

Convergence Rates of λ

October 26, 2016

Let's bound $\|\hat{\lambda} - \tilde{\lambda}\|$ instead.

This will actually get rid of the geometric mean term in Theorems 1 and 3.

We will suppose that the data is generated from the model:

$$y = g^*(x) + \epsilon$$

where ϵ are independent, sub-gaussian errors. The penalized regression models are

$$\hat{g}(\cdot|\lambda) = \arg \min_{g \in \mathcal{G}} L_T(g|\lambda)$$

Let the model class after fitting on the training data be

$$\mathcal{G}(T) = \{\hat{g}(\cdot|\lambda) : \lambda \in \Lambda\}$$

The selected penalty parameters are

$$\hat{\lambda} = \arg \min_{\lambda \in \Lambda} \|y - \hat{g}(\cdot|\lambda)\|_V^2$$

Convergence of $\hat{\lambda}$ to $\tilde{\lambda}$

Suppose that if $\|\epsilon\|_T \leq 2\sigma$, then $\mathcal{G}(T)$ satisfies the entropy condition

$$\int_0^R H^{1/2}(u, \mathcal{G}(T), \|\cdot\|_V) du \leq \psi_T(R)$$

Furthermore, suppose that

$$\frac{\psi_T(u)}{u^2}$$

is nonincreasing wrt to u for all $u > 0$.

Let $L_V^*(\lambda) = \|\hat{g}(\cdot|\lambda) - g^*\|_V^2$ be the true validation loss and let $\tilde{\lambda}$ be the global minimizer of $L_V^*(\lambda)$.

$$\tilde{\lambda} = \arg \min_{\lambda} L_V^*(\lambda)$$

Let $\tilde{\lambda}_{gen}$ be the global minimizer of the generalization error.

$$\tilde{\lambda}_{gen} = \arg \min_{\lambda} E_V [L_V^*(\lambda)] = \arg \min_{\lambda} \|\hat{g}(\cdot|\lambda) - g^*\|^2$$

- **Local strong convexity assumption:** Suppose that there is a neighborhood $N(\tilde{\lambda}_{gen})$ around $\tilde{\lambda}_{gen}$ such that the true validation loss is smooth in λ for all $\lambda \in N(\tilde{\lambda}_{gen})$ and for all $\lambda \in N(\tilde{\lambda})$, the true validation loss is m -strongly convex in λ for some $m > 0$:

$$\nabla_{\lambda}^2 L_V^*(\lambda) = \nabla_{\lambda}^2 \|\hat{g}(\cdot|\lambda) - g^*\|_V^2 \succeq mI$$

Important: m cannot shrink in n_T and n_V

- **Lipschitz assumption:** Let us also assume that for all $\lambda \in N(\tilde{\lambda}_{gen})$, the fitted functions are locally K -Lipschitz:

$$\left\| \hat{g}(\cdot|\lambda) - \hat{g}(\cdot|\tilde{\lambda}) \right\|_V \leq K \|\lambda - \tilde{\lambda}\|$$

Important: K cannot grow in n_V, n_T (I don't even think this holds for ridge regression though...)

If δ is chosen such that

$$\sqrt{n_V} \delta^2 \geq 2C [\psi_T(2\delta) \vee (2\delta)]$$

then we have

$$Pr \left(\|\hat{\lambda} - \tilde{\lambda}\| \geq \delta \wedge \|\epsilon\|_V \leq 2\sigma \wedge \|\epsilon\|_T \leq 2\sigma \right) \leq c \exp \left(-\frac{n_V \delta^2 m^2}{c^2 K^2} \right)$$

for some constant c .

