

$H: \begin{cases} \varepsilon_i \text{ are independent} \\ \mathbb{E} \varepsilon_i = 0, \mathbb{E} \varepsilon_i^2 = \sigma^2 < +\infty \end{cases}$

excess risk of predictor  $f_\theta$  ( $x \mapsto \phi(x)^T \theta$ ):

$$R(f_\theta) - R^*$$

Prop. 1:  $R^* = \sigma^2$ ;  $R(\theta) - R^* = \|\theta - \theta^*\|_{\Sigma}^2$  (Σ) forget

expected excess risk =  $\mathbb{E}[R(\hat{\theta})] - R^*$   
 of the ERM: if  $\hat{\theta}$  is random

Prop 2:  $\mathbb{E}[R(\hat{\theta})] - R^* = \text{bias} + \text{variance}$

Prop 3: bias = 0 ; var =  $\frac{\sigma^2}{n} \sum_{j=1}^n 1$

Sig V random vector.  
 $(\text{Cov}(V))_{j,k} = E[(V_j - EV_j)(V_k - EV_k)]$   
 $:= \text{Var.}$

Prop. 4:  $ER(\hat{\theta}) - \mathcal{R}^* = \frac{\sigma^2 d}{n}$

↓
↓
↓
↓

expected risk of  $\hat{\theta}$       Bayes risk      noise level      sample size

→ dimension  
 → sample size

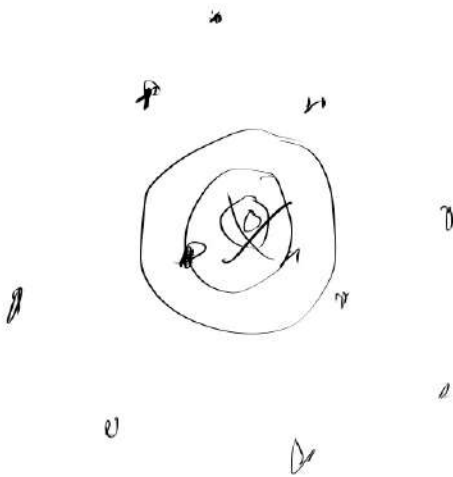
Q. no bias: in statistics, we estimate true quantities. An estimator is a random variable which depends only on  $X_1, \dots, X_n$

Def:  $\alpha^* \in \mathbb{R}$ ,  $\hat{\alpha}$  estimator of  $\alpha^*$ ,

bias( $\hat{\alpha}$ ) :=  $E[\hat{\alpha}] - \alpha^*$



zero bias "good," but not  
the end of the story



## RIDGE REGRESSION

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2}_{\text{you're lost}} + \lambda \|\theta\|^2 \right\}$$

$\lambda \|\theta\|^2$   
 regularization  
 penalization  
 penalty for

$\lambda \geq 0$ , fixed

in our notation:

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \underbrace{\frac{1}{n} \|y - X\theta\|^2}_{\text{we know this: CVX + smooth}} + \underbrace{\lambda \|\theta\|^2}_{\text{CVX + smooth } (\Psi)} \right\}$$

$F(\theta) := \text{red} + \text{blue}$  is convex, smooth  $\rightarrow \nabla F$ ?

$$\nabla F(\theta) = (\text{already done}) + \nabla(\lambda \|\theta\|^2)$$

$$\|\theta\|^2 = \theta^T \theta = \theta^T I \theta$$

$$\nabla(\theta^T I \theta) = (I + I^T) \theta = 2I\theta = 2\theta$$

$$\begin{aligned} \nabla F(\theta) &= \frac{2}{n} (\Phi^T \Phi \theta - \Phi^T y) + 2\theta \\ &= 2 \left( \left( \hat{\Sigma} + \lambda I \right) \theta - \Phi^T y \right) \end{aligned}$$

Claim 1.  $\nabla F = 0 \iff \left( \hat{\Sigma} + \lambda I \right) \theta = \Phi^T y$

equation in  $\theta$ . how to solve it?

Claim 2.  $\hat{\Sigma} + \lambda I$  is always invertible ( $\lambda > 0$ )

$$\Rightarrow \hat{\theta}_\lambda = \left( \hat{\Sigma} + \lambda I \right)^{-1} \Phi^T y$$

(original method.)

$$\hat{\Sigma} \in S_d^+$$

> per  $(\hat{\Sigma})$   
~~0~~  
 0

per absurdo:  $\hat{\Sigma} + \lambda I$  is not invertible.  
 $\det(\hat{\Sigma} + \lambda I) = 0$   
 $(-\lambda)$  is eigenvalue of  $\hat{\Sigma}$ .



Goal: expected excess risk of ridge

Q: bias-variance decomposition still holds?

