

# Non-parametric II

Reminder: Pset of Tue next Tues.

Re-cap last time.

Parameter view  
Norm  $\|\beta\|_2^2$

Function view  
 $f^T K^{-1} f$

Error  $(\beta - \beta^*)^T S (\beta - \beta^*)$

$\frac{1}{n} \|f - f^*\|_2^2$

Kernel ridge regression:

$$f = K(\lambda I + K)^{-1} y$$

$$f(x) = K_x^T (\lambda I + K)^{-1} y$$

Define RKHS:

$$\mathcal{H} = \left\{ f \mid \underbrace{\sup_T f[T]^T K[T]^{-1} f[T]}_{\|f\|_{\mathcal{H}}^2} < \infty \right\}$$

Two key results:

(1)  $\|f\|_{K[T]}$  monotonic in  $T$

(2) Isotropic extensions: can make  $\|f\|_{\mathcal{H}} = \|f\|_{K[T]}$ .

This time. Generalization

1] Bias-variance on train set (deal w/ observation noise)

2] Generalize from train to set

→ Generalization determined by eigenvalues of  $K$  (or of  $S$ )

↳  $\mu_j = n \cdot j^{-\alpha}$  (for  $K$ )  
"power law decay" or  $j^{-\alpha}$  (for  $S$ )

then error is also a power law.

$$\mu_j(S) = \frac{1}{n} \mu_j(K)$$

$\downarrow$   $\frac{1}{n} X^T X$                        $\downarrow$   $XX^T$

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

$$y_i = f^*(x_i) + \epsilon_i \quad \text{for } \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

Kernel ridge regression:

$$\hat{f}_\lambda = K(\lambda I + K)^{-1} y$$

depends on  $\epsilon$

Error:  $\mathbb{E}_\epsilon \left[ \frac{1}{n} \|\hat{f}_\lambda - f^*\|_2^2 \right]$

Bias:  $\frac{1}{n} \left\| \mathbb{E}_\epsilon [\hat{f}_\lambda] - f^* \right\|_2^2$

Vari:  $\frac{1}{n} \sum_{i=1}^n \text{Var}_\epsilon [f(x_i)]$

Bias.  $\mathbb{E}_\epsilon[\hat{f}_\lambda] = K(\lambda I + K)^{-1} \mathbb{E}_\epsilon[y]$

$$(\lambda I + K)(\lambda I + K)^{-1} f = K(\lambda I + K)^{-1} f^*$$

$$\frac{1}{n} \left\| f^* - K(\lambda I + K)^{-1} f^* \right\|_2^2$$

$$= \frac{1}{n} \left\| \lambda (\lambda I + K)^{-1} f^* \right\|_2^2$$

$$= \frac{\lambda^2}{n} \left\| (\lambda I + K)^{-1} f^* \right\|_2^2$$

Further bound.

$$f^T \underbrace{(\lambda I + K)^{-1}}_{\leq \frac{1}{\lambda}} \underbrace{(\lambda I + K)^{-1}}_{\leq K^{-1}} f$$

$$\leq \frac{1}{\lambda} f^T K^{-1} f = \frac{1}{\lambda} \|f\|_K^2$$

$$\frac{\lambda}{n} \|f^*\|_K^2$$

← often use in parametric case, but we'll see that it's too loose in nonparametric case.

Variance.

$$\hat{f}_\lambda = K(\lambda I + K)^{-1} y$$

$f^* + \epsilon$

$$\frac{1}{n} \sum_{i=1}^n \text{Var}_\epsilon[\hat{f}_\lambda(x_i)]$$

$$\frac{1}{n} \mathbb{E}_\epsilon \left[ \left\| K(\lambda I + K)^{-1} \epsilon \right\|_2^2 \right]$$

|

$$L = \frac{1}{n} \mathbb{E}_\varepsilon \left[ \varepsilon^T (\lambda I + K)^{-1} K^2 (\lambda I + K)^{-1} \varepsilon \right]$$

$$= \frac{\sigma^2}{n} \cdot \text{tr} \left( (\lambda I + K)^{-1} K^2 (\lambda I + K)^{-1} \right)$$

Further bounds

$$(\lambda I + K)^{-1} \leq \frac{1}{\lambda} I$$

$$(\lambda I + K)^{-1} \leq K^{-1}$$

$$\frac{\sigma^2}{n\lambda} \text{tr}(K) = \frac{\sigma^2}{\lambda} \text{tr} \left( \frac{K}{n} \right)$$

$$\frac{\sigma^2}{\lambda} \text{tr}(S)$$

$$\approx \frac{\sigma^2 d}{\lambda}$$

often used in parametric case, but too loose here

Naive bound.

