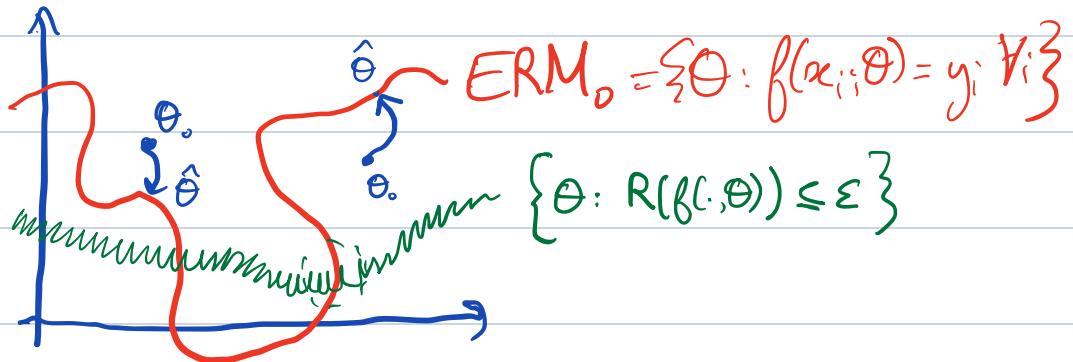


①

Lecture 9:Benign OverfittingLast week:

$\hat{\theta}$: depends on algo, initialization, parametrization etc
(step size etc.)

What about generalization of $f(\cdot; \hat{\theta})$?

→ model trained with no control on the model complexity
(until interpolation of the data)

≠ standard wisdom in statistical learning

→ models that are too complex will not generalize well

→ overfitting is bad and should be avoided

Today: "benign overfitting" → interpolating solution can generalize well

* Classical picture

* min- $\| \cdot \|_2$ interpolating solution in linear models (solution attained by GD)

↳ generalizes via "self-induced regularization"

* Example: inner-product kernel

→ RMT & exact asymptotics
→ matrix concentration

} No time for these
→ will go back to it in lecture 6

(2)

Classical Picture

Parametric family: $\mathcal{F} = \{f(\cdot; \theta) : \theta \in \mathbb{R}^p\}$

Data: $\{(y_i, x_i)\}_{i \leq n}$ $x_i \in \mathbb{R}^d \stackrel{iid}{\sim} P$ $y_i \in \mathbb{R}$

$$y_i = f(x_i; \theta_*) + \varepsilon_i \quad) \text{ noise e.g. } \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

SRM: $\hat{\theta}(\lambda) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2 + \lambda \|\theta\|_2^2 \right\}$

Excess

Test error: $R(\theta_*, \lambda) = E_\alpha \left[(f(x; \theta_*) - f(x; \hat{\theta}(\lambda)))^2 \right]$

Bias - Variance decomposition:

$$E_\varepsilon [R(\theta_*, \lambda)] = E_{\alpha, \varepsilon} \left[(f(x; \theta_*) - f(x; \hat{\theta}(\lambda)))^2 \right]$$

$$= E_\alpha \left[(f(x; \theta_*) - E_\varepsilon [f(x; \hat{\theta}(\lambda))])^2 \right] + E_{\alpha, \varepsilon} \left[(f(x; \hat{\theta}(\lambda)) - E_\varepsilon [f(x; \hat{\theta}(\lambda))])^2 \right]$$

$=: \text{BIAS } (\lambda)$

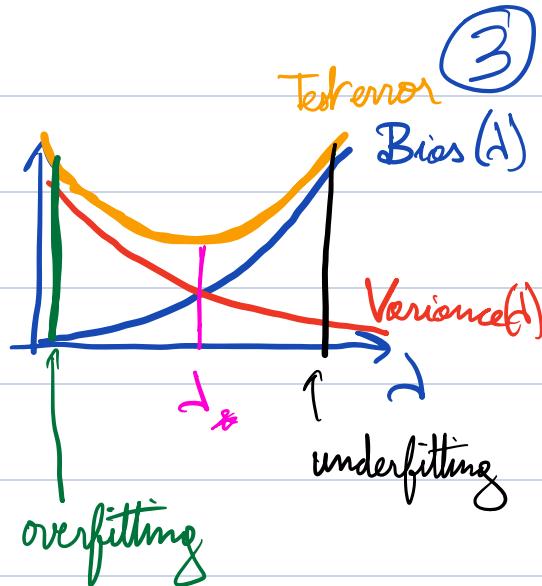
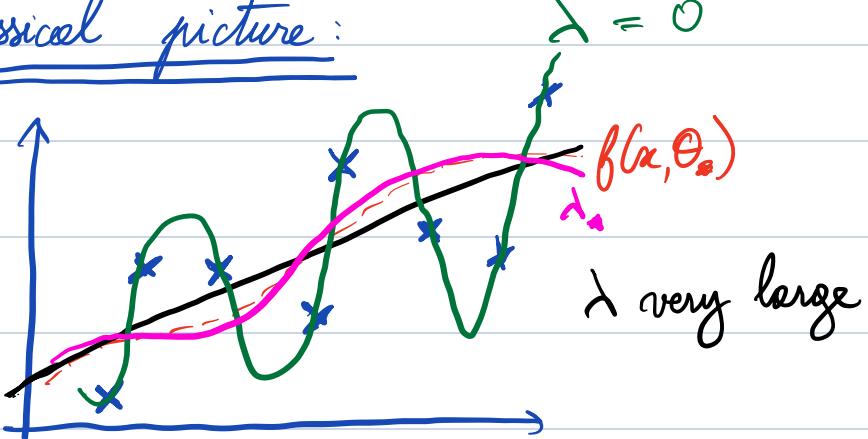
+

VARIANCE (λ)

$$= E_\alpha \left[\text{bias}(E_\varepsilon [f(x; \hat{\theta})])^2 \right]$$

$$= E_\alpha \left[\text{Var}_\varepsilon (f(x; \hat{\theta})) \right]$$

Classical picture:



Trade-off between Bias and Variance

→ need to carefully choose regularization parameter λ in order to control the complexity of the fitted model

$$f(x; \theta(\lambda))$$

⇒ too complex: small bias, large variance

OVERFITTING

⇒ too simple: large bias, small variance

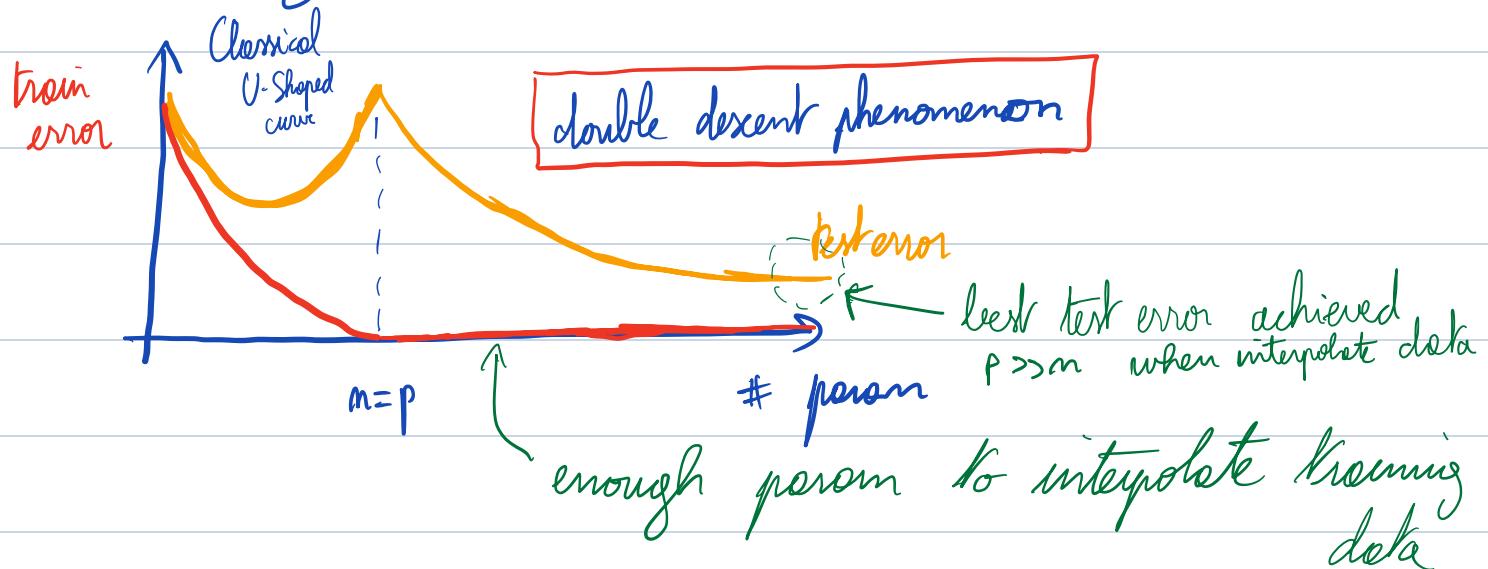
UNDERFITTING

Modern approach:

