 Let (\mathcal{H}, Z, ℓ) be a convex, ρ -Lipschitz, B -bounded (i.e. $\forall w \in \mathcal{H}, \|w\| \leq B$), and $S \stackrel{\text{iid}}{\sim} \mathcal{D}^m$, then set

$$\lambda = \sqrt{\frac{2\rho^2}{B^2 m}} \quad \text{then if } A \text{ is RLM } w, R(w) = \lambda \|w\|^2,$$

$$\mathbb{E}_S L_D(A(S)) \leq \min_{w \in \mathcal{H}} L_D(w) + \underbrace{\rho B \sqrt{\frac{8}{m}}}_{\leq \varepsilon \text{ if } m \geq \frac{8\rho^2 B^2}{\varepsilon^2}}$$

proof sketch:

$$\mathbb{E} L_D(A(S)) = \mathbb{E} \left(\underbrace{L_D(A(S)) - \hat{L}_S(A(S))}_{\textcircled{I}} \right) + \mathbb{E} \underbrace{\hat{L}_S(A(S))}_{\textcircled{II}}$$

$$\leq \frac{2\rho^2}{\lambda m} + L_D(w) + \underbrace{\lambda \|w\|^2}_{\leq \lambda B^2} \quad \forall w \in \mathcal{H}$$

$$= L_D(w) + \rho B \sqrt{\frac{8}{m}} \quad \left(\text{since } \frac{2\rho^2}{\lambda m} + \lambda m B^2 = \frac{4\rho^2}{\lambda m} \right)$$

and $\forall w \in \mathcal{H}$ so take min.

$$= \sqrt{\frac{16 B^2 \rho^4 m}{2\rho^2 m}} \quad \square$$

A similar bound, w) different assumptions

(if X bounded, these are stronger assumptions, up to constants)

Corollary 13.11 Let (\mathcal{H}, Z, l) be a convex, β -smooth, B -bounded,
 and $(\forall z \in Z) l(0, z) \leq 1$ then $(\forall \varepsilon > 0) (\forall \text{distr } \mathcal{D})$, if
 $m \geq \frac{150\beta B^2}{\varepsilon^2}$ and $\lambda = \frac{\varepsilon}{3B^2}$, then

$$\mathbb{E} L_{\mathcal{D}}(A(S)) \leq \left(\min_{w \in \mathcal{H}} L_{\mathcal{D}}(w) \right) + \left(\frac{1}{3} \right) \cdot \varepsilon$$

proof:

$$\mathbb{E} L_{\mathcal{D}}(A(S)) = \mathbb{E} \hat{L}_S(A(S)) + \underbrace{\mathbb{I}}_{\substack{\text{bdd from} \\ \text{prev. lecture}}} \leq \frac{48\beta}{\lambda m} \mathbb{E} \hat{L}_S(A(S))$$

$$(\text{and } \mathbb{E} \hat{L}_S(A(S)) \leq L_{\mathcal{D}}(w) + \lambda \|w\|^2 \quad \forall w \in \mathcal{H})$$

$$\leq \left(L_{\mathcal{D}}(w) + \lambda \|w\|^2 \right) + \frac{48\beta}{\lambda m} \left(L_{\mathcal{D}}(\bar{w}) + \lambda \|\bar{w}\|^2 \right)$$

choose $w \in \arg \min L_{\mathcal{D}}(w)$ choose $\bar{w} = 0$

$$\leq \min_{w \in \mathcal{H}} L_{\mathcal{D}}(w) + \lambda B^2 + \frac{48\beta}{\lambda m} (C + 0)$$

$\because \lambda m \geq \frac{50\beta}{\varepsilon}$

$$\leq \min_{w \in \mathcal{H}} L_{\mathcal{D}}(w) + \frac{\varepsilon}{3} + \frac{48}{50} \varepsilon \leq \frac{1}{3} \varepsilon \quad \square$$

History of regularization + stability

Andrey Tikhonov 1906-1993 (sp. Tychonoff in topology) 1943
 ubiquitous in math

Stability for learning: Roger, Wager '78 for k-NN Bousquet + Eliseef '02

Bagging (bootstrap aggregating) by Breiman '96
 is used to increase stability of unstable algo

Boosting create strong learner from weak ones
Bagging create stable learner from unstable ones

Modern example:

Train faster, generalize better: Stability of stochastic gradient descent" (Haratt, Recht, Singer ICM L'16)

Mohri et al.'s perspective ch 14 "Algorithmic Stability"

Emphasize benefit of this analysis (vs uniform conv.)
is tailored to algorithm

defines 'uniform stability' (can bound \mathbb{E})

and uses McDiarmid's ineq to get results (Thm 14.2)

Thm 14.2 loss $\leq M$, algo A is β -uniformly stable then
w.p. $> 1-\delta$,

$$L_D(A(S)) \leq \hat{L}_S(A(S)) + \beta + (2m\beta + M) \sqrt{\frac{\log(1/\delta)}{2m}}$$

and Prop. 14.4 shows if $\|x\|^2 \leq r^2 \forall x \in X$,

L is convex and ρ -admissible (similar to Lipschitz)

then RLM $\hookrightarrow R(w) = \lambda \|w\|^2$ is β -stable, $\beta \leq \frac{\rho^2 r^2}{m\lambda}$

Trickier to see λ tradeoff

$\lambda \rightarrow \infty \Rightarrow \beta \rightarrow 0$ (good)

but $\hat{L}_S(A(S))$ not as close to
 $\min_{w \in \mathcal{H}} \hat{L}_S(w)$ (=ERM)