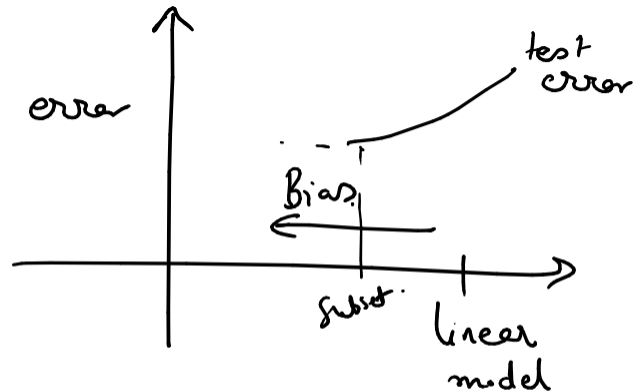# Beyond Subset selection:

↓

beyond retaining or excluding a coordinate

↳ Ridge & LASSO

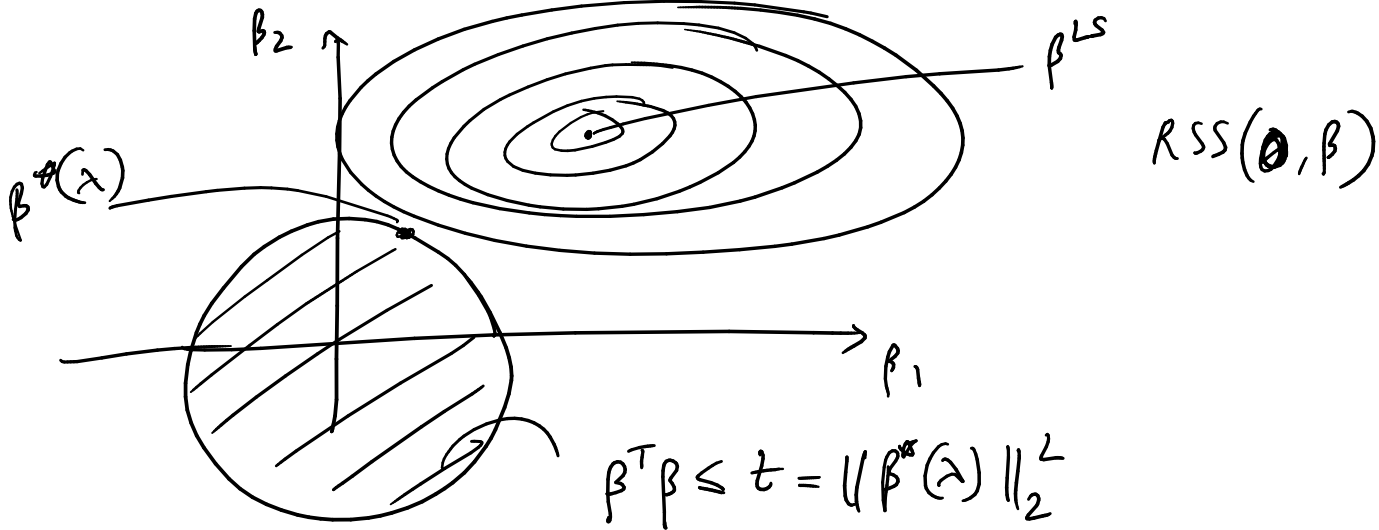## Ridge

→ loss should be low

→ coefficients $(\hat{\beta}_j)$ be low in magnitude.

$$RSS(\lambda, \beta) = \underbrace{(Y - X\beta)^T (Y - X\beta)}_{\beta^*(\lambda)} + \underbrace{\lambda \beta^T \beta}_{} \qquad \text{where } \lambda > 0$$

---

$$\min_{\beta} (Y - X\beta)^T (Y - X\beta)$$

$$\text{st} \qquad \beta^T \beta \leq \underbrace{\| \beta^*(\lambda) \|_2^2}_{} = t$$