No careful control of the complexity of the model

4

In fact, generalizing well even when trained until interpolation
of the training data $f(x_i; \hat{\theta}) = y_i = f(x_i; \theta_0) + \varepsilon;$

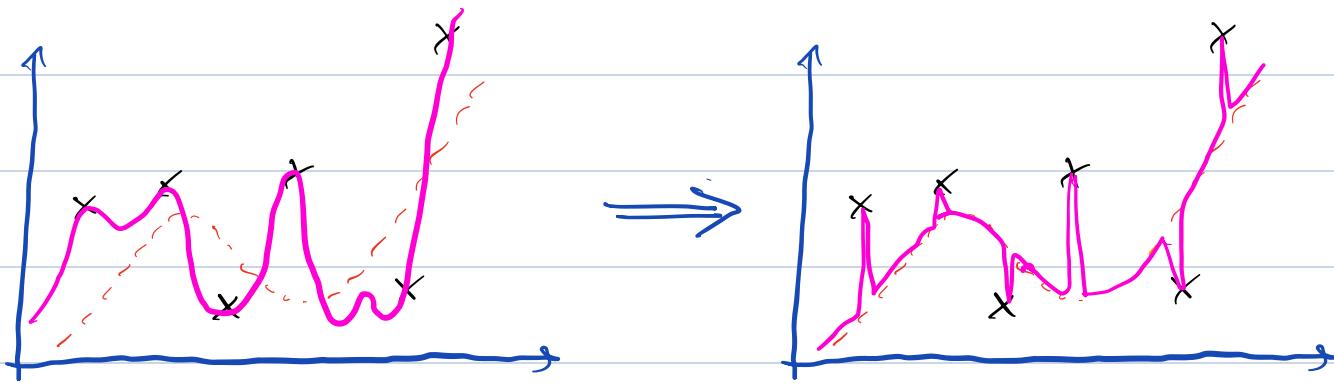


Today: I will present a scenario where the balance between the bias and the variance is achieved not by carefully tuning an explicit parameter but by a novel phenomena

Self-induced regularization

⇒ Show sufficient & necessary conditions for benign overfitting in linear models

(5)



Expected picture
for overfitting

Sometimes
benign

Benign overfitting

$$\text{estimator: } \hat{f}(x, \hat{\theta}) = \underbrace{\hat{f}_0(x, \hat{\theta}_0)}_{\text{good for prediction because smooth}} + \underbrace{\Delta(x)}_{\text{spikes that are useful for interpolation but do not harm prediction}}$$

only way to get $\hat{R}_n(\hat{\beta}) = 0$
 $R(\hat{\beta}) \ll 1$
 when noise

because $E_x [\Delta(x)^2] \ll 1$.

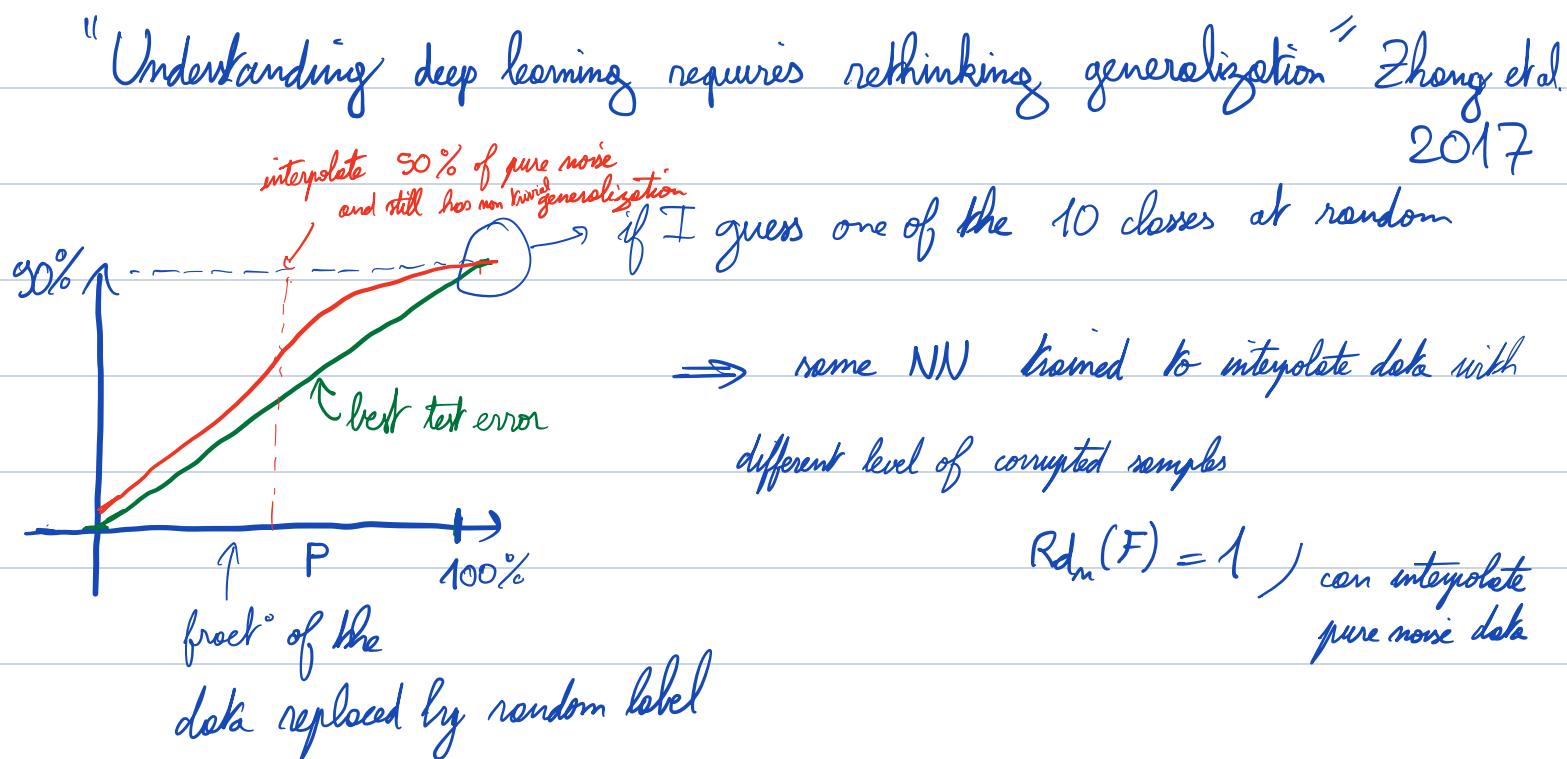
Rank 1: Double descent: we saw that # param is not a good measure of model complexity

→ w.r.t $\|\theta\|_2$ recover classical U-shaped curve

interesting phenomenon here is that we can interpolate noisy data and still generalize well.

Rank 2: Can UC explain benign overfitting?

→ bounds from lecture 2 independent of the data distribution

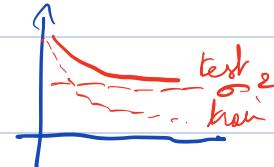


Simple rank: $y_i = f_\theta(x_i) + \varepsilon_i$

7

$$\hat{R}_m(\hat{f}) = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{f}(x_i))^2 \leq (1-\kappa) \sigma^2$$

$$R(\hat{f}) = \sigma^2 + \mathbb{E}[(\hat{f}(x) - f_\theta(x))^2]$$



$$\varepsilon_m(F) \geq |\hat{R}_m(\hat{f}) - R(\hat{f})| \geq \kappa \sigma^2$$

bounded away
from 0.

