

9 - KRR and Non-monotonic Risk Curves

- We will talk about generalization for kernel methods
- Give linear regression as an example and show 'double descent'

* Kernel Ridge Regression (KRR):

- Observe n i.i.d. samples $(x_i, y_i) \sim p(x, y)$

$$\hat{f} = \underset{\mathcal{F}}{\operatorname{arg\,min}} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{F}}^2 := \hat{R}(f) \right\}$$

↪ RHS

- By the "Representer Theorem" (last week)

$$\hat{f} = \sum_{i=1}^n \alpha_i k(x_i, \cdot) \quad \text{where } \alpha = (K + n\lambda I)^{-1} y$$

→ Model: $y_i = f_\star(x_i) + \varepsilon_i$ where $\varepsilon_i \perp\!\!\!\perp x_i$, $\mathbb{E}\varepsilon_i = 0$, $\mathbb{E}\varepsilon_i^2 = \sigma^2$.
sets $p(y|x)$ but nothing on $p(x)$ yet.

$$\Rightarrow \hat{R}(f) = \frac{1}{n} \sum_i y_i^2 + \frac{1}{n} \sum_i f(x_i)^2 - \frac{2}{n} \sum_i y_i f(x_i) + \lambda \|f\|_{\mathcal{F}}^2$$

Recall: Representer $\Psi_x = k(x, \cdot)$ $\langle f, \Psi_x \rangle = f(x)$

$$= \frac{1}{n} \sum_i y_i^2 + \frac{1}{n} \sum_i \langle f, \Psi_x \rangle_{\mathcal{F}}^2 - \frac{2}{n} \sum_i y_i \langle f, \Psi_x \rangle_{\mathcal{F}} + \lambda \|f\|_{\mathcal{F}}^2$$

Define: $\hat{\Sigma} = \frac{1}{n} \sum_i \Psi_{x_i} \otimes \Psi_{x_i}$ → self-adjoint operator

$$= \frac{1}{n} \sum_i y_i^2 + \langle f, \hat{\Sigma} f \rangle - 2 \underbrace{\left\langle \frac{1}{n} \sum_i y_i \Psi_{x_i}, f \right\rangle}_{:= b} + \lambda \langle f, f \rangle$$

$$= \frac{1}{n} \sum_i y_i^2 + \langle f, \hat{\Sigma} f \rangle - 2 \langle b, f \rangle + \lambda \langle f, f \rangle \quad (\text{quadratic in } f)$$

→ minimized at $\hat{f} = (\hat{\Sigma} + \lambda I)^{-1} b$.

- We are interested in expected excess risk: $\mathbb{E} \left[\frac{\|\hat{f} - f^*\|_{L^2(\rho)}^2}{\rho(x)} \right]$

$$\mathbb{E} \left[\|\hat{f} - f^*\|_{L^2(\rho)}^2 \right] = \mathbb{E} \left[\left\| (\hat{\Sigma} + \lambda I)^{-1} \frac{1}{n} \sum_i \gamma_i \Psi_{x_i} - f^* \right\|_{L^2}^2 \right]$$

$\downarrow = f^*(x_i) + \varepsilon_i$
 $= \langle \Psi_{x_i}, f^* \rangle + \varepsilon_i$ (!!) can we do this?

$$= \mathbb{E} \left[\left\| (\hat{\Sigma} + \lambda I)^{-1} \frac{1}{n} \sum_i \underbrace{\Psi_{x_i} \{ \langle \Psi_{x_i}, f^* \rangle + \varepsilon_i \}}_{\downarrow \perp \varepsilon_i} - f^* \right\|_{L^2(\rho)}^2 \right]$$

$$= \mathbb{E} \left[\left\| (\hat{\Sigma} + \lambda I)^{-1} \frac{1}{n} \sum_i \varepsilon_i \Psi_{x_i} \right\|_{L^2(\rho)}^2 \right] + \mathbb{E} \left[\left\| (\hat{\Sigma} + \lambda I)^{-1} \frac{1}{n} \sum_i \Psi_{x_i} \langle \Psi_{x_i}, f^* \rangle - f^* \right\|_{L^2(\rho)}^2 \right]$$