$\beta_2$

$\beta^{LS}$

$\beta^{\#}(\lambda)$

$RSS(\mathbf{0}, \beta)$

$\beta_1$

$$\beta^T \beta \leq t = \| \beta^{\#}(\lambda) \|_2^2$$

$$[\beta_1 \ \beta_2] \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

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

$$\beta^{ridge} = (X^T X + \underline{\underline{\lambda I}})^{-1} X^T Y$$

$$\underbrace{\overset{Y^T Y}{a} - b - c + \beta^T X^T X \beta}_{\beta^T (X^T X + \lambda I) \beta} + \beta^T (\lambda I) \beta$$

$$\beta^T (X^T X + \underbrace{\lambda I}_{D}) \beta$$

Goal: Relate MLE to Ridge regression. $\longrightarrow$ Lec 6 ☐

MLE eg $\quad z_1, \ldots z_N \sim N(\mu, 1)$

$$\text{Likelihood function} = P_\mu(Z_1 = z_1, Z_2 = z_2 \ldots, Z_N = z_N)$$

$$= \prod_{i=1}^{N} P_\mu(Z_i = z_i)$$

$$LL = \sum_{i=1}^{N} \log P_\mu(Z_i = z_i)$$

$$P(Z_{i=3i}) = \frac{1}{\sqrt{2\pi}\,1} \exp\left(-\frac{(3i-\mu)^2}{2\cdot 1^2}\right)$$

$$\log \quad \text{''} \quad = \quad -\frac{(3i-\mu)^2}{2} - \log\sqrt{2\pi}$$

$$\max_{\mu}(LL) = -\frac{1}{2}\sum_{i=1}^{N}(3i-\mu)^2 - N\log\sqrt{2\pi}$$

$$\frac{\partial}{\partial\mu}(\text{''}) = \sum_{i=1}^{N}(3i-\mu) = 0$$

$$\boxed{\mu = \frac{1}{N}\sum_{i=1}^{N}3i}$$

## positive semi-definite matrix

$X^T X$ is ↗ if smallest eigenvalue is $\geq 0$

$A_{p \times p}$    $\underline{\underline{\lambda_1}} \cdots \underline{\underline{\lambda_p}}$

$\lambda = v^T \boxed{A} v$  for $v$ being eigenvector.

$\longrightarrow \geq 0$

$v^T A \underline{\underline{v}} = v^T \underline{\underline{\lambda v}} = \underset{\underline{\underline{\geq 0}}}{\|v\|_2^2} \underset{\geq 0}{\lambda} \quad X^T X$

$$\lambda = \underbrace{v^T X^T X v}$$

$$= (Xv)^T (Xv)$$

$$= \|Xv\|_2^2 \geq 0$$

## Interpreting Ridge Regression:

Fact: $svd(X) = \underset{N \times p \;\; p \times p \;\;\;\; p \times p}{U \; D \; (V^T)}$

$$\boxed{y_i - \hat{\beta}^T x_i}$$

$A = U D V^T$

$|d_1| > |d_2| > \dots$

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

$$= U D V^T (V D^2 V^T)^{-1} V D U^T Y$$

$$= U D D^{-2} D U^T Y$$

$X^T X = V D^2 V^T$

$= V D U^T U D V^T$

$$\underset{N\times 1}{X\hat{\beta}} = \underset{N\times p}{U} \cdot \underbrace{\underset{p\times N}{U^T} \underset{N\times 1}{Y}}$$

$$= Uc$$
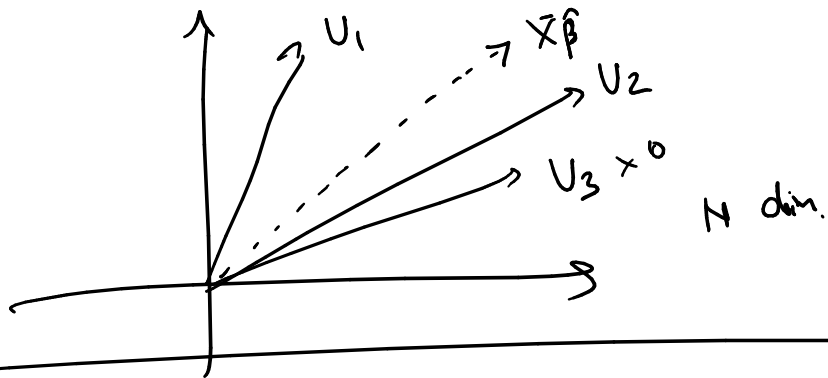
$$= \sum_{j=1}^{p} c_j \cdot U_j$$

