

Lecture 26. Nonparametric Regression

Recap so far

- Reproducing Kernel Hilbert space perspective
- Bias-variance decomposition to analyze error of Kernel ridge regression

- Power law generalization $\mu_j = j^{-\alpha} \Rightarrow$ generalization error of $n^{\frac{1}{\alpha}-1}$
($\alpha > 1$)

↑
"noise removal", not really generalization

This time.

- Look at generalization to test distⁿ
- Mercer's theorem: important analytic tool
 - ↳ Useful for computing μ_j
 - ↳ Helps draw connections b/w random features and Kernel regression

Notation.

$$x_1, \dots, x_n \sim p$$

$$y_i = \langle \beta^*, x_i \rangle + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\hookrightarrow f^*(x_i) + \varepsilon_i$$

$$S = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$$S^* = \mathbb{E}_{x \sim p} [X X^T]$$

$L^*(\beta)$: excess risk at population level

$$\Rightarrow L^*(\beta) = (\beta - \beta^*)^T S^* (\beta - \beta^*)$$

$L(\beta)$: excess risk on training samples

$$= \frac{1}{n} \sum_{i=1}^n \langle \beta - \beta^*, x_i \rangle^2 = (\beta - \beta^*)^T S (\beta - \beta^*)$$

Ridge regression estimator

$$\hat{\beta}_\lambda = \frac{1}{n} \left(\frac{\lambda}{n} I + S \right)^{-1} X^T (X \beta^* + \varepsilon)$$

$$= (\lambda I + X^T X)^{-1} X^T y$$

$$\mathbb{E}_\varepsilon [\hat{\beta}_\lambda] = \left(\frac{\lambda}{n} I + S \right)^{-1} S \beta^*$$



$$\delta = \frac{\lambda}{n} \left(\frac{\lambda}{n} I + S \right)^{-1} \beta^*$$

Bias: $\delta^T S^* \delta$, $\delta = \mathbb{E}_\varepsilon[\hat{\beta}_\lambda] - \beta^*$

$$\text{Bias}^2 = \left(\frac{\lambda}{n}\right)^2 (\beta^*)^T \left(\frac{\lambda}{n} I + S\right)^{-1} S^* \left(\frac{\lambda}{n} I + S\right)^{-1} \beta^*$$

Contrast to before:
Bias² was wrt. train set $\Rightarrow S$ instead of S^*

old formula:

$$\left(\frac{\lambda}{n}\right)^2 (\beta^*)^T \left(\frac{\lambda}{n} I + S\right)^{-1} S \left(\frac{\lambda}{n} I + S\right)^{-1} \beta^*$$

$$\lambda_n = \frac{\lambda}{n}$$

$\mathbb{E}[\text{Bias}^2]$ under Gaussian prior $\mathbb{E}[\beta^* (\beta^*)^T] = \rho^2 I$

$$\mathbb{E}[\text{Bias}^2] = \lambda_n^2 \rho^2 \text{trace} \left((\lambda_n I + S)^{-1} S^* (\lambda_n I + S)^{-1} \right)$$

$$\text{Variance} = \frac{\sigma^2}{n} \text{trace} \left((S^*)^{1/2} (\lambda_n I + S)^{-1} S (\lambda_n I + S)^{-1} (S^*)^{1/2} \right)$$

old formula: $\frac{\sigma^2}{n} \text{trace} \left((\lambda_n I + S)^{-1} S^2 (\lambda_n I + S)^{-1} \right)$

Good generalization was Bias² + Variance to be small

Issue S instead of S^* ← generally assume is nice

For bias $(\lambda_n I + S)^{-1} \approx (\lambda_n I + S^*)^{-1}$

$$(\lambda_n I + S)^{-1} \leq 2 (\lambda_n I + S^*)^{-1}$$

← Want this to be true.

⇒ can call S w/ S^* lossy only a factor of 2

For variance

$$(\lambda_n I + S)^{-1} S (\lambda_n I + S)^{-1} \leq (\lambda_n I + S)^{-1}$$

since $S (\lambda_n I + S)^{-1} \leq I$

$$\text{Variance} \leq \frac{\sigma^2}{n} \text{trace} \left((S^*)^{1/2} (\lambda_n I + S)^{-1} (S^*)^{1/2} \right)$$

Once again, as long as $(\lambda_n I + S)^{-1} \leq 2 (\lambda_n I + S^*)^{-1}$,
lose at most a factor of 2 moving to S^* .

$$\Rightarrow 2 (\lambda_n I + S) \geq \lambda_n I + S^*$$

$$\Rightarrow S^* \leq \lambda_n I + 2S$$

takes care of small eigenvalues

Need concentration on large eigenvalues

Proposition Suppose $S^* \leq \lambda_n I + 2S$.

