# Ridge Regression

Ridge penalizes the sq. err. to avoid "large" $\beta$:

$$\hat{\beta}^{(ridge)} = \underset{\beta}{\arg\min} \ \|Y - X\beta\|_2^2 + \lambda\|\beta\|_2^2$$

$\lambda \geq 0$

By adding $\lambda\|\beta\|_2^2$ if the entries of $\beta$ become large so does this penalty

$\lambda$ = penalty strength

$\lambda = 0$ gives OLS $\hat{\beta}$

$\lambda \to \infty$ we get $\hat{\beta}^{(ridge)} \to 0$

Typically, we don't include $\beta_0$ in the

Typically, we don't include $\beta_0$ in the
    penalty
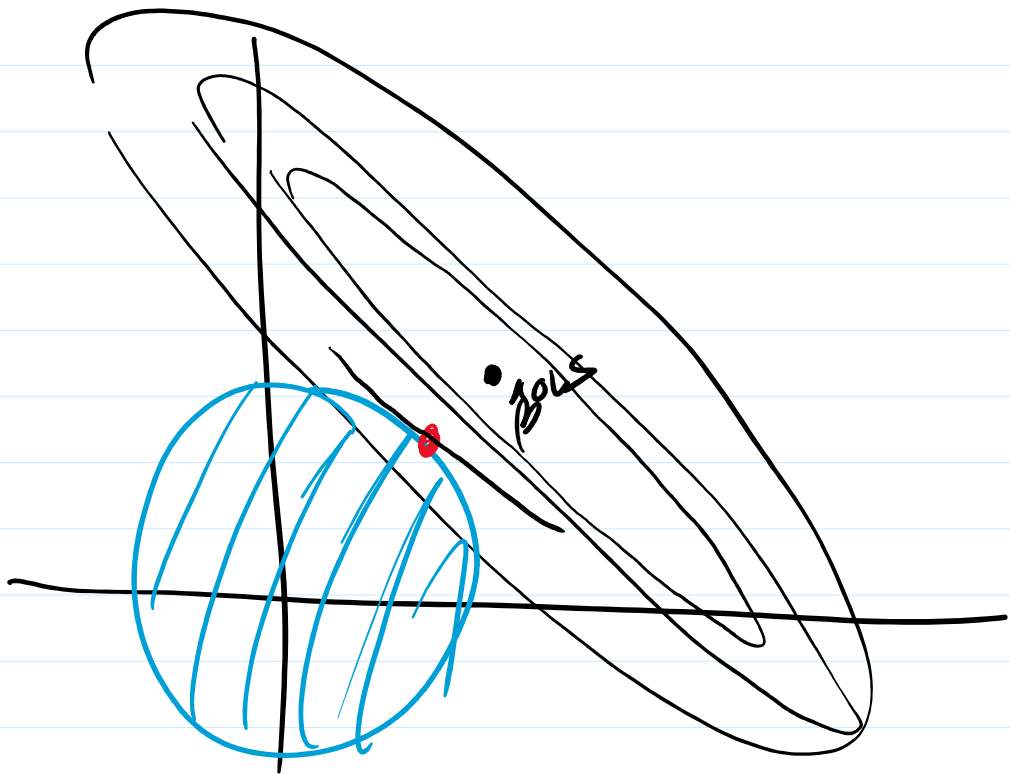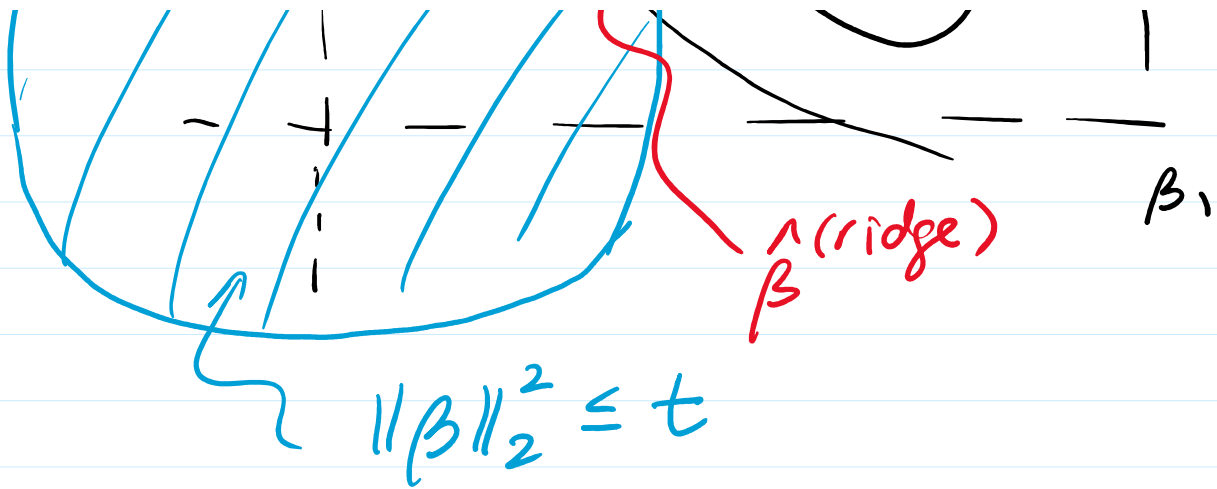
Often, we standardize vars before ridge