RIDGE REGRESSION

Setting: * Data $\{(y_i, x_i)\}_{i \leq m}$ $x_i \in \mathbb{R}^P$ $y_i \in \mathbb{R}$
 * $E[x_i] = 0$ $E[x_i x_i^T] = \Sigma$
 $x_i = \sum \frac{1}{2} \beta_j x_{ij}$

* x_i are iid $y_i = \langle \theta_*, x_i \rangle + \varepsilon_i$
 noise ε_i independent $E[\varepsilon_i] = 0$ $E[\varepsilon_i^2] = \sigma_\varepsilon^2$

Fit a linear model: $f(x, \theta) = \langle x, \theta \rangle$ $\theta \in \mathbb{R}^P$

Ridge regression: $\hat{\theta}(t) = \underset{\theta \in \mathbb{R}^P}{\operatorname{argmin}} \left\{ \|y - X\theta\|_2^2 + t \|\theta\|_2^2 \right\}$

$$X = \begin{bmatrix} & \xleftarrow{P} \\ \uparrow & \begin{bmatrix} x_1 & \cdots & x_m \end{bmatrix} \\ \downarrow & \end{bmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}$$

$$= X^T S y$$

where $S = (X X^T + t I_d)^{-1}$

Limit $t \searrow 0^+$. $\hat{\theta}(0^+) = \underset{\theta \in \mathbb{R}^P}{\operatorname{argmin}} \left\{ \|\theta\|_2 : y = X\theta \right\}$

$$= X^T S y \quad S = (X X^T)^+$$

\Rightarrow Minimum- $\|\cdot\|_2$ norm interpolating solution
 ↳ GD with $\theta^* = 0$ converges to $\hat{\theta}(0^+)$

9

Excess test error: $R(\hat{\theta}) = \mathbb{E}_x[(f(x, \hat{\theta}) - f(x, \theta_*))^2]$

$$= \mathbb{E}_x[(x, \hat{\theta} - \theta_*)^2]$$

$$= \|\hat{\theta} - \theta_*\|_{\Sigma}^2$$

$$= (\hat{\theta} - \theta_*)^T \Sigma (\hat{\theta} - \theta_*)$$

$$\text{Null}_{\Sigma} = \|\Sigma^{1/2} u\|_2$$

Bias - Variance decomposition:

$$\bar{R}(\hat{\theta}(\lambda)) = \mathbb{E}_{\varepsilon} [R(\hat{\theta}(\lambda))] = \mathbb{E}_{\varepsilon} [\underbrace{\|\hat{\theta}(\lambda) - \theta_*\|_{\Sigma}^2}]$$

$$X^T S y = X^T S X \theta_* + X^T S \varepsilon$$

$$= B(\lambda) + V(\lambda)$$

where $B(\lambda) = \|(\mathbb{I}_d - X^T S X) \theta_*\|_{\Sigma}^2$

$$V(\lambda) = \sigma^2 \text{Tr}(S X \Sigma X^T S)$$

10

Self-induced regularization: high level explanation

WLOG: assume $\Sigma = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$ $\lambda_1 \geq \lambda_2 \geq \dots$

Consider top k eigenspaces (corresponding to top k eigenvalues $\lambda_1, \dots, \lambda_k$)

write $x = \begin{bmatrix} x_0 \\ x_+ \end{bmatrix}$ $\left. \begin{array}{l} \text{first } k \text{ coordinates} \\ \text{p-k last coordinates} \end{array} \right\}$

$$X = \begin{bmatrix} -x_1- \\ \vdots \\ -x_n- \end{bmatrix} = \begin{bmatrix} \xrightarrow{k} & \xrightarrow{p-k} \\ X_0 & | & X_+ \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \xrightarrow{k} & & \\ \Sigma_0 & | & \Sigma_+ \\ \xleftarrow{p-k} & & \end{bmatrix} \quad \Theta = (\Theta_0, \Theta_+) \quad \begin{matrix} \downarrow \\ (\theta_1, \dots, \theta_k) \end{matrix} \quad \begin{matrix} \downarrow \\ (\theta_{k+1}, \dots, \theta_p) \end{matrix}$$

$$XX^\top = X_0 X_0^\top + \underbrace{X_+ X_+^\top}_{\approx \gamma \text{Id}_m}$$

(later, assume $\frac{1}{c} \text{Id} \leq M \leq c \text{Id}$)

$$\text{Recall: } \hat{\Theta}(\lambda) = X^\top S y \approx X^\top (X_0 X_0^\top + (\gamma + \lambda) \text{Id})^{-1} y$$

(11)

$$\rightarrow \hat{\theta}_0 = X_0^\top (X_0 X_0^\top + (\lambda + \gamma) \text{Id})^{-1} y$$

top k coordinates

$$\rightarrow \text{the same as } \hat{\theta}_0 = \underset{\theta_0 \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ \|y - X_0 \theta_0\|_2^2 + (\lambda + \gamma) \|\theta_0\|_2^2 \right\}$$

"self induced regularization" from the high-degree part of the features $X_+ X_+^\top$

Even when $\lambda = 0^+$: * $\hat{\theta}_0$ solution of a lower-dim problem with effective regularization $\gamma > 0$.

* $\hat{\theta}_+$ can show that it is small and does not contribute to the test error but interpolate the training data

$$f(x, \hat{\theta}) = \underbrace{\langle \hat{\theta}_0, x_0 \rangle}_{\text{"signal part of the features"} \ f_0(x)} + \underbrace{\langle \hat{\theta}_+, x_+ \rangle}_{\text{"noise part of feature" } \Delta(x)}$$

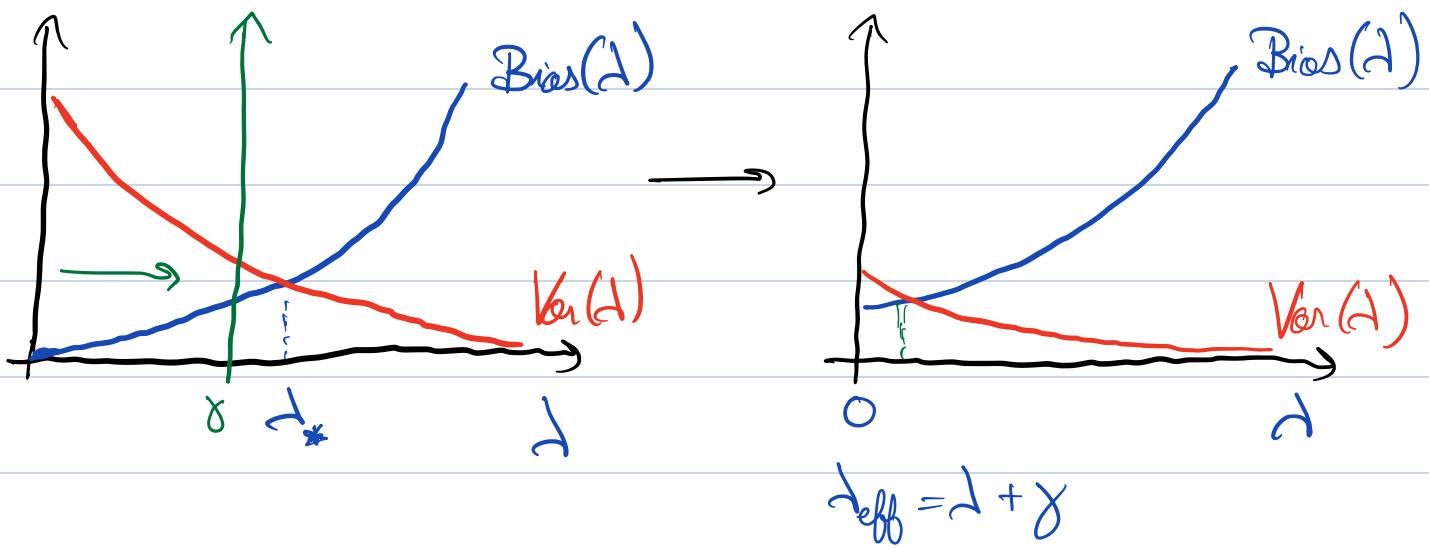
"signal part" of the features
 \rightarrow generalize well thanks to γ

"noise part" of feature
 \rightarrow interpolate data, doesn't contribute to test error

