

Typical examples:

1/ $(\varphi_j)_{j \geq 1}$ Fourier basis, $l(y, f(x)) = (y - f(x))^2$

$$\mathcal{F}_N = \{f = \sum_{j=1}^N \theta_j \varphi_j, \theta_j \in \mathbb{R}\}, \quad d(f, f^*) = \|f - f^*\|_{L^2}^2$$

$$y_i = f^*(x_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2).$$

• Approximation error: $f^* = \sum_{j \geq 1} \theta_j^* \varphi_j$

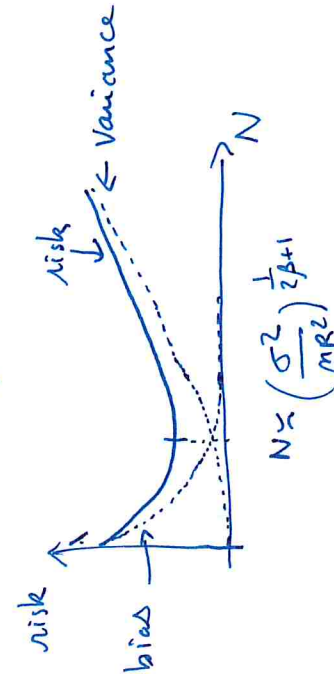
$$\min_{f \in \mathcal{F}_N} \|f - f^*\|_{L^2}^2 = \sum_{j > N+1} \theta_j^{*2} \quad \mathbb{E}[d(f, f^*)] \leq \sum_{j > N+1} \sigma_j^2 + \frac{N}{M} \sigma^2.$$

• Stochastic error: $\frac{N}{M} \sigma^2$

case when $f^* \in W_{\beta}^{p,q}(\mathbb{R}^d) = \{f: \sum_{j \geq 1} j^{2\beta} \theta_j^2 \leq R^2\}$.

$$\rightarrow \sum_{j > N+1} \theta_j^2 \leq \frac{1}{N^{2\beta}} R^2$$

$$\Rightarrow \text{risk} \leq \frac{R^2}{N^{2\beta}} + \frac{N}{M} \sigma^2 \quad \text{"optimal" } N \Rightarrow N \asymp \left(\frac{\sigma^2}{MR^2}\right)^{\frac{1}{2\beta+1}}$$



Surprises in High-Dimensional Ridgeless Least Squares Interpolation.

Hastie, Pantanari, Rosset, Tibshirani.

1 Classical approximation / variance trade-off.

Classical learning problem: $(x_i, y_i) \stackrel{iid}{\sim} \mathbb{D}$,

learn $f^*(x) = \mathbb{E}[Y | X=x]$:

$$\mathcal{F}_N \in \operatorname{argmin}_{f \in \mathcal{F}_N} \frac{1}{M} \sum_{i=1}^M l(y_i, f(x_i))$$

Risk: approximation / variance trade-off

$$\mathbb{E}[d(f, f^*)] = \underbrace{\min_{f \in \mathcal{F}_N} d(f, f^*)}_{\text{approximation error}} + \underbrace{\mathbb{E}[d(f, f^*) - \min_{f \in \mathcal{F}_N} d(f, f^*)]}_{\text{stochastic error}}$$

when $d(f, f^*) = \|f - f^*\|_H^2 \leftarrow$ Hilbert norm.

we have the bias/variance decomposition

$$\mathbb{E}[\|f - f^*\|_H^2] = \underbrace{\mathbb{E}[\|f - \mathbb{E}[f]\|_H^2]}_{\text{Variance term}} + \underbrace{\|\mathbb{E}[f] - f^*\|_H^2}_{\text{bias term}}$$

2/ $\mathcal{F}_N = \text{ball } B_{\mathcal{H}}(0, N)$ in a RKHS.