Also, typically choose $\lambda$ via x-val.

- - - - - - - - - - - - -

Second interpretation: ridge is equivalent
to

$$\hat{\beta}^{(ridge)} = \underset{\beta}{\arg\min} \|Y - X\beta\|_2^2$$

$$\text{s.t. } \|\beta\|_2^2 \leq t$$

$t$

1-1 corresp
w/ $\lambda$

Ex.

$\beta_2$

$\hat{\beta}$

$L(\beta)$

$\hat{\beta}$ (ridge)

$\beta_1$

$\|\beta\|_2^2 \leq t$

$\bullet \ \hat{\beta}_{OLS}$

How do we get $\hat{\beta}$ (ridge)?

Because $\lambda\|\beta\|^2$ is quadratic and so is

$\left( \quad \quad \sum_{i=1}^{p} \beta_i^2 \right.$

$$\longrightarrow \sum_{j=1}^{p} \beta_j^2$$

$\|Y - X\beta\|_2^2$ then there is a closed form soln for $\hat{\beta}^{(ridge)}$.

OLS: $\frac{\partial L}{\partial \beta} = 0 \Rightarrow$ Solve $(X^T X)\beta = X^T Y$

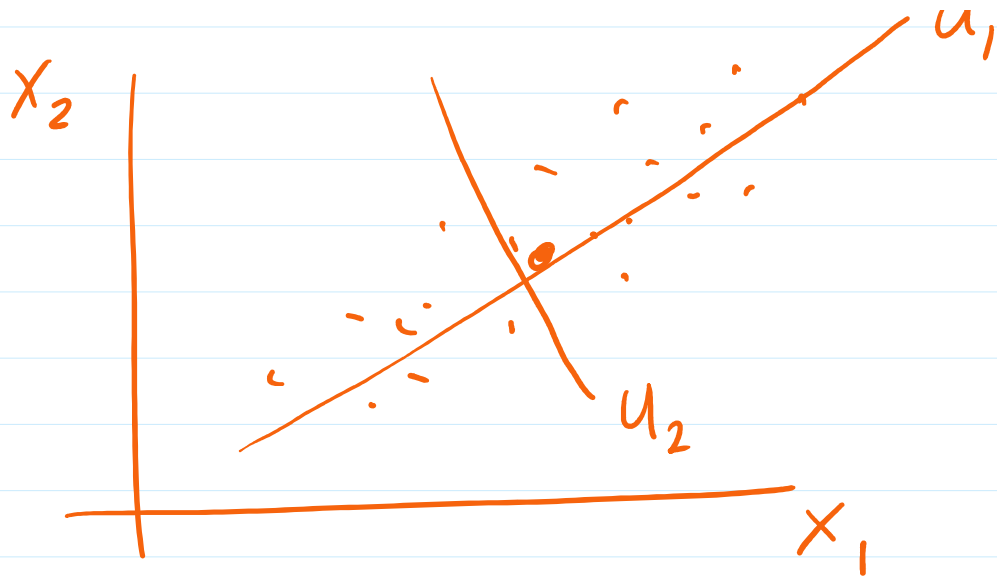Ridge: $\frac{\partial L}{\partial \beta} = 0 \Rightarrow$ Solve $(X^T X + \lambda I)\beta = X^T Y$

for $\lambda > 0$ $\quad X^T X + \lambda I$ is invertible

So $\boxed{\hat{\beta}^{(ridge)} = (X^T X + \lambda I)^{-1} X^T Y .}$

For OLS the sens. of $\hat{\beta}^{(OLS)}$ depended on

$$K(X^T X)$$

For ridge the sens. of $\hat{\beta}^{(ridge)}$ deps.

