

Situation:  $Y = X\beta + \varepsilon$   $E(\varepsilon) = 0, \text{Var}(\varepsilon) = \sigma^2 I$   
 $N \times 1$   $N \times p$   $p \times 1$   $N \times 1$   
 when  $X^T X = I$  (and also  $(X^T X)^{-1} = I$ )  
 $p \times p$

(as before we have centered  $x$ 's and  $y$ , so  $\beta_0$  is not needed)

## 1) least squares (OLS)

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T Y = X^T Y$$

$$E(\hat{\beta}_{LS}) = X^T E(Y) = \underbrace{X^T X}_I \beta = \beta \quad (\text{we knew that})$$

$$\text{Cov}(\hat{\beta}_{LS}) = X^T \text{Cov}(Y) X = (X^T X) \sigma^2 I = \sigma^2 I$$

$$\begin{aligned} \text{MSE}(\hat{\beta}_{LS}) &= \underbrace{(E(\hat{\beta}_{LS}) - \beta)^T (E(\hat{\beta}_{LS}) - \beta)}_0 + \underbrace{\text{tr}(\text{Var}(\hat{\beta}_{LS}))}_{\sigma^2 I \text{ } p \times p} \\ &\quad \text{general formula} \\ &= p \sigma^2 \end{aligned}$$

## 2) Ridge regression

$$\begin{aligned} \hat{\beta}(\lambda) &= (\underbrace{X^T X}_I + \lambda I)^{-1} X^T Y = \begin{bmatrix} 1+\lambda & & 0 \\ & \ddots & \\ 0 & & 1+\lambda \end{bmatrix}^{-1} X^T Y \\ &= \left( \frac{1}{1+\lambda} \right) \cdot X^T Y \\ &= (1+\lambda)^{-1} \hat{\beta}_{LS} \end{aligned}$$

$p \times p$   $p \times N$   $N \times 1$   
 $p \times 1$

Can use results for  $\hat{\beta}_{LS}$  to get  $E, \text{Cov}, \text{MSE}$  for  $\hat{\beta}(\lambda)$

$$E(\hat{\beta}(\lambda)) = (1+\lambda)^{-1} \underbrace{E(\hat{\beta}_{LS})}_{\beta} = (1+\lambda)^{-1} \beta$$

$$\text{Cov}(\hat{\beta}(\lambda)) = (1+\lambda)^{-1} \text{Cov}(\hat{\beta}_{LS}) (1+\lambda)^{-1}$$

$$= (1+\lambda)^{-2} \cdot \sigma^2 \mathbf{I} = \frac{\sigma^2}{(1+\lambda)^2} \frac{\mathbf{I}}{p \times p}$$

$$\text{MSE}(\hat{\beta}(\lambda)) = ((1+\lambda)^{-1} \beta - \beta)^T ((1+\lambda)^{-1} \beta - \beta) + \text{tr} \left( (1+\lambda)^{-2} \sigma^2 \frac{\mathbf{I}}{p \times p} \right)$$

$$= \underbrace{\left( \frac{1}{1+\lambda} - 1 \right)}_{\frac{1-1-\lambda}{1+\lambda}} \beta^T \beta \left( \frac{1}{1+\lambda} - 1 \right) + (1+\lambda)^{-2} p \sigma^2$$

$$= \lambda^2 (1+\lambda)^{-2} \beta^T \beta + p (1+\lambda)^{-2} \sigma^2$$

When is  $\text{MSE}(\hat{\beta}(\lambda))$  minimal wrt  $\lambda$ ?  $\frac{\lambda^2}{(1+\lambda)^2}$

$$\frac{d \text{MSE}(\hat{\beta}(\lambda))}{d\lambda} = \frac{2\lambda \cdot (1+\lambda)^2 - \lambda^2 \cdot 2(1+\lambda)}{(1+\lambda)^4} \beta^T \beta$$

$$- 2p (1+\lambda)^{-3} \sigma^2 = 0$$

$$\frac{2\lambda(1+\lambda) - \lambda}{(1+\lambda)^3} \beta^T \beta - \frac{2p \sigma^2}{(1+\lambda)^3} = 0$$

$$\lambda \beta^T \beta - p \sigma^2 = 0 \Leftrightarrow \underline{\underline{\lambda = \frac{p \sigma^2}{\beta^T \beta}}}$$

This means that  $\lambda_{\text{opt}} = \sigma^2 \cdot \frac{p}{p^T p}$

Is this a minimum? Yes - could have checked!

See plots in L7.html under solutions

Theobald (1974): There exists  $\lambda > 0$  such that  $\text{MSE}(\hat{\beta}(\lambda)) < \text{MSE}(\hat{\beta}(0))$

Proof in WNVW show  $\lambda < 2\sigma^2(p^T p)^{-1}$  is sufficient for

$\text{MSE}(\hat{\beta}(\lambda)) < \text{MSE}(\hat{\beta}(0))$ , but for us the optimal  $\lambda$  is larger than this upper limit.