\Rightarrow Noise part acts as an effective ridge regularization

$$\text{Effective regularization } \lambda \rightarrow \lambda_{\text{eff}} = \lambda + \gamma$$

12



"Shifted Bias-Variance curve" $\lambda=0$ close to optimal regularization

"Benign Overfitting in Ridge Regression"

(13)

- Bartlett, Trigler, 2020

$$x_i = \sum_{j=1}^n z_j: \quad z_j \text{ iid sub-Gaussian (for simplicity)}$$

more general

even for non linear kernels

Recall: $S = (\lambda \text{Id} + X X^\top)^{-1} = (\lambda \text{Id} + X_+ X_+^\top + X_0 X_0^\top)^{-1}$

$$S_+ = (\lambda \text{Id} + X_+ X_+^\top)^{-1} \quad S = (S_+^{-1} + X_0 X_0^\top)^{-1}$$

THM [Bartlett, Trigler, '20] Assume w.h.p $\frac{\lambda_{\max}(S_+)}{\lambda_{\min}(S_+)} \leq L$

Then w.h.p :

$$\text{Bias}(\lambda) \lesssim L^4 \left[\|\theta_0^*\|_{\Sigma_0^{-1}}^2 \left(\frac{\lambda + \sum_{i>L} \lambda_i}{m} \right)^2 + \|\theta_+^*\|_{\Sigma_+}^2 \right]$$

$$\text{Var}(\lambda) \lesssim \epsilon_\varepsilon^2 L^2 \left[\frac{k}{m} + \frac{m \sum_{i>L} \lambda_i^2}{(\lambda + \sum_{i>L} \lambda_i)^2} \right]$$

- RMK:
- * Matching lower bound (up to multiplicative constants / minimax)
 - * In particular, interpolators $\lambda = 0$ can be near optimal.
 - * Does not fit θ_+^* at all
 - * Does not depend on p (if sums $\sum \lambda_i$ cv : applies to $p=\infty$)

15

B.O: small variance: need to choose k so that

$$(1) \quad k \ll m$$

$$(2) \quad \frac{\left(\sum_{i>k} d_i\right)^2}{\sum_{i>k} d_i^2} \gg m$$

E.g. $d_{k+1} := \# \text{ eigenvalues of } \Sigma \text{ in } \left[\frac{d_{k+1}}{2}, d_{k+1} \right]$

and assume eigenvalues $< \frac{d_{k+1}}{2}$ are negligible

$$\Rightarrow (2) \text{ requires } \frac{(d_{k+1} \cdot d_{k+1})^2}{d_{k+1} \cdot d_{k+1}^2} = d_{k+1} \gg m$$

and we have Variance $\lesssim \frac{k}{m} + \frac{m}{d_{k+1}}$

Variance of fitting θ_0^*
 (actually this is the parametric rate!)

effect of overfitting
 ↓
 negligible

Rmk: $\lambda_i = \frac{1}{i^\alpha (1 + \log i)^\beta}$ decay of eigenvalues

Test error ($\lambda = 0^+$) $\rightarrow 0$ as $n \rightarrow \infty$

IFF $\alpha = 1$ $\beta > 1$

"consistency of interpolating solution"

To get benign overfitting here:

(1) Overparametrization: $p \gg n$

directions in parameter space unimportant for prediction \Rightarrow sample size

(2) Smallest eigenvalues of Σ must decay slowly

$$\lambda_{k+1} \gg n$$

Self-induced regularization

$$\sum z_i^2$$

z_i iid c -sub-G coordinates

hold way more generally

Lemma: $\exists c > 0$ s.t. with prob at least $1 - 2 \exp(-\frac{m}{c})$

$$\frac{\lambda_{\max}(X_+ X_+^T)}{\lambda_{\min}(X_+ X_+^T)} \leq c \frac{\sum_{i \leq k} \lambda_i + m \lambda_{k+1}}{\sum_{i \leq k} \lambda_i - c m \lambda_{k+1}} \leq L$$

$m \lambda_{k+1} \ll \sum_{i \leq k} \lambda_i$

Example: Kernel regression with inner product kernel on the sphere

$$u_i \underset{iid}{\sim} \text{Unif}(S^{d-1}(\sqrt{d}))$$

$$y_i = f_\infty(u_i) + \varepsilon_i$$

→ general squared integrable function

Linear model: $\hat{f}(u, \theta) = \langle \theta, \phi(u) \rangle$

$$\phi(u) = (\zeta_0, \zeta_1 u_1, \dots, \zeta_d u_d)$$

$$\zeta_k = O(d^{-\frac{k}{2}})$$

$$\zeta_2 Y_{21}(u), \dots, \zeta_2 Y_{2d^2}(u)$$

⋮

$$\zeta_k Y_{k1}(u), \dots, \zeta_k Y_{kO(d^k)}(u) \quad) \text{ basis of degree } k \text{ polynomials}$$

in u .

$$O(d^k)$$

$$x_i = \phi(u_i) = \sum z_i$$

$$\mathbb{E}[z_i z_i^\top] = \text{Id}_\infty$$

$$\Sigma = \begin{pmatrix} \zeta_0^2 & & & \\ & \ddots & & \\ & & \zeta_1^2 & \\ & & & \zeta_2^2 \\ & & & & \ddots & & \end{pmatrix}$$

d
 $O(d^2)$

(17)

$$\text{In } d^l \ll n \ll d^{l+1}$$

$$\text{Take } k := d^l \quad X_+ X_+^T \approx \left(\sum_{s=l+1}^{\infty} B_{d,s} \mathbb{Z}_s^2 \right) \text{Id}_n$$

$$f_*(u) = \underbrace{\langle \alpha_0, \Theta_{*,0} \rangle}_{\text{best degree-}l \text{ polynomial approx}} + \langle \alpha_+, \Theta_{*,+} \rangle$$

$$\hat{\Theta}_0 \approx \Theta_{*,0} \quad \mathbb{E}[\hat{\Theta}_+, \alpha_+] = \|\hat{\Theta}_+\|_{\Sigma_+}^2 \approx 0$$

[needs very different tools to analyze \rightarrow not sub-Gaussian]

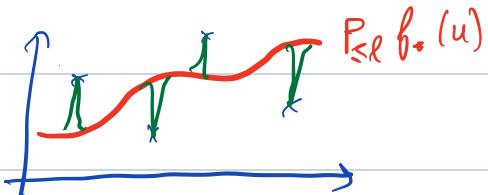
Thm [M., Ghorbani, Mei, Montanari, '19]

If $d^{l+\delta} \leq n \leq d^{l+1-\delta}$ as $n, d \rightarrow \infty$, then

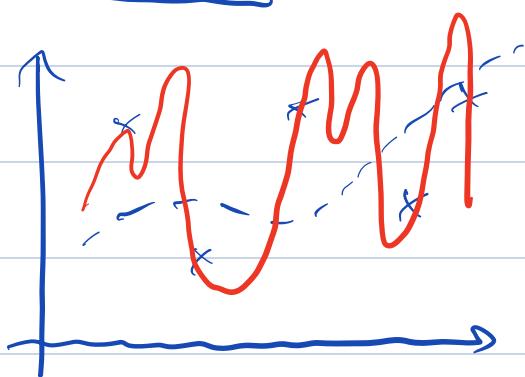
$$\text{Test error} = \|P_{\geq l} f_*\|_{L^2}^2 + o(1)$$

$$\hat{f}(u) = P_{\leq l} f_*(u) + \Delta(u)$$

spiky part (high degree polynomial)



Conclusion:

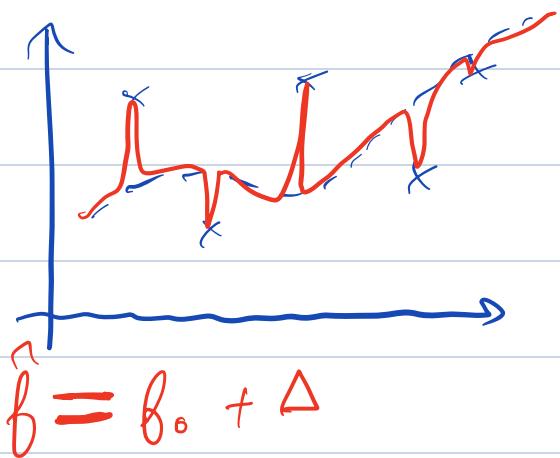


classical picture

BAD OVERFITTING



Modern picture



sometimes

BENIGN OVERFITTING

$$\hat{f} = f_0 + \Delta$$

Above "low-dim" ^{classical} picture misleading

→ much easier to interpolate benignly
in high dim

(Of course, it required a lot of work to realize it)