$$y = f^* + \varepsilon$$

$$\text{Error} \leq \frac{\lambda}{n} \|f^*\|_{\mathcal{H}}^2 + \frac{\sigma^2}{\lambda} \text{tr} \left( \frac{K}{n} \right)$$

$$\lambda \approx \sqrt{n}$$

$$\lambda = \frac{\sigma}{\|f^*\|_{\mathcal{H}}} \cdot \sqrt{\text{tr}(K)}$$

observable

$$\text{Error} \lesssim \frac{\sigma \|f^*\|_{\mathcal{H}}}{\sqrt{n}} \cdot \sqrt{\text{tr}(K/n)} \approx \frac{1}{\sqrt{n}}$$

Interpretation.

$$\text{Error} = \frac{1}{n} \|\hat{f}_\lambda - f^*\|_2^2$$

Parametric settings  $\frac{1}{n}$

Non-parametric settings  $\frac{1}{\sqrt{n}}$

doesn't depend on geometry of  $\mathcal{H}$

↑  
Something fishy here

---

Issue.  $\|f^*\|_{\mathcal{H}} < \infty$  is a bad assumption

Ridge regression: Bayesian interpretation  
optimal Bayesian inference under Gaussian  
prior on  $\beta$ .

$$\beta \sim N(0, \rho^2 \cdot I)$$

$\Rightarrow$  ridge regression is optimal

$$\mathbb{E}[\|\beta\|_2^2] = \rho^2 \cdot p \quad \leftarrow \text{pretty bad if } p \rightarrow \infty$$

Kernel ridge regression

↙ Gaussian process prior

$$f \sim N(0, \rho^2 \cdot K)$$

$$\|f\|_{\mathcal{H}}^2 = f^T K^{-1} f \quad \mathbb{E}[f^T K^{-1} f] = \rho^2 \cdot \text{tr}(K \cdot K^{-1}) \\ = \rho^2 \cdot n \rightarrow \infty.$$

---

$K =$  Gaussian kernel

$\hookrightarrow$  smooth functions ( $\hat{f}_\lambda$  always smooth)

---

Analysis Strategy.

Assume  $f \sim N(0, \sigma^2 \cdot K)$  and analyze expected error under this prior on  $f$ .

$$\begin{aligned}
 \text{Bias}^2 &= \mathbb{E}_{f \sim N(0, \sigma^2 \cdot K)} \left[ \frac{\lambda^2}{n} \left\| (\lambda I + K)^{-1} f^* \right\|_2^2 \right] \\
 &= \frac{\lambda^2}{n} \text{tr} \left( \mathbb{E}_{f^*} \left[ (\lambda I + K)^{-1} f^* (f^*)^T (\lambda I + K)^{-1} \right] \right) \\
 &= \frac{\lambda^2 \sigma^2}{n} \text{tr} \left( (\lambda I + K)^{-1} K (\lambda I + K)^{-1} \right) \\
 \text{Var} &= \frac{\sigma^2}{n} \cdot \text{tr} \left( (\lambda I + K)^{-1} K^2 (\lambda I + K)^{-1} \right)
 \end{aligned}$$

Let  $\mu_j = j^{\text{th}}$  eigenvalue of  $K$

$$\begin{aligned}
 \text{Error} &= \frac{\lambda \sigma^2}{n} \sum_j \frac{\lambda \mu_j}{(\lambda + \mu_j)^2} + \frac{\sigma^2}{n} \sum_j \frac{\mu_j^2}{(\lambda + \mu_j)^2} \\
 &\approx 1 \text{ if } \lambda = \mu_j \qquad \qquad \qquad \approx 1 \text{ if } \mu_j > \lambda \\
 &\approx \frac{\mu_j}{\lambda} \text{ if } \mu_j < \lambda \qquad \qquad \approx \frac{\mu_j^2}{\lambda^2} \text{ if } \mu_j < \lambda \\
 &\approx \frac{\lambda}{\mu_j} \text{ if } \mu_j \gg \lambda \qquad \leq 1 \qquad \approx \min\left(1, \frac{\mu_j}{\lambda}\right) \\
 &\qquad \qquad \qquad \qquad \qquad \qquad \qquad \approx \min\left(1, \frac{\mu_j^2}{\lambda^2}\right)
 \end{aligned}$$

$$\approx \min\left(1, \frac{\mu_j}{\lambda}\right)$$

$$\text{Error} \leq \left(\frac{\lambda p^2 + \sigma^2}{n}\right) \sum_j \min\left(1, \frac{\mu_j}{\lambda}\right)$$