Variance $\triangleq V(\lambda)$

Bias $\triangleq B(\lambda)$

Let $\Sigma = \mathbb{E} \hat{\Sigma}$ (or $\mathbb{E} [\Psi_x \otimes \Psi_x]$) and observe

$$\begin{aligned} \|g\|_{L^2(\rho)}^2 &= \int g(x)^2 d\rho(x) = \int \langle g, \Psi_x \rangle^2 d\rho(x) = \left\langle g, \int \Psi_x \otimes \Psi_x d\rho(x) g \right\rangle \\ &= \langle g, \Sigma g \rangle. \end{aligned}$$

$$V(\lambda) \stackrel{?}{=} \mathbb{E} \left[\left\langle (\hat{\Sigma} + \lambda I)^{-1} \frac{1}{n} \sum_i \varepsilon_i \Psi_{x_i}, \sum (\hat{\Sigma} + \lambda I)^{-1} \frac{1}{n} \sum_i \varepsilon_i \Psi_{x_i} \right\rangle \right]$$

$$= \frac{1}{n^2} \mathbb{E} \left[\text{Tr} \left((\hat{\Sigma} + \lambda I)^{-1} \sum (\hat{\Sigma} + \lambda I)^{-1} \underbrace{\sum_i \varepsilon_i^2 \Psi_{x_i} \otimes \Psi_{x_i}}_{\mathbb{E}_{\varepsilon} \rightarrow \hat{\Sigma} \cdot n \cdot \sigma^2} \right) \right]$$

$$= \frac{\sigma^2}{n} \mathbb{E} \left[\text{Tr} \left((\hat{\Sigma} + \lambda I)^{-1} \underbrace{\sum (\hat{\Sigma} + \lambda I)^{-1} \sum}_{\leq I} \right) \right] \quad \text{since } (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} \leq I$$

$$\leq \frac{\sigma^2}{n} \mathbb{E} \left[\text{Tr} \left((\hat{\Sigma} + \lambda I)^{-1} \sum \right) \right]$$

To simplify calculations, we assume $f^* \in \mathcal{F} \Rightarrow \langle f^*, \psi_x \rangle = f^*(x)$

$$\begin{aligned}
 \mathbb{B}(\lambda) &:= \mathbb{E} \left[\left\| (\hat{\Sigma} + \lambda I)^{-1} \frac{1}{n} \sum_i \psi_{x_i} \langle \psi_{x_i}, f^* \rangle - f^* \right\|_{L^2(p)}^2 \right] \\
 &= \mathbb{E} \left[\left\| (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} f^* - f^* \right\|_{L^2(p)}^2 \right] \\
 \text{by } ! &= \mathbb{E} \left[\left\langle \left\{ (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} - I \right\} f^*, \sum \left\{ (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} - I \right\} f^* \right\rangle_{\mathcal{F}} \right] \\
 &\quad (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} - I = (\hat{\Sigma} + \lambda I)^{-1} (\hat{\Sigma} + \lambda I - \lambda I) - I \\
 &\quad = I - \lambda (\hat{\Sigma} + \lambda I)^{-1} - I \\
 &= \mathbb{E} \left[\left\| \lambda \hat{\Sigma}^{\frac{1}{2}} (\hat{\Sigma} + \lambda I)^{-1} f^* \right\|_{\mathcal{F}}^2 \right]
 \end{aligned}$$

Thus: Excess Risk $= V(\lambda) + \mathbb{B}(\lambda)$

$$\leq \frac{\sigma^2}{n} \mathbb{E} \left[\text{Tr} \left((\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} \right) \right] + \lambda^2 \mathbb{E} \left[\left\| \hat{\Sigma}^{\frac{1}{2}} (\hat{\Sigma} + \lambda I)^{-1} f^* \right\|_{\mathcal{F}}^2 \right]$$

- $\hat{\Sigma}$ concentrates around Σ (Bach p185) for constant d. Need $\|\psi_x\|_{\mathcal{F}} \leq R$

- First term $\approx \frac{\sigma^2}{n} \text{Tr} \left((\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} \right) = O \left(\frac{\sigma^2}{n\lambda} \right)$