(19)

Proof of Bartlett / Trigler

$$S = (X_0 X_0^T + X_+ X_+^T + I)^{-1} \quad S_+ = (X_+ X_+^T + I)^{-1}$$

$$\frac{\lambda_{\max}(S_+)}{\lambda_{\min}(S_+)} \leq L$$

Here only prove Variance term ($\sigma_\epsilon^2 = 1$ WLOG)

$$V(X, I) = \text{Tr}(SX\Sigma XS) = V_0 + V_+$$

$$V_0 := \text{Tr}(SX_0 \Sigma_0 X_0^T S)$$

$$V_+ := \text{Tr}(SX_+ \Sigma_+ X_+^T S)$$

1st Term: $V_0 = \text{variance along top } k \text{ directions}$

$$\lesssim L^2 \frac{k}{m}$$

Proof: $SX_0 = (X_0 X_0^T + S_+^{-1})^{-1} X_0$

$$= S_+ X_0 (I + X_0^T S_+ X_0)^{-1}$$

Introduce $M := \sum_0^{\frac{1}{2}} (I + X_0 S_+ X_0^T)^{-1} \sum_0^{\frac{1}{2}}$

$V_0 = \text{Tr}(S_+ X_0 \sum_0^{\frac{1}{2}} M^2 \sum_0^{\frac{1}{2}} X_0^T S_+)$

(20)

$$V_0 \leq \|M\|_{op}^2 \operatorname{Tr}(S_+ Z_0 Z_0^T S_+)$$

$$\begin{aligned} \|M\|_{op} &= \left\| \left(\Sigma_0^{-1} + \Sigma_0^{-\frac{1}{2}} X_0^T S_+ X_0 \Sigma_0^{-\frac{1}{2}} \right) \right\|_{op} \leq \frac{1}{\lambda_{min}(Z_0^T S_+ Z_0)} \\ &\leq \frac{1}{\lambda_{max}(S_+) \lambda_{min}(Z_0^T Z_0)} \end{aligned}$$

Hence

$$V_0 \leq \frac{\lambda_{max}(S_+)^2 \operatorname{Tr}(Z_0^T Z_0)}{\lambda_{min}(S_+)^2 \cdot \lambda_{min}(Z_0^T Z_0)^2}$$

$$\frac{\operatorname{Tr}(A)}{\lambda_{min}(A)} \leq \frac{\lambda_{max}(A)}{\lambda_{min}(A)} \cdot \operatorname{rank}(A)$$

$$\leq L^2 \frac{\lambda_{max}(Z_0^T Z_0)}{\lambda_{min}(Z_0^T Z_0)} \times \frac{k}{\lambda_{min}(Z_0^T Z_0)}$$

$$\left. \begin{aligned} Z_0 &\in \mathbb{R}^{m \times p} \text{ iid } O(1) - \text{subGaussian} \quad k \leq c m \\ \lambda_{min}(Z_0), \lambda_{max}(Z_0) &\asymp \sqrt{m} \end{aligned} \right]$$

$$V_0 \leq L^2 \cdot \frac{k}{m}$$

□

2nd term

$V_+ = \text{variance due to overfitting}$

$$\leq L^2 \frac{\sum_{i>L} d_i^2}{(\lambda + \sum_{i>L} d_i)^2}$$

Rank:

$$\frac{\lambda_{\max}(S_+)}{\lambda_{\min}(S_+)} \leq L \Rightarrow \lambda_{\max}(S_+) \leq \frac{L}{\lambda + \sum_{i>L} d_i}$$

$$\lambda_{\max}(S_+) \leq L \lambda_{\min}(S_+) \leq L \frac{n}{\text{Tr}(S_+^{-1})}$$

$$\frac{1}{n} \text{Tr}(S_+^{-1}) = \frac{1}{n} \text{Tr}(\lambda \text{Id} + X_+ X_+^T) = \lambda + \frac{1}{n} \sum_{i=1}^n \|x_{i,+}\|_2^2$$

WLLN

$$\asymp \lambda + \text{Tr}(\Sigma_+)$$

Proof: $V_+ = \text{Tr}(S X_+ \Sigma_+ X_+^T S) \leq \lambda_{\max}(S)^2 \text{Tr}(\Sigma_+ X_+^T S)$

$$\leq \lambda_{\max}(S_+)^2 \sum_{i=1}^n \langle z_{i,+}, \Sigma_+ z_{i,+} \rangle$$

$$\leq \lambda_{\max}(S_+)^2 n \cdot \text{Tr}(\Sigma_+^2)$$

Using the above remark

$$V_+ \leq L^2 \frac{n \cdot h(\sum_+^2)}{(2 + \sum_{i>h} d_i)^2}$$

Sharp asymptotics using RMT

So far we only used matrix concentration
 ↳ in particular: bounds hold for $n \geq c$.

⇒ lots of work have considered to derive sharp formulas for the bias and variance as $n, p \rightarrow \infty$ using RMT

Below a very brief presentation of these results

[Mestre, Montanari, Rosset, Tibshirani, 2013]

(Follow Andree's presentation)

- Setting: * $\{(x_i, y_i)\}_{i \leq m}$ $y_i \in \mathbb{R}$ $x_i \in \mathbb{R}^p$ iid
- * $y_i = \langle \theta_*, x_i \rangle + \varepsilon_i$ $\mathbb{E} \varepsilon_i = 0$ $\mathbb{E} \varepsilon_i^2 = \sigma_\varepsilon^2$
- * $x_i = \sum z_i$ $\mathbb{E} z_i = 0$ $\mathbb{E}[z_i z_i^\top] = I_p$
- * z_i independent coordinates $\mathbb{E}[|z_{ij}|^k] \leq C_k < \infty \quad \forall k$

Ridge solution $\hat{\theta}(\lambda) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \frac{1}{m} \|y - X\theta\|_2^2 + \lambda \|\theta\|^2 \right\} = \frac{1}{m} X^\top (\lambda + \frac{1}{m} X X^\top)^{-1} y$

$$R(\hat{\theta}) = \|\hat{\theta}(\lambda) - \theta_*\|_2^2 = \beta_m + V_m$$

$$B_m = \lambda^2 \langle \theta_0, S_\lambda \sum S_\lambda \theta_0 \rangle \quad \hat{\Sigma} = \frac{1}{n} X^T X$$

$$V_m = \frac{\sigma^2}{n} \text{Tr}(\Sigma \hat{\Sigma} S_\lambda^{-2}) \quad S_\lambda = (\hat{\Sigma} + \lambda I)^{-1}$$