Furthermore, this completely removes the geometric term since we also have that

$$\begin{aligned} Pr \left(\left\| \hat{g}(\cdot|\lambda) - \hat{g}(\cdot|\tilde{\lambda}) \right\|_V \geq \frac{\delta}{K} \wedge \|\epsilon\|_V \leq 2\sigma \wedge \|\epsilon\|_T \leq 2\sigma \right) &\leq Pr \left(\|\hat{\lambda} - \tilde{\lambda}\| \geq \delta \wedge \|\epsilon\|_V \leq 2\sigma \wedge \|\epsilon\|_T \leq 2\sigma \right) \\ &\leq c \exp \left(-\frac{n_V \delta^2 m^2}{c^2 K^2} \right) \end{aligned}$$

Proof

Let $\hat{\lambda}$ be the global minimizer of the validation loss. Therefore

$$\left\| y - \hat{g}(\cdot|\hat{\lambda}) \right\|_V^2 \leq \left\| y - \hat{g}(\cdot|\tilde{\lambda}) \right\|_V^2$$

Recall the basic inequality

$$L_V^*(\hat{\lambda}) - L_V^*(\tilde{\lambda}) = \left\| g^* - \hat{g}(\cdot|\hat{\lambda}) \right\|_V^2 - \left\| g^* - \hat{g}(\cdot|\tilde{\lambda}) \right\|_V^2 \leq 2 \left\langle \epsilon, \hat{g}(\cdot|\hat{\lambda}) - \hat{g}(\cdot|\tilde{\lambda}) \right\rangle_V$$

Suppose that $\hat{\lambda} \in N(\tilde{\lambda}_{gen})$. Using the mean value theorem and the local strong convexity assumption, there is some $\alpha \in (0, 1)$ such that

$$\begin{aligned} L_V^*(\hat{\lambda}) - L_V^*(\tilde{\lambda}) &= (\hat{\lambda} - \tilde{\lambda})^\top \nabla_{\tilde{\lambda}}^2 L_V^*(\lambda)|_{\lambda=\alpha\tilde{\lambda}+(1-\alpha)\hat{\lambda}} (\hat{\lambda} - \tilde{\lambda}) \\ &\geq \left(\min_{\alpha} \nabla_{\tilde{\lambda}}^2 L_V^*(\lambda)|_{\lambda=\alpha\tilde{\lambda}+(1-\alpha)\hat{\lambda}} \right) \|\hat{\lambda} - \tilde{\lambda}\|_2^2 \\ &\geq m \|\hat{\lambda} - \tilde{\lambda}\|_2^2 \end{aligned}$$

Therefore we get

$$m \|\hat{\lambda} - \tilde{\lambda}\|_2^2 \leq 2 \left\langle \epsilon, \hat{g}(\cdot|\hat{\lambda}) - \hat{g}(\cdot|\tilde{\lambda}) \right\rangle_V$$

Anyhow, if we assume local strong convexity and the Lipschitz condition, we can proceed with a peeling argument

$$\begin{aligned} Pr \left(\|\hat{\lambda} - \tilde{\lambda}\| \geq \delta \right) &= \sum_{s=0}^{\infty} Pr \left(2^s \delta \leq \|\hat{\lambda} - \tilde{\lambda}\| \leq 2^{s+1} \delta \right) \\ &= \sum_{s=0}^{\infty} Pr \left(\|\hat{\lambda} - \tilde{\lambda}\|^2 \geq 2^{2s} \delta^2 \wedge \|\hat{\lambda} - \tilde{\lambda}\| \leq 2^{s+1} \delta \right) \end{aligned}$$

$$\begin{aligned}
&= \sum_{s=0}^{\infty} Pr \left(2 \left\langle \epsilon, \hat{g}(\cdot|\hat{\lambda}) - \hat{g}(\cdot|\tilde{\lambda}) \right\rangle_V \geq m 2^{2s} \delta^2 \wedge \left\| \hat{g}(\cdot|\hat{\lambda}) - \hat{g}(\cdot|\tilde{\lambda}) \right\|_V \leq 2^{s+1} \delta K \right) \\
&\leq \sum_{s=0}^{\infty} Pr \left(\sup_{\lambda, \lambda' \in N(\tilde{\lambda}_{gen}): \|\hat{g}(\cdot|\lambda) - \hat{g}(\cdot|\lambda')\|_V \leq 2^{s+1} \delta K} \left\langle \epsilon, \hat{g}(\cdot|\lambda) - \hat{g}(\cdot|\lambda') \right\rangle_V \geq m 2^{2s-1} \delta^2 \right)
\end{aligned}$$

To apply the lemma based on vandegeer corollary 8.3 (see below), we must check all the conditions are satisfied.