$$\begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{bmatrix} = c$$

$$U = \begin{bmatrix} U_1 & \cdots & U_p \\ | & & | \end{bmatrix}_{N\times p}$$

$U_1$

$X\hat{\beta}$

$U_2$

$U_3 \times 0$

N dim.

<u>Ridge</u>

$$X\beta^{ridge} = UD(\underline{\underline{D^2 + \lambda I})^{-1}}DU^TY$$

$$= \sum_{j=1}^{P} \left(\frac{d_j^2}{d_j'^2 + \lambda}\right)(U_j^TY) \cdot U_j$$

What does $d_j^2$ mean?

$$X = U D V^T_{p \times p}$$

$$\frac{1}{N} X^T X_{p \times p}$$

$V$ : the Columns are Called pricipal Component directions.

Centering :
$X$
· take Column means
· Subtract from Column entries

$$X^T \underset{p \times p}{X} = \underset{p \times p}{V} \underset{p \times p}{D^2} \underset{p \times p}{V^T} \qquad : \text{eigen decomposition.}$$

$$= \sum_{j=1}^{p} d_j^2 \, v_j v_j^T \qquad v_1 \begin{bmatrix} d_1^2 & & 0 \\ & \ddots & \\ 0 & & d_p^2 \end{bmatrix}_{p \times p}$$

Fact

$$d_1^2 = Var(X v_1)$$

$$\begin{array}{ccccc} 1 & 2 & \underline{3} & 4 & 5 \\ -2 & -1 & 0 & 1 & 2 \end{array}$$

$$\underline{\underline{\text{Var}\left(\bar{X}v_1\right)}} = v_1^T \bar{X}^T \bar{X} v_1$$

$$= v_1^T \left( \underline{\underline{\sum_{j=1}^{p} d_j^2 \, v_j v_j^T}} \right) v_1$$

$$= \sum_{j=1}^{p} d_j^2 \left( v_1^T v_j \right) \left( v_j^T v_1 \right)$$

$$= \underline{\underline{d_1^2 \cdot 1 \cdot 1}}$$

$$\underbrace{\bar{X}v_1}_{N \times 1 \text{ dim}}$$

$$\text{Svd}(X^T X)$$
$$= V D^2 V^T$$
$$= \text{pca}(\bar{X})$$
$$= \text{eigen decom}$$
$$\text{position}$$
$$(X^T X)$$

Missing argument for $\beta^*(\lambda)$:

1. We know $\beta^*(\lambda)$ is feasible. i.e, $\beta^*(\lambda)^T \beta^*(\lambda) \leq \|\beta^*(\lambda)\|_2^2$ by definition.

. Also for any other $\beta$ we have

$$Y - X\beta^*(\lambda) \Big)^T \Big( Y - X\beta^*(\lambda) \Big) + \lambda \|\beta^*(\lambda)\|_2^2$$

$$\leq (Y - X\beta)^T (Y - X\beta) + \lambda \|\beta\|_2^2$$

$$Y - X\beta^*(\lambda)^T (Y - X\beta^*(\lambda)) \leq (Y - X\beta)^T (Y - X\beta) + \lambda \underbrace{\Big( \|\beta\|_2^2 - \|\beta^*(\lambda)\|_2^2 \Big)}_{< 0}$$

$\therefore$ $\beta^*(\lambda)$ is the minimizer for

$$\min \ (Y - \bar{X}\beta)^T (Y - \bar{X}\beta)$$

Subject to $\beta^T \beta \leq \|\beta^*(\lambda)\|_2^2$.