# Lecture 15:

<u>OLS</u>: $\hat{\beta}^{OLS} = (X^TX)^{-1}X^TY$

$X = UDV^T$  where  $U, V$ are orthogonal

$$U^TU = I, \quad V^TV = I$$

$$D = \left[\begin{array}{cc|c} \sigma_1 & & 0 \\ & \ddots & \sigma_r & \\ \hline & 0 & & 0 \end{array}\right]$$

$$\boxed{(AB)^T = B^T A^T}$$

$$X^TX = (UDV^T)^T UDV^T = VD^TU^TUDV^T$$

$$= VD^TDV^T \qquad D = \left[\begin{array}{c|c} D_* & 0 \\ \hline 0 & 0 \end{array}\right]$$

$$\overbrace{(X^TX)^{-1}}^{} = V(D^TD)^{-1}V^T$$

$$D^TD = \left[\begin{array}{c|c} D_*^T & 0 \\ \hline 0 & 0 \end{array}\right]\left[\begin{array}{c|c} D_* & 0 \\ \hline 0 & 0 \end{array}\right]$$

<span style="color:red">assume rank $X = P = \#$ Cols</span>

<span style="color:red">then $D^TD = D_*^2$</span>

$$= \left[\begin{array}{c|c} D_*^T D_* & 0 \\ \hline 0 & 0 \end{array}\right]$$

$$= VD_*^{-2}V^T$$

<span style="color:red">no zero blocks</span>

$D_* = \text{diag}(\sigma_1, \ldots, \sigma_r)$

$= D_*^T$

$$D_* = \left[\begin{array}{ccc} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{array}\right]$$

$$= \left[\begin{array}{c|c} D_*^2 & 0 \\ \hline 0 & 0 \end{array}\right]_{N \times P}$$

then $D_*^2 = \left[\begin{array}{ccc} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_r^2 \end{array}\right]$ and $D_*^{-2} = \left[\begin{array}{ccc} 1/\sigma_1^2 & & \\ & \ddots & \\ & & 1/\sigma_r^2 \end{array}\right]$

$$\tilde{Y}^{OLS} = X\hat{\beta}^{GLS} = \overbrace{UDV^T}^{X} \overbrace{VD_*^{-2}V^T}^{(X^TX)^{-1}} \overbrace{VD U^T Y}^{X^T}$$

$$= UDD_*^{-2}D^TU^TY \qquad D = \begin{bmatrix} D_* \\ \hline 0 \end{bmatrix}$$

$$\begin{bmatrix} I & | & 0 \\ \hline 0 & | & 0 \end{bmatrix}$$

If I ignore $DD_*^{-2}D$ then $\hat{Y} = UU^TY$

$$UU^T = \sum_{j=1}^{N} U_j U_j^T \quad \text{when } U_j = j^{th} \text{ col of } U$$

I can't actually ignore $DD_*^{-2}D$ so instead I get

$$U\begin{bmatrix} I & | & 0 \\ \hline 0 & | & 0 \end{bmatrix}U^T = \sum_{j=1}^{P} U_j U_j^T \qquad (P \leq N)$$

So

$$\boxed{\hat{Y}^{OLS} = \cdots = \sum_{j=1}^{P} U_j U_j^T Y}$$

proj. onto $U_j$

① project $Y$ onto $U_j$ : $U_j U_j^T Y$

② Sum up these components

# Ridge:

$$\hat{y}^{\,ridge} = X\hat{\beta}^{\,ridge}$$

$$= X\left(X^TX + \lambda I\right)^{-1}X^TY$$

$$= UDV^T\left(VD_*^2V^T + \lambda I\right)^{-1}VD^TU^TY$$

$$= UD\left(V^T(VD_*^2V^T + \lambda I)V\right)^{-1}D^TU^TY$$

$$= UD\left(\cancel{V^T}VD_*^2\cancel{V^T}V + \lambda\cancel{V^T}V\right)^{-1}D^TU^TY$$

$$= UD\left(\underbrace{D_*^2 + \lambda I}\right)^{-1}D^TU^TY$$

$$D_*^2 + \lambda I = diag\left(\sigma_i^2 + \lambda\right)$$

$$\left(D_*^2 + \lambda I\right)^{-1} = diag\left(\frac{1}{\sigma_i^2 + \lambda}\right)$$

$$D\left(D_*^2 + \lambda I\right)^{-1}D^T = \begin{bmatrix} \sigma_i^2/\sigma_i^2 + \lambda & & 0 & \\ & \ddots & & 0 \\ 0 & & \sigma_r^2/\sigma_i^2 + \lambda & 0 \\ \hline & 0 & & 0 \end{bmatrix}$$