- Second term $\approx \lambda^2 \left\langle f^*, \underbrace{(\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma}}_{\leq I} \underbrace{(\hat{\Sigma} + \lambda I)^{-1} f^*}_{\leq \frac{1}{\lambda}} \right\rangle = O \left(\lambda \|f^*\|_{\mathcal{F}}^2 \right)$

$$\Rightarrow \text{Excess Risk} \approx \frac{\sigma^2}{n\lambda} + \lambda \|f^*\|_{\mathcal{F}}^2$$

- choosing $\lambda = \frac{1}{\sqrt{n}}$ $\approx \frac{1}{\sqrt{n}}$ \Rightarrow generalization

Theorem: Let $y_i = f^*(x_i) + \varepsilon_i$ for $i=1 \dots n$ for $f^* \in \mathcal{F}$ and $\hat{f} = \underset{\mathcal{F}}{\operatorname{arg\,min}} \hat{R}(f)$ for $\lambda = \frac{1}{\sqrt{n}}$. Then, if $\|\psi_x\|_{\mathcal{F}} \leq R \forall x$, we have

$$\mathbb{E} \left[\|\hat{f} - f^*\|_{L^2(p)}^2 \right] \lesssim \frac{1}{n}$$

* Linear regression: We consider linear funcs:

- RKHS: $\mathcal{F} = \left\{ f_\theta(x) = \langle \theta, x \rangle : \theta \in \mathbb{R}^d \right\}$
- $\langle \cdot, \cdot \rangle_{\mathcal{F}}$: $\langle f_\theta, f_\omega \rangle_{\mathcal{F}} = \langle \theta, \omega \rangle_{\mathbb{R}^d}$
- Representer: $\langle \Psi_x, f_\theta \rangle_{\mathcal{F}} = f_\theta(x) = \langle \theta, x \rangle$
 $\Psi_x(y) = \langle x, y \rangle$ (or $\Psi_x = f_x$)
- Model: $y = \langle \theta_*, x \rangle + \varepsilon$ $\varepsilon \sim N(0, \sigma^2)$
- $x \sim N(0, I)$ } $E \langle \theta_*, x \rangle^2 = 1$
 $\theta_* \sim N(0, \frac{1}{n} I)$ }
- $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \theta, x_i \rangle)^2 + \frac{\lambda}{2} \|\theta\|_2^2$
 $= \left(\frac{1}{n} X^\top X + \lambda I \right)^{-1} \frac{1}{n} X^\top y$
- Excess Risk: $E \left[\|\hat{\theta} - \theta_*\|^2 \right] \triangleq ER(\lambda)$

$$ER(\lambda) = B(\lambda) + V(\lambda) \quad \boxed{\text{where}}$$

$$\begin{aligned} B(\lambda) &= \lambda^2 E \left[\langle \theta^*, (\hat{\Sigma} + \lambda I)^{-2} \theta_* \rangle \right] & V(\lambda) &= \frac{\sigma^2}{n} E \left[\text{Tr} \left((\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \right) \right] \\ &= \frac{\lambda^2}{d} E \left[\text{Tr} \left((\hat{\Sigma} + \lambda I)^{-2} \right) \right] & &= \sigma^2 \frac{1}{n} E \left[\frac{1}{d} \sum_{i=1}^d \frac{\lambda_i}{(\lambda_i + \lambda)^2} \right] \\ &= \lambda^2 E \left[\frac{1}{d} \sum_{i=1}^d \frac{1}{(\lambda_i + \lambda)^2} \right] \end{aligned}$$

where λ_i 's are the eigenvalues of $\hat{\Sigma}$.

— Marchenko-Pastur Law: Let $d, n \rightarrow \infty$ and $\frac{d}{n} \rightarrow \gamma$.

Let $X \in \mathbb{R}^{n \times d}$ st. X_{ij} are iid mean, variance 1.

Then, for any reasonable fnc ϕ and $\hat{\Sigma} = \frac{1}{n} X^T X$

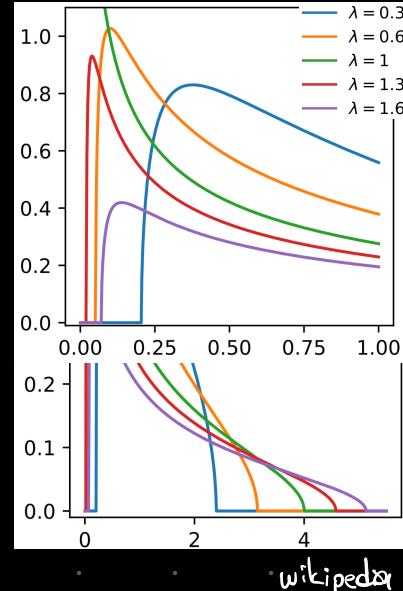
$$\frac{1}{d} \sum_{j=1}^d \phi(\lambda_j(\hat{\Sigma})) \xrightarrow{\text{a.s.}} \int \phi \, d\mu$$