Let μ_j^* = j^{th} largest eigenvalue of S^* .

Then, generalization error satisfies

$$\mathbb{E}_{\beta^*} \left[\mathbb{E}_{\varepsilon} \left[L^*(\beta_{\lambda}) - L^*(\beta^*) \right] \right]$$

$$\leq 4\lambda_n \rho^2 \sum_j \frac{\lambda_n \mu_j^*}{(\lambda_n + \mu_j^*)^2} + \frac{2\sigma^2}{n} \sum_j \frac{\mu_j^*}{\lambda_n + \mu_j^*}$$

$$\leq \left(4\lambda_n \rho^2 + \frac{2\sigma^2}{n} \right) \sum_j \min\left(1, \frac{\mu_j^*}{\lambda_n}\right).$$

$$\frac{\lambda_n}{\lambda_n + \mu_j^*} \leq 1 \Rightarrow \frac{\lambda_n \mu_j^*}{(\lambda_n + \mu_j^*)^2} \leq \frac{\mu_j^*}{\lambda_n + \mu_j^*} \leq \min\left(1, \frac{\mu_j^*}{\lambda_n}\right)$$

$$\mathbb{E}[\text{Bias}^2] = \lambda_n^2 \rho^2 \text{trace} \left((\lambda_n I + S)^{-1} S^* (\lambda_n I + S)^{-1} \right)$$

$$\text{Variance} = \frac{\sigma^2}{n} \text{trace} \left((S^*)^{1/2} (\lambda_n I + S)^{-1} S (\lambda_n I + S)^{-1} (S^*)^{1/2} \right)$$

$$\mathbb{E}[\text{Bias}^2] \leq 4 \lambda_n^2 \rho^2 \text{trace} \left((\lambda_n \mathbf{I} + S^*)^{-1} S^* (\lambda_n \mathbf{I} + S^*)^{-1} \right)$$

$$\parallel$$

$$\sum_j \frac{\mu_j^*}{(\lambda_n + \mu_j^*)^2}$$

$\mu_j = n_j^{-\alpha}$ K

Corollary? If we plug $\mu_j^* = j^{-\alpha}$, (S) get same $n^{\frac{1}{\alpha}-1}$ rate as before.

Upside: Really do just need to care about when $S^* \leq \lambda_n \mathbf{I} + 2\mathcal{B}$.

$$S \geq \frac{1}{2} (S^* + \lambda_n \mathbf{I}) - \lambda_n \mathbf{I}$$

$$\underbrace{(S^* + \lambda_n \mathbf{I})^{-1/2}}_{\text{simple}} S (S^* + \lambda_n \mathbf{I})^{-1/2} \geq \frac{1}{2} \mathbf{I} - \lambda_n (S^* + \lambda_n \mathbf{I})^{-1}$$

looks like a covariance matrix Σ

find something like Σ^*

S : sample covariance matrix of X_i

Σ : sample covariance matrix of $(S^* + \lambda_n I)^{-1/2} X_i$

$$\Sigma^* = (S^* + \lambda_n I)^{-1/2} S^* (S^* + \lambda_n I)^{-1/2} \\ = S^* (S^* + \lambda_n I)^{-1}$$

$$\frac{1}{2} I = I - \frac{1}{2} I$$

↓

$$I = (S^* + \lambda_n I) (S^* + \lambda_n I)^{-1}$$

$$\cancel{(S^* + \lambda_n I)} (S^* + \lambda_n I)^{-1} - \cancel{\lambda_n I} (S^* + \lambda_n I)^{-1} - \frac{1}{2} I$$

$$= S^* (S^* + \lambda_n I)^{-1} - \frac{1}{2} I = \Sigma^* - \frac{1}{2} I$$

Upshot Condition is equivalent to asking that

$$\hat{\Sigma} \succeq \Sigma^* - \frac{1}{2} I$$

Suffices to show that $\|\Sigma^* - \hat{\Sigma}\|_{op} \leq \frac{1}{2}$.