remark: in a RKHS $\|f\|_{\mathcal{H}} \Leftrightarrow$ regularity as

$$|f(x) - f(y)| = |\langle f, k(x, \cdot) - k(y, \cdot) \rangle_{\mathcal{H}}|$$

$$\leq \|f\|_{\mathcal{H}} \|k(x, \cdot) - k(y, \cdot)\|_{\mathcal{H}}.$$

Lagrangian version

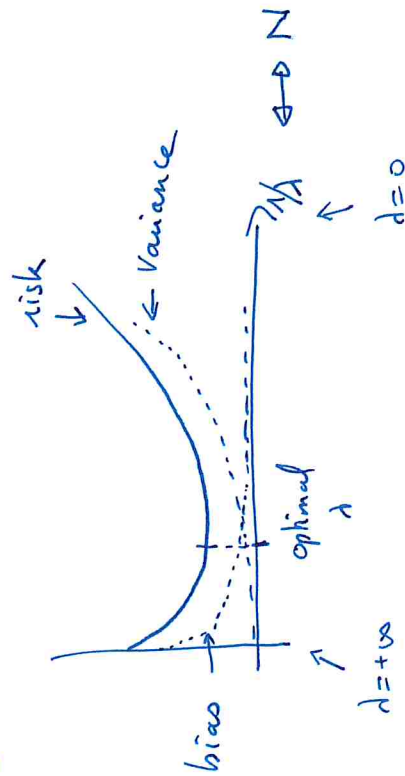
$$\hat{f} \in \underset{f \in \mathcal{H}}{\text{argmin}} \frac{1}{M} \sum_{i=1}^M \ell(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2$$

solution: $\hat{f} = \sum_{j=1}^M \hat{\beta}_j k(x_j, \cdot)$ where

$$\hat{\beta} \in \underset{\beta \in \mathbb{R}^M}{\text{argmin}} \frac{1}{M} \sum_{i=1}^M \ell(y_i, (K\beta)_i) + \lambda \beta^T K \beta$$

where $K = (k(x_i, x_j))_{i,j=1,\dots,M}$.

again we have



Yet: in deep-learning or Random Forests

→ number p of parameters $\gg N = \text{sample size}$

→ no regularisation (interpolation of data points)

↳ work empirically very well!

Why? Something to do with optimisation algorithm?

② Simple least square regression

$$Y \in \mathbb{R}^m = X \beta + \varepsilon, \text{ with } \mathbb{E}[\varepsilon] = 0, \text{cov}(\varepsilon) = \sigma^2 I_m.$$

$$\ell(y, y') = (y - y')^2$$

$$\text{So } \mathcal{L}(\beta) = \frac{1}{M} \sum_{i=1}^M \ell(y_i, x_i^T \beta) = \frac{1}{M} \|Y - X\beta\|^2$$

$$\Rightarrow \hat{\beta} \in \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \|Y - X\beta\|^2.$$

classical case: $p \leq M$ and $\text{rank}(X) = p$ (full rank)

$$\text{Then } \hat{\beta} = (X^T X)^{-1} X^T Y$$

• $\mathbb{E}[\hat{\beta}] = \beta^*$ (no bias)

$$\cdot \text{cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

Hence

$$\mathbb{E}[\|\hat{\beta} - \beta\|^2] = \|\hat{\beta} - \mathbb{E}[\hat{\beta}]\|^2 + \text{tr cov}(\hat{\beta}) = \sigma^2 \text{tr} (X^T X)^{-1}$$

Typically diverge when $p \rightarrow M$.

Interpolation case: $p \geq m$, $\text{rank}(X) = m$.

$\rightarrow \min_{\beta} \|Y - X\beta\|^2 = 0$

\rightarrow no unique solution to $Y = X\beta$

Lemma: $\mathcal{Y} = \{Y : Y = X\beta\}$ is equal to $\hat{\beta}_0 + \text{ker}(X)$
 where $\hat{\beta}_0 = X^T(XX^T)^{-1}Y$

$X\hat{\beta}_0 = Y$ and $X(\beta - \hat{\beta}_0) = 0 \quad \forall \beta \in \mathcal{Y}$

Remark: $X: \mathbb{R}^p = \underbrace{\text{ker}(X)}_{\dim=p-m} \oplus \underbrace{\text{Im}(X^T)}_{\dim=m} \rightarrow \mathbb{R}^m = \text{Im}(X)$

So $\hat{\beta}_0 \perp \text{ker}(X)$.

question: if we minimise $L(\beta) = \frac{1}{2} \|Y - X\beta\|^2$ by gradient descent started at 0, where do we converge?

GD: $\beta^{(k+1)} = \beta^{(k)} - \frac{\lambda D}{m} X^T(X\beta^{(k)} - Y)$ as $\nabla L(\beta) = \frac{2}{m} X^T(X\beta - Y)$
 if $\beta^{(0)} = 0$ then $\beta^{(k)} \in \text{Im}(X^T) \quad \forall k \geq 0$.

So $\beta^{(k)} \xrightarrow[k \rightarrow \infty]{} \hat{\beta}_0!$

Message: Gradient Descent select a specific solution.