where μ is the M-P law given as

$$\frac{d\mu}{dx} = \begin{cases} (1 - \gamma^{-1}) \delta_0(x) + \nu(x) & \text{if } \gamma > 1 \\ \nu(x) & \text{if } \gamma \in [0, 1] \end{cases}$$

and $\nu(x) = \begin{cases} \frac{1}{2\pi} \frac{\sqrt{(\gamma_+ - x)(x - \gamma_-)}}{\gamma x} & x \in [\gamma_-, \gamma_+] \\ 0 & \text{else} \end{cases}$

with $\gamma_{\pm} = (1 \pm \sqrt{\gamma})^2$



wikipedia

* Stieltjes transform:

$$s(z) = \int \frac{1}{x-z} \, d\mu(x)$$

of M-P law:

$$s(-z) = \frac{-1 + \gamma - z + \sqrt{(1-\gamma+z)^2 + 4\gamma z}}{2\gamma z} \quad \text{for } z > 0.$$

* For linear regression: $ER(\lambda) = V(\lambda) + B(\lambda)$

\checkmark : $V(\lambda) = \sigma^2 \frac{1}{n} \mathbb{E} \left[\frac{1}{d} \sum_{i=1}^d \frac{\lambda_i}{(\lambda_i + \lambda)^2} \right]$

$$\rightarrow \sigma^2 \gamma \int \frac{x}{(\lambda+x)^2} \, d\mu(x) = \sigma^2 \gamma \left\{ \int \frac{1}{x+\lambda} \, d\mu(x) - \int \frac{\lambda}{(\lambda+x)^2} \, d\mu(x) \right\}$$

$$= \sigma^2 \gamma \left\{ s(-\lambda) - \lambda s'(-\lambda) \right\}.$$

B : $B(\lambda) \rightarrow \lambda^2 \int \frac{1}{(\lambda+x)^2} \, d\mu(x) = \lambda^2 s'(-\lambda)$

Theorem: Let $\gamma_i = \langle \theta_*, x_i \rangle + \epsilon_i$ for $x \sim N(0, I)$ and $\epsilon \sim N(0, \sigma^2)$

and $\theta_* \sim N(0, \frac{1}{n} I)$. Then, the ridge regression solution

$$\hat{\theta}_\lambda = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (\gamma_i - \langle \theta, x_i \rangle)^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

satisfies $E[\|\hat{\theta}_\lambda - \theta_*\|^2] \triangleq ER(\lambda) = B(\lambda) + V(\lambda)$ where as $\frac{d}{n} \rightarrow \infty$

$$B(\lambda) \rightarrow \lambda^2 s(-\lambda)$$

$V(\lambda) \rightarrow \sigma^2 \gamma \{s(-\lambda) - \lambda s'(-\lambda)\}$ almost surely.

Remarks:

- "Ridgeless" case ($\lambda \downarrow 0$) \Rightarrow Gradient descent can find this!
Minimum norm solution when $d > n$.
Implicit regularization (A3)

$$\lim_{\lambda \downarrow 0} B(\lambda) = B(0_+) = \lim_{\lambda \downarrow 0} \lambda^2 s(-\lambda) = \begin{cases} 0 & \gamma < 1 \\ 1 - \frac{1}{\gamma} & \gamma \geq 1 \end{cases} \quad (A3)$$

$$\lim_{\lambda \downarrow 0} V(\lambda) = V(0_+) = \lim_{\lambda \downarrow 0} \sigma^2 \gamma \{s(-\lambda) - \lambda s'(-\lambda)\} = \sigma^2 \begin{cases} \frac{\gamma}{1-\gamma} & \gamma < 1 \\ \frac{1}{\gamma-1} & \gamma \geq 1 \end{cases}$$

* Variance diverges
only when $\lambda = 0_+$.