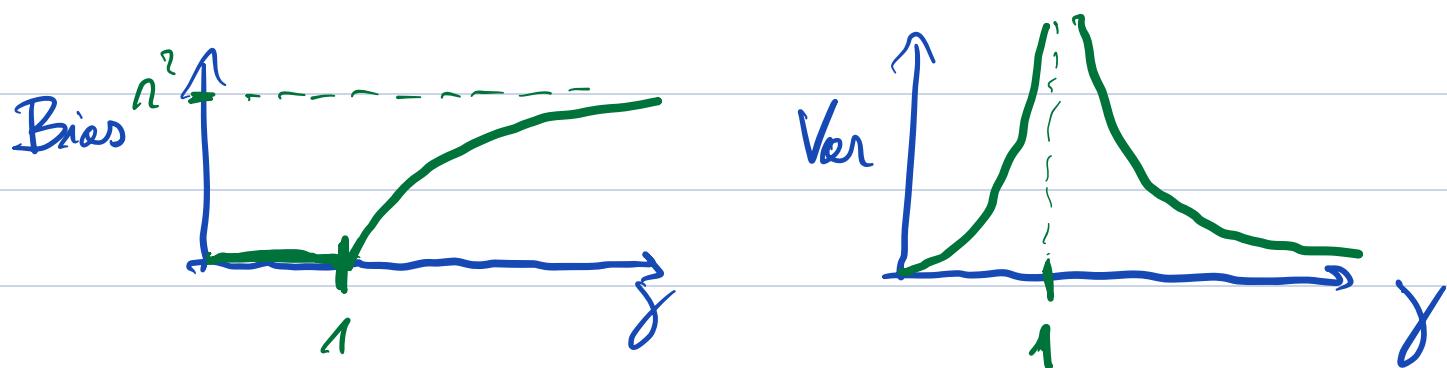
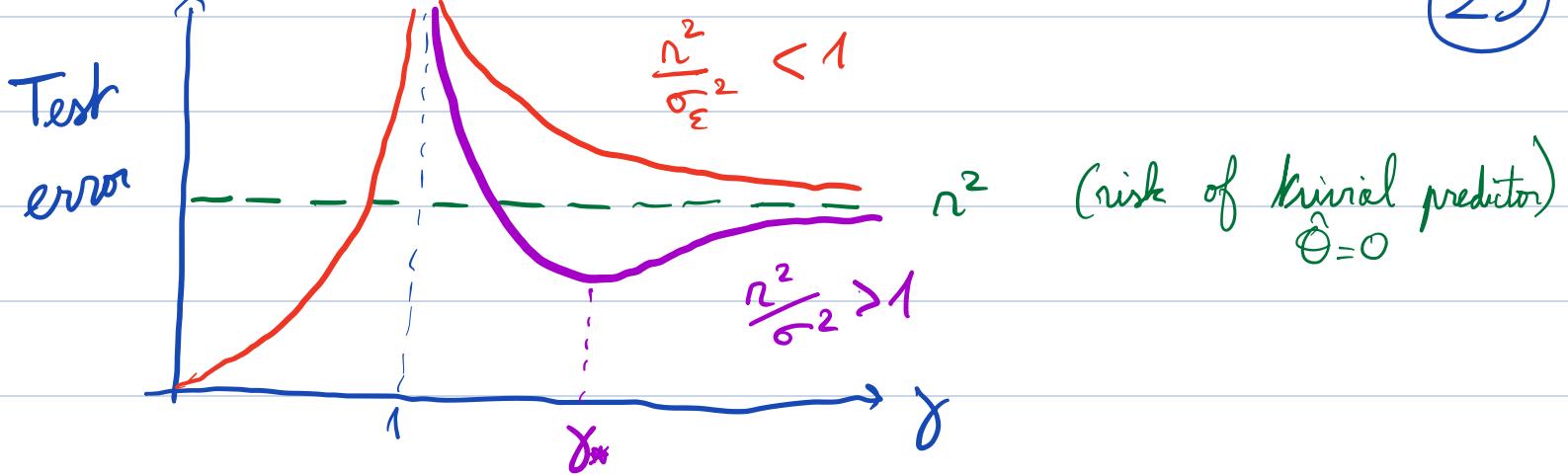
Then [Mortie et al, 19] $\Sigma = I \quad \lambda = 0^+ \text{ (min-norm)}$
 $m, p \rightarrow \infty \quad f_m \rightarrow \gamma \in (0, \infty) \quad \| \theta_0 \|_2^2 \rightarrow n^2$

Then $\lim_{m, p \rightarrow \infty} B_m = n^2 \left(1 - \frac{1}{\gamma}\right)_+$

$$\lim_{m, p \rightarrow \infty} V_m = \sigma^2 \frac{1}{\left(\frac{1}{\gamma} \vee \gamma - 1\right)}$$

In particular,

$$\lim_{m, p \rightarrow \infty} R(\hat{\theta}) = \begin{cases} \frac{\gamma \sigma^2}{1-\gamma} & \text{if } \gamma < 1 \\ n^2 \left(1 - \frac{1}{\gamma}\right) + \frac{\sigma^2}{\gamma-1} & \text{if } \gamma > 1 \end{cases}$$



- Rank:
- * Prove way more stuff in their paper
 - * Double descent
→ however not minimizer $\gamma \rightarrow \infty$
 - * Not a good model $p = d$
↳ true fit $f_{*}(x) = \langle \theta_{*}, \alpha \rangle$ becomes more complex as p grows
 - * Better models
[Mei, Montanari, '19]

Crash course Random Matrix Theory (RMT)

26

$$\hat{\Sigma} = \frac{XX^T}{n} \in \mathbb{R}^{n \times n}$$

empirical spectral distribution (ESD)

$$\mu_m = \frac{1}{P} \sum_{i=1}^P \delta_{\lambda_i(\hat{\Sigma})}$$

(weak cr)

Marchenko - Pastur law: $\mu_m \xrightarrow[m, p \rightarrow \infty]{m/p \rightarrow \gamma} \mu_\infty$

To prove it: Stieltjes transform $z \in \mathbb{C}$

$$M_m(z) = \frac{1}{P} \sum_{i=1}^P \frac{1}{\lambda_i(\hat{\Sigma}) - z} = \frac{1}{P} h((\hat{\Sigma} - zI_p)^{-1})$$

Thm: [MP] let $z \mapsto m(z)$ be the analytic function $z \in \mathbb{C}_+$

s.t. $|m(z)| \leq \frac{C}{|z|}$ as $|z| \geq C$ and

solution of

$$\frac{m}{1+z m} = 1 + \gamma m$$

Then a.s. $M_m(z) \rightarrow m(z)$ as $\frac{n, p \rightarrow \infty}{n/p \rightarrow \gamma}$

Stieltjes inversion:

$$\mu([a, b]) = \lim_{\eta \rightarrow 0} \int_a^b \frac{1}{\pi} \operatorname{Im}(m(\lambda + i\eta)) d\lambda$$

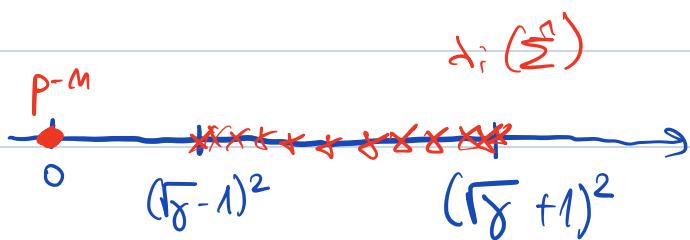
Remark: $m(z) = \frac{1-z-\gamma - \sqrt{(1-z-\gamma)^2 - 4\gamma z}}{2\gamma z}$

by solving the
fixed pt equat°

Theorem: [Bei-Yin law] $\gamma > 1 \quad \operatorname{rank}(\hat{\Sigma}) = m \quad \text{w.h.p.}$

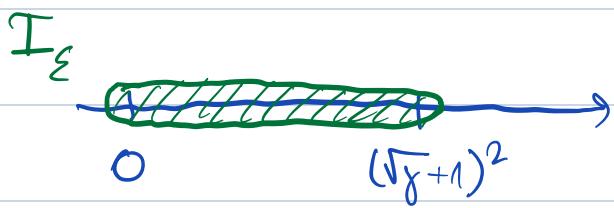
$$\left\{ \begin{array}{l} \lambda_{\max}(\hat{\Sigma}) = (\sqrt{\gamma} + 1)^2 + o_p(1) \\ \lambda_{\min}(\hat{\Sigma}) = (\sqrt{\gamma} - 1)^2 + o_p(1) \end{array} \right.$$

For $\gamma \leq 1 \quad \lambda_i(\hat{\Sigma}) = 0 \quad p-m$



Rmk: $M_m(z)$ is unif. cont. (differentiable)
on $\mathbb{C} \setminus I_\varepsilon$

$$I_\varepsilon = \left\{ z : \operatorname{Re}(z) \in [-\varepsilon, (\sqrt{\gamma} + 1)^2 + \varepsilon] \quad \left| \operatorname{Im}(z) \right| \leq \varepsilon \right\}$$



All derivatives unif. bounded
outside I_ε

(convince yourself using expression of $M_m(z)$)

$$M_m(z) \xrightarrow{a.s.} m(z) \quad \forall z \notin I_\varepsilon.$$

and some for derivatives

Take $z = -\lambda$ $M_m(-\lambda) = \frac{1}{P} \operatorname{Tr}\left((\hat{\Sigma} + \lambda)^{-1}\right) \rightarrow m(-\lambda)$

