Ridge regression $\hat{\beta}_\lambda^R = (X^T X + \lambda I)^{-1} X^T y$

**Theorem 1** Suppose rank$(X) = p$ and $\lambda > 0$ is suff. small (depending on $\beta^0, \sigma^2$)

Have $E(\hat{\beta}^{OLS} - \beta^0)(\hat{\beta}^{OLS} - \beta^0)^T - E(\hat{\beta}_\lambda^R - \beta^0)(\hat{\beta}_\lambda^R - \beta^0)^T \overset{(*)}{\phantom{=}}$ is pos. def.

**Proof** Bias $E\hat{\beta} - \beta^0 = (X^T X + \lambda I)^{-1} X^T X \beta^0 - \beta^0$

$\hat{\beta}_\lambda^R \nearrow$
$\qquad = (X^T X + \lambda I)^{-1} \{ X^T X \beta^0 - X^T X \beta^0 - \lambda \beta^0 \}$

$\qquad = -\lambda (X^T X + \lambda I)^{-1} \beta^0$

$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}$

$(X^T X)^{-1} = (X^T X + \lambda I)^{-1} \underbrace{(X^T X + \lambda I)(X^T X)^{-1}(X^T X + \lambda I)}(X^T X + \lambda I)^{-1}$
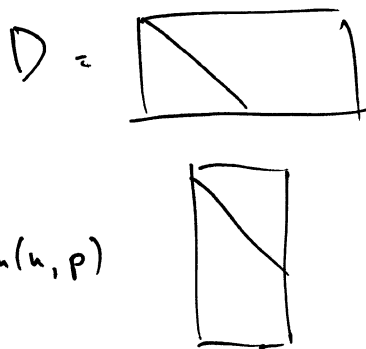
$$X^T X + 2\lambda I + \lambda^2 (X^T X)^{-1}$$

(*) pos def iff

$\sigma^2 (X^T X + 2\lambda I + \lambda^2 (X^T X)^{-1}) - \sigma^2 X^T X - \lambda^2 \beta^0 \beta^{0T}$ is pos def

$\lambda (2\sigma^2 I + \sigma^2 \lambda (X^T X)^{-1} - \lambda \beta^0 \beta^{0T})$

But this if pos def for $0 < \lambda < \dfrac{2\sigma^2}{\|\beta^0\|_2^2}$ $\square$

**The singular value decomposition (SVD)**

Can factorise any $X \in \mathbb{R}^{n \times p}$ into its SVD

$X = \underset{\substack{n \times n \\ \text{orthogonal}}}{U} \quad \underset{n \times p}{D} \quad \underset{\substack{p \times p \\ \text{orthogonal}}}{V^T}$
$\qquad\qquad\qquad\qquad\qquad D = $ 

$D$ has $D_{11} \geq D_{22} \geq \dots \geq D_{mm} \geq 0$ where $m = \min(n, p)$
and all other entries of $D$ are 0.

Alternative form (when $n > p$) 
   Replace $U$ by its first $p$ cols
   $\qquad D$ by its first $p$ rows
Have $X = \underset{n \times p}{U} \underset{p \times p}{D} \underset{p \times p}{V^T}$ $\qquad U$ has orthonormal cols so $U^T U = I$ but $U U^T \neq I$.

Suppose $n > p$. Fitted vals of ridge regression

$$X\hat{\beta}_\lambda^R = X(X^TX + \lambda I)^{-1} X^T Y$$

$$= UDV^T(VD^2V^T + \lambda I)^{-1} VD U^T Y$$

Now $VD^2V^T + \lambda I = V(D^2 + \lambda I)V^T$

and $\{V(D^2 + \lambda I)V^T\}^{-1} = V(D^2 + \lambda I)^{-1}V^T$

$$X\hat{\beta}_\lambda^R = UD(D^2 + \lambda I)^{-1} DU^T Y$$

$$= \sum_{j=1}^{p} U_j \frac{D_{jj}^2}{D_{jj}^2 + \lambda} U_j^T Y$$

OLS fitted vals

$$X\hat{\beta}^{OLS} = \sum_{j=1}^{p} U_j U_j^T Y \qquad \text{(provided } D_{pp} > 0 \text{ so } X \text{ has full col rank)}$$

Both OLS and ridge regression compute coords of $Y$ w.r.t. the cols of $U$.
Ridge shrinks those coords by factor $D_{jj}^2/(D_{jj}^2 + \lambda)$ whilst OLS leaves these unchanged.

Cols of $U$ correspond to the principal components of $X$.
Take $u \in \mathbb{R}^p$, $\|u\|_2^2 = 1$. Sample variance $Xu \in \mathbb{R}^n$ is $\frac{1}{n}\|Xu\|_2^2 = \frac{1}{n} u^T X^T X u$

$$= \frac{1}{n} u^T V D^2 V^T u$$

Let $W = V^T u$. Then $\|w\|_2 = 1$

$$\frac{1}{n} u^T V D^2 V^T u = \frac{1}{n} W^T D^2 W = \frac{1}{n} \sum_j w_j^2 D_{jj}^2 \leq \frac{1}{n} D_{11}^2$$

with equality when $W = \begin{pmatrix} 1 \\ 0 \\ \vdots \end{pmatrix}$, so $u = V_1$.

Thus $V_1$ gives the coeffs of the linear combination of cols of $X$ with largest sample variance, when coeffs are constrained to have $l_2$-norm 1.

$$XV_1 = U_1 D_{11} \text{ is the 1st principal component of } X.$$
$$UDV^T V_1$$

Subsequent principal components obey the same optimality conditions subject to being normal to all prior principal components. Can show $r^{th}$ principal comp is $D_{rr} U_r$.

## Conclusion:

Ridge regression works best when most of the signal $EY$ is in the direction of the large principal compt of $X$.