Lemma: if $X \stackrel{\text{SVD}}{=} \sum_{k=1}^m \sigma_k U_k V_k^T$

$\hat{\beta}_0 = \sum_{k=1}^m \frac{1}{\sigma_k} V_k U_k^T Y$

$= \underset{\beta \in \mathcal{Y}}{\text{argmin}} \| \beta \|^2$

$= \lim_{\lambda \gg 0} \underbrace{(X^T X + \lambda I)^{-1} X^T Y}_{\text{argmin}_{\beta} \{ \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \}}$

Proof: • for $\hat{\beta} \in \mathcal{Y}$, $\hat{\beta} = \hat{\beta}_0 + v \in \text{ker}(X)$

as $\hat{\beta}_0 \perp \text{ker}(X)$: $\|\hat{\beta}\|^2 = \|\hat{\beta}_0\|^2 + \|v\|^2 \geq \|\hat{\beta}_0\|^2$

so $\hat{\beta}_0 \in \underset{\beta \in \mathcal{Y}}{\text{argmin}} \|\beta\|^2$.

• $X^T(XX^T)^{-1} = \sum_k \sigma_k V_k U_k^T \times \sum_k \frac{1}{\sigma_k^2} U_k U_k^T$

$= \sum_k \frac{1}{\sigma_k} V_k U_k^T$

• $(X^T X + \lambda I)^{-1} X^T = \left(\sum_{k=1}^m \frac{1}{\sigma_k^2 + \lambda} V_k V_k^T + \frac{1}{\lambda} P^\perp \right) \sum_{k=1}^m \sigma_k V_k U_k^T$

$= \sum_{k=1}^m \frac{\sigma_k}{\sigma_k^2 + \lambda} V_k U_k^T$

$\xrightarrow{\lambda \gg 0} \sum_{k=1}^m \frac{1}{\sigma_k} V_k U_k^T$



Paper: analyse $\hat{\beta}_0$ in different settings:

linear: $X_i = \Sigma^{1/2} Z_i$, Z_i iid $\mathbb{E}[Z_i] = 0$
 $\text{cov}(Z_i) = I_p \in \mathbb{R}^p$

non-linear: $X_i = \varphi(WZ_i)$, Z_i iid $\mathcal{N}(0, I_d)$
 $Z = \begin{pmatrix} x \\ \vdots \\ x \end{pmatrix} \begin{matrix} \xrightarrow{w} \\ \xrightarrow{y} \\ \xrightarrow{y} \\ \xrightarrow{y} \end{matrix} X$
 W with $w_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \frac{1}{d})$

mis-specified: some explanatory variables are missing.

asymptotic setting: $p \ll m$ with $m \rightarrow \infty$
 and $d \ll p$

Focus on: sections 2, 3 and 5 (linear setting).

risk: $R_X(\hat{\beta}, \beta) = \mathbb{E}[(x_0^T \hat{\beta} - x_0^T \beta)^2 | X]$ where $x_0 \perp X, Y$.
 $= \mathbb{E}[\| \hat{\beta} - \beta \|^2 | X]$ where $\| \beta \|^2 = \beta^T Z \beta$.

Lemma: $\mathbb{E}[\|Z\|_\Sigma^2] = \| \mathbb{E}[Z] \|^2_\Sigma + \text{Tr}(\Sigma \text{cov}(Z))$

$$\begin{aligned} \# \quad Z^T \Sigma Z &= \mathbb{E}[Z]^T \Sigma \mathbb{E}[Z] + \underbrace{(Z - \mathbb{E}[Z])^T \Sigma (Z - \mathbb{E}[Z])}_{= \text{cov}(Z)} \\ &= \langle \Sigma, (Z - \mathbb{E}[Z])(Z - \mathbb{E}[Z])^T \rangle \\ &+ 2 \underbrace{\mathbb{E}[Z]^T \Sigma (Z - \mathbb{E}[Z])}_{\mathbb{E}[\dots]} = 0 \end{aligned}$$

#

0.54
4

$$\text{So } R_X(\hat{\beta}, \beta) = \underbrace{\| \mathbb{E}[\hat{\beta}] - \beta \|^2_\Sigma}_{B_X(\hat{\beta}, \beta)} + \underbrace{\text{Tr}(\Sigma \text{cov}(\hat{\beta} | X))}_{V_X(\hat{\beta}, \beta)}$$

③ Linear case (sections 2 and 3)

X is full rank: $\text{rank}(X) = \min(m, p)$.