Proof of Thm:

$$\begin{aligned}
 \boxed{\text{Variance}} \quad & \frac{1}{\sigma^2} V_m = \frac{1}{m} \mathbb{E} \left(\sum (\hat{\Sigma} + \lambda I)^{-2} \right) \\
 &= \frac{1}{m} \left[\frac{1}{P} \mathbb{E} \left((\hat{\Sigma} + \lambda I)^{-1} \right) - \frac{1}{P} \lambda \mathbb{E} \left((\hat{\Sigma} + \lambda I)^{-2} \right) \right] \\
 &= \frac{1}{m} \left[M_m(-\lambda) + \lambda \frac{\partial}{\partial \lambda} M_m(-\lambda) \right] \\
 &\rightarrow \gamma \left[m(-\lambda) + \lambda \frac{\partial}{\partial \lambda} m(-\lambda) \right] \\
 &= \gamma \frac{\partial}{\partial \lambda} (\lambda m(-\lambda))
 \end{aligned}$$

$$m(-\lambda) = \frac{-1 - \lambda - \gamma + \sqrt{(1+\lambda-\gamma)^2 + 4\lambda\gamma}}{2\gamma\lambda}$$

Ridgeless limit: $\lambda \downarrow 0$ $\gamma > 1$ $m(-\lambda) = \frac{\gamma-1}{\gamma\lambda} + \frac{1}{\gamma(\gamma-1)} + O(\lambda)$

$$\frac{\partial}{\partial \lambda} (\lambda m(-\lambda)) = \frac{1}{\gamma(\gamma-1)} + O(\lambda)$$

$$\lim_{\lambda \rightarrow 0} \lim_{m, P \rightarrow \infty} V_m(\lambda) = \frac{\sigma^2}{\gamma-1} \quad \gamma > 1$$

Rank: * need to invert the limits (doable with some work)
 * same for $\gamma < 1$

Bias

$$B_m = \lambda^2 \langle \theta_*, (\hat{\Sigma} + \lambda)^{-2} \theta_* \rangle$$

α isotropic: $B_m = \lambda^2 \|\theta_*\|_2^2 \frac{1}{p} \text{Tr}((\hat{\Sigma} + \lambda)^{-2}) + o(1)$

[e.g. prove it for $\alpha \sim N(0, \text{Id}_d)$]

$$B_m = \lambda^2 \|\theta_*\|_2^2 \left(-\frac{\partial}{\partial \lambda} M_m(-\lambda) \right)$$

$$\rightarrow n^2 (1 - \frac{1}{\gamma})_+ \quad \text{with some calculus}$$

□

31

Anisotropic case

$$\Sigma \neq I$$

risk will depend on both Σ, θ_*

Interpolating solution $\lambda = 0^+$

"Effective regularization" $= \lambda_*$ unique positive solution of

$$\frac{1}{\lambda} = \frac{1}{p} \operatorname{Tr}(\Sigma(\Sigma + \lambda_* I)^{-1})$$

$$\overline{B}_m := \frac{\lambda_*^2 \langle \theta_*, (\Sigma + \lambda_*)^{-2} \Sigma \theta_* \rangle}{1 - \frac{1}{m} \operatorname{Tr}(\Sigma^2 (\Sigma + \lambda_*)^{-2})}$$

$$\overline{V}_m := \frac{\sigma^2}{m} \frac{\operatorname{Tr}(\Sigma^2 (\Sigma + \lambda_*)^{-2})}{1 - \frac{1}{m} \operatorname{Tr}(\Sigma^2 (\Sigma + \lambda_*)^{-2})}$$

Thm [Mortie et al, '19]

$$\left\{ \begin{array}{l} \lim_{n, p \rightarrow \infty} |B_n - \overline{B}_m| = 0 \\ \lim_{n, p \rightarrow \infty} |V_n - \overline{V}_m| = 0 \end{array} \right.$$

a.s.

Interpretation: test error:

(32)

Ridge regression \longleftrightarrow gaussian sequence model

$$y^s = \sum^{\frac{1}{2}} \theta_* + \frac{\omega}{\sqrt{m}} g \quad g \sim N(0, I_p)$$

$$\hat{\theta}^s = \underset{\theta}{\operatorname{argmin}} \left\{ \|y^s - \sum^{\frac{1}{2}} \theta\|_2^2 + \lambda_* \|\theta\|_2^2 \right\}$$

$$\text{w fixed pt of: } \omega^2 = \frac{\sigma_\varepsilon^2}{m} + \mathbb{E}_g \left[\|\hat{\theta}^s - \theta_*\|_{\Sigma}^2 \right]$$

$$R_m^s = \mathbb{E} \left[\|\hat{\theta}^s - \theta_*\|_{\Sigma}^2 \right] = \underbrace{B_m^s}_{\parallel} + \underbrace{V_m^s}_{\parallel}$$

$$\overline{B}_m \quad \overline{V}_m$$

	Original problem	Sequence model
design	X random	$\sum^{\frac{1}{2}}$ deterministic
reg.	$\lambda = 0^+$	$\lambda_* > 0 !!$
noise	σ_ε^2	ω^2

self induced regularization!! more regularization
in the model than what is merely expected from $\lambda=0$.

Appendix:Concentration of $X_+ X_+^\top$

I'll present a nice decoupling argument that applies under quite general conditions
 (much beyond sub-Gaussian e.g. inner-product kernels)

$$\{u_i\}_{i \in [m]} \text{ i.i.d. } K = (K(u_i, u_j))_{i, j \in [m]} \in \mathbb{R}^{m \times m}$$

$K: \mathcal{U} \times \mathcal{U} \rightarrow \mathbb{R}$ p.s.d. fct \rightarrow p.s.d operator $\|K\|$
 (lecture on kernels)

$$K(u_i, u_j) = \sum_{k=1}^{\infty} \beta_k^2 \phi_k(u_i) \phi_k(u_j)$$

EIGENDECOMPOSITION

$$\beta_1^2 \geq \beta_2^2 \geq \dots$$

eigenvalues ↓
 eigenvectors $\mathbb{E}[\phi_\ell(u)\phi_\ell(u)] = \delta_{\ell,\ell}$

$$\|K\|_{op} = \beta_1^2$$

$$\begin{aligned} \mathrm{Tr}(K) &= \sum_{k=1}^{\infty} \beta_k^2 < \infty \quad (\text{"trace-class"}) \\ &= \mathbb{E}[K(u, u)] \end{aligned}$$

BS

Assumptions:

(A) Diagonal: $\mathbb{E} \left[\max_{i \in [n]} |K(u_i, u_i) - T_K(K)| \right] \leq C \sqrt{n \|K\|_{op} T_K(K)}$

(B) Hypercontractivity: $\forall p \text{ integer } \geq 2, \exists C_p \text{ s.t.}$

$$\forall v \in \mathbb{R}^n \quad \mathbb{E} \left[\left(\sum v_k \phi_k(u) \right)^p \right]^{\frac{1}{p}} \leq C_p \|v\|_2$$

Thm: [M., Mei, Montanari, 2021] For all $\delta > 0$, $\exists C$ s.t.

$$\mathbb{E} [\|K - T_K(K)\|_q] \leq C \sqrt{n^{1+\delta} \|K\|_{op} T_K(K)}$$

Rmk: * Bounds on expectation: use Markov to get w.h.p
 ↳ to get better dependency in prob → sub-G: union bound
 can also improve $n^{1+\delta} \rightarrow n$ ↳ hyper: moment method

* $(\phi_k(u_i))_{k \geq 1} \rightarrow z_i$ i.i.d sub-Gaussian entries

→ easy to check (A) & (B)

* $\lambda_{\max}(K) \leq T_K(K) \left(1 + C \sqrt{\frac{n^{1+\delta} \|K\|_{op}}{T_K(K)}} \right)$

$\lambda_{\min}(K) \geq T_K(K) \left(1 - C \sqrt{\frac{n^{1+\delta} \|K\|_{op}}{T_K(K)}} \right)$

We will use the following result in random matrix concentration

Prop [Vershynin 2010: Theorem 5.48] Let $A \in \mathbb{R}^{n \times k}$ with

$A = [\alpha_1, \dots, \alpha_m]^T$ where $\alpha_i \in \mathbb{R}^k$ are indep random vectors with common second moment matrix $\Sigma = \mathbb{E}[\alpha_i \alpha_i^T]$.