Theorem (Koltchinskii & Lounici)

If x_i are Gaussian w/ covariance matrix Σ^* , then

$$\|\hat{\Sigma} - \Sigma^*\| \leq O\left(\max\left(\frac{\text{tr}(\Sigma^*)}{n}, \sqrt{\frac{\text{tr}(\Sigma^*) \|\Sigma^*\|}{n}}\right)\right)$$

~~$$\text{tr}(\Sigma^*) \leq p \cdot \|\Sigma^*\|$$~~

~~$$\Rightarrow \sqrt{\frac{p}{n}} \|\Sigma^*\|$$~~

$p = \infty$

Recall. $\Sigma^{\#} = S^{\#} (\lambda_n I + S^{\#})^{-1}$

$$\Rightarrow \|\Sigma^{\#}\| \leq 1$$

$$\text{trace}(\Sigma^{\#}) = \sum_j \frac{\mu_j^{\#}}{\lambda_n + \mu_j^{\#}}$$

\Rightarrow As long as

$$\sum_j \frac{\mu_j^{\#}}{\lambda_n + \mu_j^{\#}} \ll n,$$

$$\|\Sigma^{\#} - \hat{\Sigma}\|_{\text{op}} \leq \frac{1}{2}$$

$$\Rightarrow S^{\#} \leq 2S + \lambda_n I$$

If $\mu_j^* = j^{-\alpha}$, Condition
 works out to $\lambda_n \gg \frac{\alpha}{\alpha-1} n^{-\alpha}$

Before, optimal $\lambda_n = \frac{\sigma^2}{np^2}$

$\frac{1}{n}$

Upshot Generally fine in terms
 of generalization, UNLESS
 σ^2 is really small.

If we kind propagate lower bound

on λ_n , regularity bound on

generalization error is

$$\mathbb{E}[\text{Error}] \leq O\left(\frac{\alpha}{\alpha-1} \cdot \max\left(n^{\frac{1}{\alpha}-1} (\sigma^2)^{\frac{\alpha-1}{\alpha}} (\rho^2)^{\frac{1}{\alpha}}, n^{1-\alpha} \rho^2\right)\right)$$

what we had before

new term that accounts for generalization

Mercer's theorem

Definition, Eigenvalues of S, K span

Important for understanding generalization.

$$K = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

What about eigenvalues of K ?

Define operator T_K

$$T_K f(x) = \int_{\mathcal{X}} k(x, s) f(s) d\nu(s)$$

arbitrary
finite, positive
measure ν

$$"(T_K f)_x = \sum_s k_{xs} f_s"$$

Theorem (Mercer)

Suppose $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a positive definite kernel, k is continuous, and \mathcal{X} is compact w/ finite positive Borel measure ν supported on all of \mathcal{X} .

Then, there is an orthonormal basis

(.)

(wrt $\langle f, g \rangle = \int f(x)g(x) dV(x)$)

of eigenfunctions e_1, e_2, \dots of T_k , w/ corresponding eigenvalues μ_1, μ_2, \dots , s.t.

$$K(x, y) = \sum_{l=1}^{\infty} \mu_l e_l(x) e_l(y).$$

Observations,

3 inner products

$$\langle a, a \rangle_H,$$

$$\langle a, a \rangle_{L^2}$$

(function error)

$$\langle a, a \rangle_{L^2(V)}$$

← $\sqrt{\quad}$ arbitrary
useful for
computation

Consequences:

• Any f can be represented as

$$f(x) = \sum_l c_l e_l(x)$$

$$c_l = \int_{\mathcal{X}} f(x) e_l(x) d\nu(x)$$

$$= \langle f, e_l \rangle_{L^2(\nu)}$$

$$\bullet \|f\|_{\mathcal{H}}^2 = \sum_{l=1}^{\infty} \frac{c_l^2}{\mu_l} \quad (\text{vs. } \|f\|_{L^2(\nu)}^2 = \sum_{l=1}^{\infty} c_l^2)$$

INDEPENDENT OF ν

$$\bullet \text{ If we take } \nu = \rho, \quad \langle \cdot, \cdot \rangle_{\mathcal{H}} = \langle \cdot, \cdot \rangle_{L^2(\rho)}$$