classical regime: $m \gg p$ i.e. $\delta \ll 1$

$$R_X(\hat{\beta}, \beta) = 0 + \sigma^2 \text{Tr}(\Sigma (X^T X)^{-1}) \xrightarrow{\text{Thm 1}} \sigma^2 \frac{\delta}{1-\delta}$$

interpolation regime: $p \geq m$ i.e. $\delta \geq 1$

- $\mathbb{E}[\hat{\beta}_0 | X] = \mathbb{E}[X^T (X X^T)^{-1} Y | X]$
 $Y = X\beta + \varepsilon \quad \underbrace{X^T (X X^T)^{-1} X \beta}_{P_{\text{Im}(X^T)}} + 0$

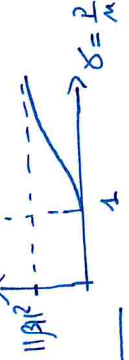
So $\beta - \mathbb{E}[\hat{\beta}_0 | X] = \underbrace{(I - P_{\text{Im}(X^T)}) \beta}_{= P_{\text{ker}(X)}}$

$$\begin{aligned} | B_X(\hat{\beta}, \beta) &= \mathbb{E}[\| P_{\text{ker}(X)} \beta \|^2_\Sigma] \\ \cdot \text{cov}(\hat{\beta}_0 | X) &= \sum_{k=1}^m \frac{1}{\sigma_k^2} v_k v_k^T \text{cov}(Y | X) \left(\sum_{k=1}^m \frac{1}{\sigma_k^2} v_k v_k^T \right)^T \\ &= \sigma^2 \sum_{k=1}^m \frac{1}{\sigma_k^2} v_k v_k^T \end{aligned}$$

$$V_x(\beta, \beta) = \text{Tr}(\Sigma \text{cov}(\beta_0 | X))$$

$$= \sigma^2 \sum_{k=1}^m \frac{1}{\sigma_k^2} V_k^T \Sigma V_k$$

Lemma 2: $\Sigma = I_p, \delta > 1$ increase with δ !!



$$B_x(\beta, \beta) \rightarrow \|\beta\|^2 (1 - \frac{1}{\delta})$$

idea: $\mu := \frac{\beta}{\|\beta\|}$

$$B_x(\beta, \beta) = \|\text{Proj}_{\text{ker}(X)} \beta\|^2 = \|\beta\|^2 \|\text{Proj}_{\text{ker}(X)} \mu\|^2$$

$$\text{ker}(X) = \text{Span}\langle w_1, \dots, w_{p-m} \rangle$$

↑ random vectors (unit norm)

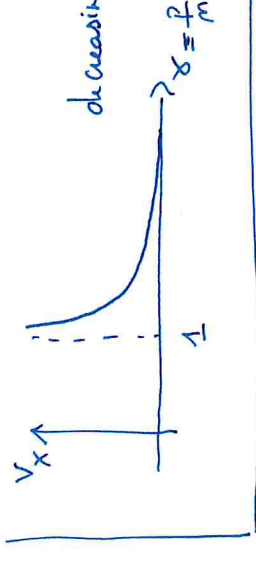
$$\|\text{Proj}_{\text{ker}(X)} \mu\|^2 = \sum_{k=1}^{p-m} \langle w_k, \mu \rangle^2$$

we have $\mathbb{E}[\langle w_k, \mu \rangle^2] = \frac{1}{p}$ as $\sum_{k=1}^{p-m} \langle w_k, \mu \rangle^2 = 1$ \uparrow $\|\mu\|^2 = 1$

so $\|\text{Proj}_{\text{ker}(X)} \mu\|^2 \approx \frac{p-m}{p} = 1 - \frac{1}{\delta}$

#

Lemma 3: $V_x(\beta, \beta) = \sigma^2 \sum_{k=1}^{p-m} \frac{1}{\sigma_k^2} \rightarrow \frac{\sigma^2}{\delta-1}$



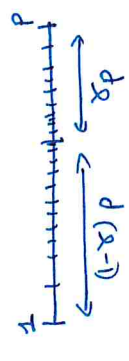
decreasing with δ !!

idea: Marcenko-Pastur:

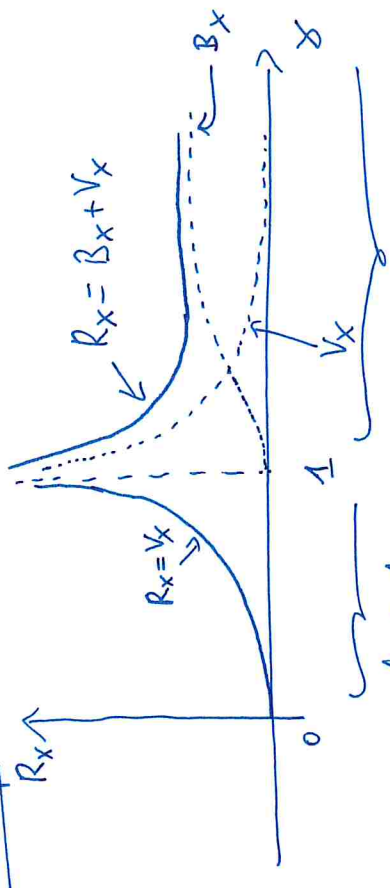
$$\frac{1}{p} \sum_{k=1}^p \delta \frac{\sigma_k^2}{\sigma_k^2/p} \rightarrow \mu_\delta(dx) = \frac{\sqrt{((1+\sqrt{\delta})^2 - x)(x - (1-\sqrt{\delta})^2)}}{2\pi x} dx$$

so $\sum_{k=1}^{p-m} \frac{1}{\sigma_k^2} = \frac{1}{p} \sum_{k=1}^{p-m} \frac{1}{\sigma_k^2/p} \rightarrow \int_{Q_{1-\delta}}^{+\infty} \frac{1}{x} \mu_\delta(dx) = \frac{1}{\delta-1}$

where $F_{\mu_\delta}(Q_{1-\delta}) = 1-\delta$



Sum-up



classical regime

interpolation regime

#

④ Dis-specified case (section 5)

observe (x_i, y_i) iid with

$$y_i = x_i^T \beta + w_i^T \theta + \epsilon_i$$

unobserved, iid.

nisk: compare $x_0^T \beta$ to $E[y_0 | x_0, w_0] = x_0^T \beta + w_0^T \theta$.

$$R_x(\hat{\beta}, \beta) := E[(x_0^T \hat{\beta} - E[y_0 | x_0, w_0])^2 | X]$$

$$\stackrel{\text{Pythagore}}{=} E[(x_0^T \hat{\beta} - E[y_0 | x_0])^2 | X] +$$

$$E[\underbrace{(E[y_0 | x_0] - E[y_0 | x_0, w_0])^2}_{\text{approximation bias}} | X]$$

approximation bias

$$= \sigma_\theta^2$$

σ_θ^2 complex in general.

When all entries of x_i, w_i are iid:

$$E[y_0 | x_0] = x_0^T \beta \quad \text{and} \quad E[y_0 | x_0, w_0] = x_0^T \beta + w_0^T \theta$$

$$\Rightarrow \sigma_\theta^2 = E[(w_0^T \theta)^2 | X] = \| \theta \|^2$$

$$= r^2 (1 - r)$$

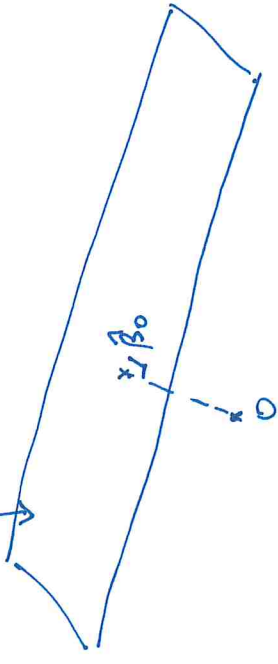
$$\uparrow \frac{\| \theta \|^2 + \| \beta \|^2}{r^2}$$

fraction of signal explained by x_0 .

signal strength

intuition?

$$y = \hat{\beta}_0 + \text{ker}(X) \quad \dim(y) = p - m = p(1 - \frac{1}{r})$$



when $p \rightarrow$, y becomes larger, so more

solutions to $y = X\beta$

$$\Rightarrow \| \hat{\beta}_0 \|^2 \gg 0$$

\Rightarrow $\left\{ \begin{array}{l} - \text{bias} \rightarrow \\ - \text{variance} \rightarrow \end{array} \right.$

Yet: this simple setting miss the approximation / variance trade-off: best at $r \equiv 0$

\leadsto misspecified setting.

more complex behaviours

→ formula theorem 4

→ case $1 - K(\delta) = (1 + \delta)^{-\alpha}$ (→ Fourier)

Figure 3 and 4.

Sometimes best R_x for $\delta > 1$.

Take Home Message:

- optimisation algorithms select some specific solutions in the interpolation domain
- can lead to opposite behavior of in the

classical regime

→ some limits yet

i) $R_x(\beta_0) \geq \min_{\lambda \geq 0} R_x(\beta_\lambda)$ Ridge.

ii) full optimisation is not performed (usually) in the interpolation domain.

↳ early stopping ↔ ^{implicit} regularisation

iii) Neural Networking much more complex as features are learnt at the same time as β → the behavior might be very different.