Let $\Gamma = \mathbb{E}\left[\max_{i \in [m]} \|\alpha_i\|_2^2\right]$. Then there exists a universal constant C such that

$$\mathbb{E}\left[\|A\|_{op}^2\right]^{\frac{1}{2}} \leq \left(\|\Sigma\|_{op} \cdot m\right)^{\frac{1}{2}} + C \cdot (\Gamma \cdot \log(m \wedge k))^{\frac{1}{2}}$$

Proof of theorem: $K = \text{ddiag}(K) + \Delta$

where $\text{ddiag}(K) = \text{diag}(K(u_i, u_i))_{i \in [m]}$

$$\Delta_{ij} = \begin{cases} K(u_i, u_j) & i \neq j \\ 0 & \text{if } i = j. \end{cases}$$

so that $\mathbb{E}\left[\|K - \text{Tr}(K) \cdot I_m\|_{op}\right]$

$$\leq \mathbb{E}\left[\|\text{ddiag}(K) - \text{Tr}(K) \cdot I_m\|_{op}\right] \stackrel{(I)}{=} + \mathbb{E}\left[\|\Delta\|_{op}\right] \stackrel{(II)}{=}$$

Using assumption (A) :

$$(I) = \mathbb{E} \left[\max_{i \in [m]} |K(u_i, u_i) - T(K)| \right] \leq \sqrt{m \|K\|_{op}} T(K)$$

For the second term: note that Δ does not have independent rows or columns and we can't use the above proportion directly.

We will use a standard decoupling argument:

Lemma:

$$\text{Let } \Delta_{T_1, T_2} = (\Delta_{ij})_{i \in T_1, j \in T_2} \in \mathbb{R}^{|T_1| \times |T_2|}$$

$$\mathbb{E} [\|\Delta\|_{op}] \leq 4 \max_{T \subseteq [m]} \mathbb{E} [\|\Delta_{T, T^c}\|_{op}]$$

Proof of lemma: For $\|v\|_2 = 1$

$$v^T \Delta v = \sum_{i \neq j} v_i \Delta_{ij} v_j = 4 \mathbb{E}_T \left[\sum_{\substack{i \in T \\ j \in T^c}} v_i \Delta_{ij} v_j \right]$$

where T is a random subset of $\{1, \dots, m\}$ where each element is selected with proba $1/2$)

$$\text{Therefore } \mathbb{E}\{\|\Delta\|_{\text{op}}\} = \mathbb{E}\left[\sup_{v \in S^m} v^T \Delta v\right]$$

$$\leq 4 \mathbb{E}_T \mathbb{E}\left[\sup_{v \in S^m} \sum_{\substack{i \in T \\ j \in T^c}} v_i \Delta_{ij} v_j\right] \leq 4 \sup_{T \subseteq [n]} \mathbb{E}\{\|\Delta_{T, T^c}\|\}$$

□

Using this lemma, it suffices to bound $\|\Delta_{T, T^c}\|_{\text{op}}$ for all $T \subseteq [n]$

$$\Delta_{T, T^c} = (\Delta_{ij})_{i \in T, j \in T^c}$$

Each row is $\Delta_{i, T^c} = (K(u_i, u_j))_{j \in T^c}$

so conditional on $(u_j)_{j \in T^c}$, Δ_{T, T^c} has iid rows!!

Denote \mathbb{E}_T the expectation over $(u_i)_{i \in T}$ conditional on $(u_j)_{j \in T^c}$. We can use the above general concentration bound

$$\mathbb{E}_T\{\|\Delta_{T, T^c}\|_{\text{op}}\} \leq \left\{\|\Sigma_T\|_{\text{op}} |T|\right\}^{1/2} + C \left\{\Gamma_T \cdot \log(|T| \wedge |T^c|)\right\}^{1/2}$$

where we denoted

$$\Sigma_T := \mathbb{E}_{u_i} [\Delta_{i, T^c} \Delta_{i, T^c}^T] \quad \Gamma_T := \mathbb{E}_T \left[\max_{i \in [n]} \|\Delta_{i, T^c}\|_2^2 \right]$$

Mence

$$\boxed{\mathbb{E}[\|\Delta\|_{op}] \leq \sup_{T \subseteq [m]} \left\{ (\mathbb{E}_{T^c}[\|\Sigma_T\|_{op}] \cdot m)^{\frac{1}{2}} + C(\mathbb{E}_{T^c}[T] \cdot \log m)^{\frac{1}{2}} \right\}}$$

Bounding $\mathbb{E}_{T^c}[\|\Sigma_T\|_{op}]$

$$\begin{aligned} \|\Sigma_T\|_{op} &= \|\mathbb{E}_{u_i}[\Delta_{iT^c} \Delta_{iT^c}^T]\|_{op} = \sup_{\|v\|_2=1} \sum_{ij \in T^c} \sum_k \mathbb{E}_{u_i} \phi_k(u_i) \phi_k(u_j) v_i v_j \\ &\leq \|K\|_{op} \sup_{\|v\|_2=1} \sum_{ij \in T^c} \mathbb{E}_{u_i} \phi_k(u_i) \phi_k(u_j) v_i v_j \\ &\leq \|K\|_{op} \left\| (K(u_i, u_j))_{ij \in T^c} \right\|_{op} \\ &\leq \|K\|_{op} (\| \text{ddiag } K \|_{op} + \|\Delta\|_{op}) \end{aligned}$$

By hypercontractivity:

$$\begin{aligned} \mathbb{E}[\| \text{ddiag } K \|_{op}] &\leq \mathbb{E}\left[\max_{i \in [m]} K(u_i, u_i)\right]^{\frac{1}{p}} \leq \mathbb{E}\left[\sum_{i=1}^m K(u_i, u_i)^p\right]^{\frac{1}{p}} \\ &\leq m^{\frac{1}{p}} \mathbb{E}[K(u_i, u_i)^p]^{\frac{1}{p}} \leq C_p m^{\frac{1}{p}} \text{Tr}(K) \end{aligned}$$

Hence

$$\mathbb{E}_{T^c}[\Sigma_T] \leq C_p m^{\frac{1}{p}} \|K\|_{op} \text{Tr}(K) + \|K\|_{op} \mathbb{E}[\|\Delta\|_{op}]$$

Bounding $\mathbb{E}_{T^c}[\Gamma_T]$

By hypercontractivity,

$$\mathbb{E}_{T^c}[\Gamma_T] = \mathbb{E}[\max_{i \in T} \|\Delta_{iT^c}\|_2^2]$$

$$\leq n \cdot \mathbb{E}[\max_{i \in T, j \in T^c} \Delta_{ij}^2]$$

$$\leq n \cdot \mathbb{E} \left[\sum_{i,j \in [m]} \Delta_{ij}^{2p} \right]^{\frac{1}{p}}$$

$$\leq n^{1+\frac{2}{p}} \cdot \mathbb{E}[\Delta_{ij}^{2p}]^{\frac{1}{p}}$$

$$\leq n^{1+\frac{2}{p}} C_p^2 \mathbb{E}[\Delta_{ij}^2]$$

$$\mathbb{E}_{T^c}[\Gamma_T] \leq n^{1+\frac{2}{p}} C_p^2 \|K\|_{op} \text{Tr}(K)$$

Putting everything together:

40

$$\mathbb{E}[\|\Delta\|_{\text{op}}] \leq K_p \left\{ \|K\|_{\text{op}} \text{Tr}(K) m^{1+\frac{2}{p}} \log m \right\}^{\frac{1}{2}} \\ + K_p \left\{ m \cdot \|K\|_{\text{op}} \mathbb{E}[\|\Delta\|_{\text{op}}] \right\}^{\frac{1}{2}}$$

Use that $\alpha^2 - \varepsilon_1 \alpha - \varepsilon_2 \leq 0 \Rightarrow \alpha \leq \frac{\varepsilon_1 + (\varepsilon_1^2 + 4\varepsilon_2)^{\frac{1}{2}}}{2}$

$$\leq (\varepsilon_1^2 + 4\varepsilon_2)^{\frac{1}{2}}$$

$$\mathbb{E}[\|\Delta\|_{\text{op}}] \leq K_p \left\{ m \|K\|_{\text{op}} + \left[\|K\|_{\text{op}} \text{Tr}(K) m^{1+\frac{2}{p}} \log m \right]^{\frac{1}{2}} \right\}$$