We choose δ such that

$$\begin{aligned}
\frac{\sqrt{n_V}}{8} &\geq \frac{C}{4\delta^2} [\psi_T(2\delta) \vee (2\delta)] \\
&\geq \frac{C}{2^{2s+2}\delta^2} [\psi_T(2^{s+1}\delta) \vee (2^{s+1}\delta)]
\end{aligned}$$

where the second line follows from the assumption that $\psi_T(u)/u^2$ is nonincreasing wrt u . Hence we have satisfied the condition in corollary 8.3. So for all $s = 0, 1, \dots$ since

$$\sqrt{n_V} 2^{2s-1} \delta^2 \geq C [\psi_T(2^{s+1}\delta) \vee (2^{s+1}\delta)]$$

we have

$$Pr \left(\sup_{\lambda, \lambda' \in N(\tilde{\lambda}_{gen}): \|\hat{g}(\cdot|\lambda) - \hat{g}(\cdot|\lambda')\|_V \leq 2^{s+1} \delta K} \left\langle \epsilon, \hat{g}(\cdot|\lambda) - \hat{g}(\cdot|\lambda') \right\rangle_V \geq m 2^{2s-1} \delta^2 \wedge \|\epsilon\|_V \leq 2\sigma \wedge \|\epsilon\|_T \leq 2\sigma \right) \leq \exp \left(-n_V \frac{2^{4s}}{4C^2 (2^{s+1}\delta)^2} \right)$$

Hence we have

$$\begin{aligned}
Pr \left(\|\hat{\lambda} - \tilde{\lambda}\| \geq \delta \wedge \|\epsilon\|_V \leq 2\sigma \wedge \|\epsilon\|_T \leq 2\sigma \right) &\leq C \sum_{s=0}^{\infty} \exp \left(-n_V \frac{2^{4s-2} \delta^4 m^2}{4C^2 (2^{s+1}\delta)^2 K^2} \right) \\
&\leq c \exp \left(-\frac{n_V \delta^2}{c^2} \right)
\end{aligned}$$

for some constant c .

Example

If \mathcal{G} is a parametric family, we must choose

$$\delta \geq R \left(n_V^{-1/2} \right)$$

for some constant R , so we have an asymptotic convergence rate for

$$Pr \left(\|\hat{\lambda} - \tilde{\lambda}\| \geq R n_V^{-1/2} \wedge \|\epsilon\|_V \leq 2\sigma \wedge \|\epsilon\|_T \leq 2\sigma \right) \leq c \exp \left(-\frac{R^2 m^2}{c^2 K^2} \right)$$

Jean's questions

- Locally m -strongly convex where m doesn't shrink with n_T or n_V seems like a strong assumption.
- Fitted functions are locally K -Lipschitz where K doesn't change with n_T seems like a strong assumption too.
- Counterexample: In ridge regression, if p grows with n , I believe the Lipschitz constant is on the order of λ_{min}^{-2} . But if λ_{min} is shrinking with n , then K is changing with n .

- We need to ensure that $N(\tilde{\lambda}_{gen})$ contains $\tilde{\lambda}$ and $\hat{\lambda}$ with high probability - can we ensure this?
Thoughts:

- We might be able to bootstrap results to show that $N(\tilde{\lambda}_{gen})$ contains $\tilde{\lambda}$ and $\hat{\lambda}$ with high probability.
- We know that the difference

$$\|\hat{g}_{\hat{\lambda}} - g^*\|_V - \|\hat{g}_{\tilde{\lambda}} - g^*\|_V \leq \delta$$

with “high probability”. So if the global minimizer of $\|\hat{g}_{\tilde{\lambda}} - g^*\|_V$ is more than δ smaller than all other local minimas of $L_V^*(\lambda)$, then $\hat{\lambda}$ will be located in the same region as $\tilde{\lambda}$. By definition, this region must be quasi-convex.

- If we can show that this region is $N(\tilde{\lambda}_{gen})$ and it is strongly convex, we’d be done. The problem is if this region has crazy behavior (e.g. the loss is very flat wrt λ)