→ proof true for any estimator; in particular for  $\hat{\theta}_\lambda$

$$\begin{aligned} \mathbb{E}[(d-\alpha)^2] &= \mathbb{E}[(d - \mathbb{E}d + \mathbb{E}d - \alpha^*)^2] \\ &= \mathbb{E}[(d - \mathbb{E}d)^2 + 2(d - \mathbb{E}d)(\mathbb{E}d - \alpha^*) + (\mathbb{E}d - \alpha^*)^2] \\ &= \cancel{\mathbb{E}[(d - \mathbb{E}d)^2]} + 2 \underbrace{\mathbb{E}[(d - \mathbb{E}d)(\mathbb{E}d - \alpha^*)]}_{\mathbb{E}d - \alpha^*} + \mathbb{E}[(\mathbb{E}d - \alpha^*)^2] \end{aligned}$$

To do: bias and variance of  $\hat{\theta}_\lambda$

Prop 5:  $\text{bias}(\hat{\theta}_\lambda) = -d(\hat{\Sigma} + dI)^{-1} \theta^*$

Rk:  $d=0$ :  $\hat{\theta}_\lambda = \hat{\theta}$ , bias = 0 as expected

Proof:  $\mathbb{E}\hat{\theta}_\lambda = \mathbb{E}\left[\frac{1}{n}(\hat{\Sigma} + dI)^{-1} \Phi^T Y\right]$

$$= \mathbb{E}\left[\frac{1}{n}(\hat{\Sigma} + dI)^{-1} \Phi^T (\Phi \theta^* + \varepsilon)\right]$$

$$= \mathbb{E}\left[\frac{1}{n}(\hat{\Sigma} + dI)^{-1} \underbrace{\Phi^T \Phi}_{\hat{\Sigma}} \theta^*\right] + \cancel{\mathbb{E}[\dots \varepsilon]} \quad \text{since } \mathbb{E}\varepsilon_i = 0$$

$$\mathbb{E}[\hat{\theta}_\lambda] = (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma} \theta^*$$

$$\mathbf{I} = (\hat{\Sigma} + \lambda \mathbf{I})^{-1} (\hat{\Sigma} + \lambda \mathbf{I}) = (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma} + \lambda (\hat{\Sigma} + \lambda \mathbf{I})^{-1}$$

$$(\hat{\Sigma} + \lambda \mathbf{I})^{-1} = \mathbf{I} - \lambda (\hat{\Sigma} + \lambda \mathbf{I})^{-1}$$

$$\Rightarrow \mathbb{E}[\hat{\theta}_\lambda] - \theta^* = -\lambda (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \theta^* \quad \square$$

$$\underbrace{\quad}_{\frac{1}{1 - \lambda \frac{\lambda}{\sigma^2}}}$$

$$\text{Prop 6: } \text{Var}(\hat{\theta}_\lambda) = \frac{\sigma^2}{n} \text{trace}(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda \mathbf{I})^{-2})$$

$$\begin{aligned} \text{Proof: } \hat{\theta}_\lambda - \mathbb{E}[\hat{\theta}_\lambda] &= (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \Phi^T \mathcal{Y} - (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \Phi^T \mathbb{E}[\mathcal{Y}] \\ &= (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \Phi^T \underbrace{\mathcal{Y} - \mathbb{E}[\mathcal{Y}]}_{\varepsilon} \end{aligned}$$

$$\begin{aligned} \varepsilon &\in \mathbb{R}^n \\ \varepsilon &= \text{trace}(\varepsilon) \end{aligned}$$

$$\mathbb{E}[\|\hat{\theta}_\lambda - \mathbb{E}[\hat{\theta}_\lambda]\|_{\frac{1}{2}}^2] = \mathbb{E}\left[\frac{1}{n} (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \Phi^T \varepsilon\right]_{\frac{1}{2}}^2$$

$$= \frac{1}{n^2} \mathbb{E}\left[\varepsilon^T (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \Phi^T \varepsilon\right] \in \mathbb{R}$$

$$= \frac{1}{n^2} \mathbb{E}\left[\text{trace}(\dots)\right]$$

$$\begin{aligned}
\text{Var} &= \frac{1}{n^2} \mathbb{E}[\text{trace}(\varepsilon^T \Phi (\hat{\Sigma} + dI)^{-1} \hat{\Sigma} (\hat{\Sigma} + dI)^{-1} \Phi^T \varepsilon)] \\
&\quad \left[ \text{cyclic property : } \text{trace}(AB) = \text{trace}(BA) \right] \\
&= \frac{1}{n^2} \mathbb{E}[\text{trace}(\varepsilon \varepsilon^T \frac{\hat{\Sigma}}{(\hat{\Sigma} + dI)^2})] \\
&= \frac{1}{n^2} \text{trace}(\mathbb{E}[\varepsilon \varepsilon^T] \frac{\hat{\Sigma}}{(\hat{\Sigma} + dI)^2}) \\
&= \frac{\sigma^2}{n^2} \text{trace}(\mathbb{E}[(\hat{\Sigma} + dI)^{-1} \hat{\Sigma} (\hat{\Sigma} + dI)^{-1}] \sigma^2) \\
&= \frac{\sigma^2}{n} \text{trace}(\hat{\Sigma} (\hat{\Sigma} + dI)^{-1} \hat{\Sigma} (\hat{\Sigma} + dI)^{-1}) \\
&= \frac{\sigma^2}{n} \text{trace}(\hat{\Sigma}^2 (\hat{\Sigma} + dI)^{-2}) \quad \square
\end{aligned}$$

Prop 7:  $\mathbb{E}R(\hat{\theta}_d) - R^* = \frac{\sigma^2}{n} \theta^{*T} (\hat{\Sigma} + dI)^{-2} \hat{\Sigma} \theta^*$   
 $+ \frac{\sigma^2}{n} \text{trace}(\hat{\Sigma}^2 (\hat{\Sigma} + dI)^{-2})$

$$\frac{d^2}{d^2 + 3}$$

$$\approx \frac{1}{d^2}$$

Rh:  $d=0$  has NO reason to be the best.