$$(X^T X + \lambda I)$$

$$K(X^TX + \lambda I)^{-1}$$



For OLS we showed that if

$$X = UDV^T \quad (\text{full rank})$$

then

$$\hat{Y} = X\hat{\beta}^{(OLS)} = \sum_{j=1}^{P} u_j u_j^T Y$$

① proj. $Y$ onto $u_j$

② Sum up these contribs.

$$X_{\cdot} \qquad \qquad \qquad / u_1$$

For ridge : can show that

$$\hat{Y} = X\hat{\beta}^{\wedge(ridge)} = \sum_{j=1}^{P} \left( \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \right) u_j u_j^T Y$$

① project $Y$ onto $u_j$

② rescale each by $\dfrac{\sigma_j^2}{\sigma_j^2 + \lambda} \leq 1$

③ sum up contribs

}

scale comps assoc. w/ smaller $\sigma$: more towards

$\overline{\sigma_j}$ more towards zero

## Degrees of Freedom

For OLS: $df = P$ if $rank(X) = P$

For ridge: $df = \sum_{j=1}^{P} \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \leq P$
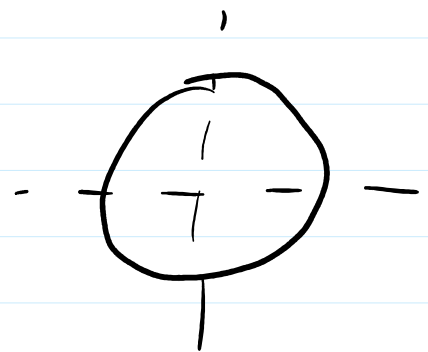
$$df \to 0 \quad as \quad \lambda \to \infty$$
$$df \to P \quad as \quad \lambda \to 0$$

## Norms

Euclidean Norm:
$$\|X\|_2 = \sqrt{\sum_{j=1}^{P} x_j^2}$$

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

Can generalize:

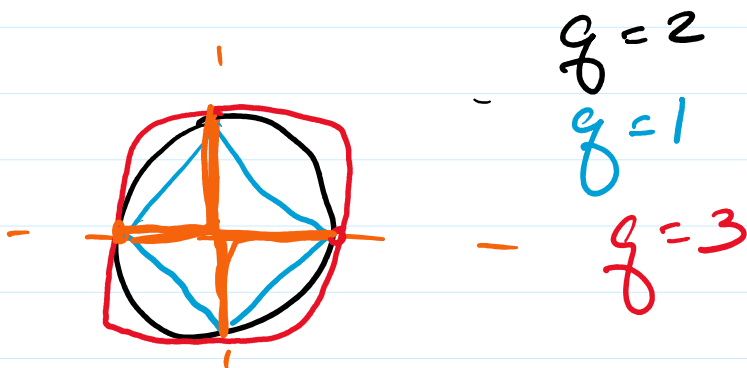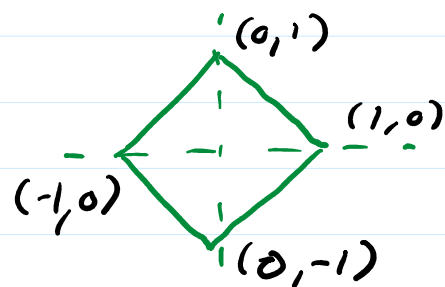$q$-norm: $\|x\|_q = \left( \sum_{j=1}^{P} |x_j|^q \right)^{1/q}$

$$q\text{-norm}: \|x\|_q = \left( \sum_{j=1} |x_j|^q \right)$$

When $q = 2$, I get euclidean norm
$$(L2 \text{ norm})$$

If $q = 1$, get L1 norm

$$\|x\|_1 = \sum_{j=1}^{p} |x_j|$$

Consider $\{x : \|x\|_1 = 1\}$

(0,1)

(1,0)

(-1,0)

(0,-1)

$q = 2$

$q = 1$

$q = 3$

as $q \to \infty$ I get $\|x\|_q \to \max_j |x_j|$
$$= \|x\|_\infty$$

$q \to 0$ I get $\|x\|_q \to \#$ of non-zero

$$q \to 0 \quad I \text{ set } \|X\|_q \to \# \text{ of non-zero elements}$$
$$= \|X\|_0$$

Variable selection is like zeroing out some of my $\hat{\beta}s$:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \cdots$$

$\hookleftarrow$ set $\hat{\beta}_2 = 0$

now: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \cdots$