Define  
 $J = \max\{j \mid \mu_j \geq \lambda\}$

$$= \left(\frac{\lambda p^2 + \sigma^2}{n}\right) \left(\underline{J} + \frac{1}{\lambda} \sum_{j > J} \mu_j\right)$$

Example.

$$\mu_j = j^{-\alpha} \cdot n, \quad \alpha > 1$$

$$J^{-\alpha} \cdot n = \lambda$$

$$J = \left(\frac{\lambda}{n}\right)^{-\frac{1}{\alpha}}$$

$$J = \left(\frac{\lambda}{n}\right)^{-\frac{1}{\alpha}}$$

$$\frac{1}{\lambda} \sum_{j > J} \mu_j \approx \int_J^{\infty} x^{-\alpha} dx$$

$$= \frac{J^{1-\alpha}}{\alpha-1} = \frac{\lambda/n}{\alpha-1} \left(\frac{\lambda}{n}\right)^{-\frac{1}{\alpha} - 1 - \alpha}$$

$$= \frac{1}{\alpha-1} \left(\frac{\lambda}{n}\right)^{-\frac{1}{\alpha}}$$

$$\text{Error} \leq \left(\frac{\lambda p^2 + \sigma^2}{n}\right) \cdot \left(\frac{\alpha}{\alpha-1} \left(\frac{\lambda}{n}\right)^{-\frac{1}{\alpha}}\right)$$

$$\lambda \gg \frac{\sigma^2}{p^2}$$

$$\text{Error}(\lambda^*) = \frac{\sigma^2}{n} \cdot \frac{\alpha}{\alpha-1} \left( \frac{\sigma^2}{np^2} \right)^{-\alpha}$$

$$= \frac{\alpha}{\alpha-1} n^{\frac{1}{\alpha}-1} (\sigma^2)^{1-\frac{1}{\alpha}} (p^2)^{\frac{1}{\alpha}}$$

$\rightarrow n^{\frac{1}{\alpha}-1}$  rate of decay

- $\alpha \rightarrow 1$ : very poor error
- $\alpha = 2$ :  $1/\sqrt{n}$
- $\alpha \rightarrow \infty$ : approach  $1/n$  error from parametric setting

## 2) Generalizing from train test

$$\text{Train: } \frac{1}{n} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2$$

$$\text{Test: } \mathbb{E}_{x \sim p} [(f(x) - f^*(x))^2]$$

Return to parameter view

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

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

Ridge regression:  $\hat{\beta}_\lambda = \left( \frac{\lambda}{n} I + S \right)^{-1} \left( \sum_{i=1}^n x_i y_i \right)$

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

$$= (\lambda_n I + S)^{-1} \left( S \beta^* + \sum_{i=1}^n x_i \varepsilon_i \right)$$



Error:  $(\hat{\beta}_\lambda - \beta^*)^T S (\hat{\beta}_\lambda - \beta^*)$

$$\mathbb{E}[\hat{\beta}_\lambda] = (\lambda_n I + S)^{-1} S \beta^*$$

Bias<sup>2</sup>:  $\delta = \mathbb{E}[\hat{\beta}_\lambda] - \beta^*$ ,  $\delta = \lambda_n (\lambda_n I + S)^{-1} \beta^*$

$$\delta^T S^* \delta = \lambda_n^2 (\beta^*)^T (\lambda_n I + S)^{-1} S^* (\lambda_n I + S)^{-1} \beta^*$$

Under Gaussian priori

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

Key observation: Almost the same as bias if we ran ridge regression on test dist<sup>n</sup>.

$$\text{trace} \left( (\lambda_n I + S^*)^{-1} S^* (\lambda_n I + S^*)^{-1} \right)$$

In particular, as long as

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

test error  $\leq$  C. train error

Proposition

Assume

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

Then

$$\mathbb{E}[\text{Error}] \leq 4 \left( \lambda_n \rho^2 + \frac{\sigma^2}{n} \right) \sum_j \min \left( 1, \frac{\mu_j^*}{\lambda_n} \right)$$

$\frac{\mu_j^*}{\lambda_n}$

whose  $\mu_j^*$  is  $j^{\text{th}}$  largest eigenvalue of  $S^*$ .

$n$  times smaller  
b/c  $S$  vs  $K$

Take-away As long as we pick  $\lambda_n$  s.t.  
 $S^* \leq 2S + \lambda_n I$  holds,  
we get good generalization.

TODO: Show when this is the case.