and so similar to previously

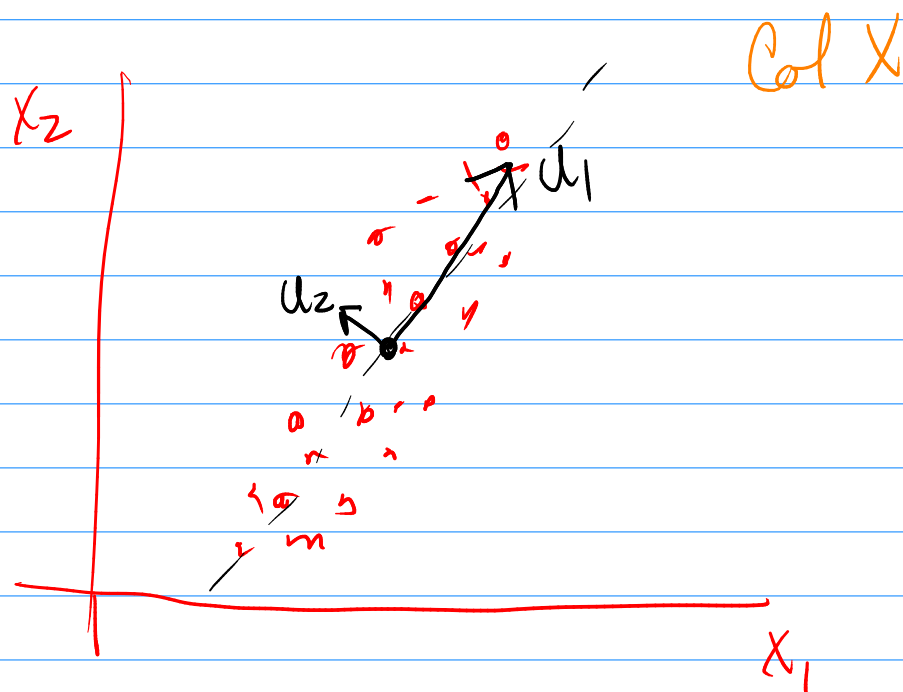$$\boxed{\hat{y}^{\,ridge} = \sum_{j=1}^{r}\left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right)U_jU_j^TY}$$

① proj. $Y$ onto $U_j$ : $U_jU_j^TY$

② weighting by $\sigma_i^2/\sigma_i^2 + \lambda < 1$

③ sum up these components

$X = UDV^T$

$(AB)^{-1} = B^{-1}A^{-1}$

(1) If $\lambda = 0$ we get OLS
   b/c $\sigma_i^2 / \sigma_i^2 + \lambda = 1$

(2) Large $\lambda$ really "shrinks" the contribution
   of the $j^{th}$ component

(3) Shrinks contrib. of smaller $\sigma_i$ faster



---

Degrees of Freedom

For OLS: $df = P$

For Ridge: $df = \sum_{j=1}^{r} \frac{\sigma_i^2}{\sigma_i^2 + \lambda} < P$

$df \to 0$ as $\lambda \to \infty$
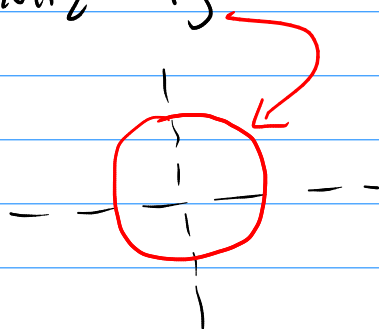
<span style="color:red">As we increase $\lambda$, flexibility decreases, my bias increases, var dec.</span>

---

Aside: Norms

Ex. Euclidean Norm:

$$\|x\|_2 = \sqrt{\sum_{i=1}^{P} x_i^2}$$

Consider: $\{x \mid \|x\|_2 = 1\}$

Can generalize

q-norm: $\|x\|_q = \left(\sum_{i=1}^{P} |x_i|^q\right)^{1/q}$

notice when $q=2$ I get Euclidean Norm.

When $q=1$ I get the L1-norm

$$\|x\|_1 = \sum_{i=1}^{P} |x_i|$$

consider $\{x \mid \|x\|_1 = 1\}$

<span style="color:red">(0,1)</span>
<span style="color:red">(1,0)</span>
<span style="color:red">(-1,0)</span>
<span style="color:red">(0,-1)</span>