# 6 Regularization

Recall the multiple regression model:

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon$$

## Note

- If $p \gg 0$ but $p < n$ then the OLS estimator

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad \text{exhibits high variance}$$

- If $p > n$ then the OLS estimates are not identifiable b/c

$$(X^T X)\beta = X^T y$$

<span style="color:red">rank deficient $\Rightarrow$ not invertible</span>

So there are infinitely many solutions.

- Do subset selection

- Include all predictors, but regularize their effects to shrink them toward zero.

## Motivation

Given samples $(X_1, Y_1), \ldots, (X_n, Y_n)$

OLS minimizes

$$(Y - X\beta)^T (Y - X\beta)$$

regularization adds a penalty ⟹

$$(Y - X\beta)^T (Y - X\beta) + \lambda P(\beta)$$

where $P(\beta)$ is a penalty/regularization term that grows with the size of $\beta$, and shrinks to zero when $\beta = 0$

$\lambda \geq 0$ is a smoothing | shrinkage | complexity | regularization parameter.

**Note** Why does $P(\vec{\beta})$ control the size of the model?

- Case 1: $\vec{\beta} = \underline{0}$, so $P(\vec{\beta}) = 0 \Rightarrow Y$ does not depend on $\underline{x}$ at all, so we get a small model

- Case 2: $\vec{\beta}$ big in every feature, $P(\vec{\beta}) \gg 0$

  $\Rightarrow Y$ depends strongly on all features.

## 6.1 Ridge Regression

$\beta_0$ only measures the avg value of $Y$, so should not be penalized.

If we use __centered features__     $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$

$$Y_i = \beta_0 + \beta_1 (x_i - \bar{x}) + \varepsilon_i \qquad i = 1, \ldots, n$$

What is $\hat{\beta_0}$?

$$\frac{\partial}{\partial \beta_0} \sum_{i=1}^{n} \left( Y_i - \beta_0 - \beta_1 (x_i - \bar{x}) \right)^2$$

$$= -2 \sum_{i=1}^{n} (Y_i - \beta_0) + 2\beta_1 \sum_{i=1}^{n} (x_i - \bar{x})$$

$$= -2n\bar{Y} + 2n\beta_0$$

$$\overset{set}{\Rightarrow} \boxed{\bar{Y} = \hat{\beta_0}} \text{ is OLS.}$$

(red annotations, right side:)

$$\longrightarrow 0 \checkmark$$

$$\sum (x_i - \bar{x})$$
$$= \sum x_i - \sum \bar{x}$$
$$= \sum x_i - n\bar{x}$$
$$= \sum x_i - \sum x_i = 0$$

$\Rightarrow$ Throughout the section we will center features + response

$\boxed{\text{Note}}$ "The basic idea behind ridge regression is

$$P(\vec{\beta}) = \|\vec{\beta}\|_2^2 = \vec{\beta}^T \vec{\beta} = \sum_{j=1}^{P} \beta_j^2$$

where now $\quad \vec{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_P \end{pmatrix}$ $\longrightarrow$ <span style="color:red">no $\beta_0$ term b/c we'll use centered features.</span>

$\boxed{\text{Problem}}$ $\quad P = 2 \qquad x_1 = $ budget of movie in \$s ($\sim 1000000$s ish)

$\qquad\qquad\qquad\qquad x_2 = $ rating of movie $\qquad (\sim 1\text{-}10$ ish$)$

$$\|\vec{\beta}\|_2^2 = \beta_1^2 + \beta_2^2 \quad \text{units?}$$

Nonsensical $\Rightarrow$ need to remove the units of features.

## Warning / Note / Convention

For the rest of the chapter, we will assume:

- $\hat{\beta}_0 = \bar{Y}$, $\quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}$ $\quad$ p features, no $\beta_0$

- observations are centered: $\quad Y_i \rightarrow Y_i - \bar{Y}$

- Features have been centered & scaled:

$$x_i \longrightarrow \frac{x_i - \bar{x}}{\sqrt{\frac{1}{n} \sum_{j=1}^{n} (x_j - \bar{x})^2}} \qquad [x \text{ is unitless}]$$

# Implications of assumptions

- $Y$ is mean zero when $(X_1, \ldots, X_p) = (0, \ldots, 0)$

- $X_i$ s are unitless

- $\displaystyle\sum_{i=1}^{n} X_i^2 = \sum_{i=1}^{n} \left[ \dfrac{x_i - \bar{x}}{\sqrt{\frac{1}{n} \sum_{j=1}^{n}(x_j - \bar{x})^2}} \right]^2$    $\rightarrow$ old xs

   new xs

$$= \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{\frac{1}{n} \sum_{j=1}^{n}(x_j - \bar{x})^2} = n$$

**DEF** The <u>ridge regression</u> estimator for $\beta$ minimizes

$$\sum_{i=1}^{n} \left( y_i - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2$$

$$= (y - X\beta)^\top (y - X\beta) + \lambda \|\beta\|_2^2$$

<span style="color:red">$\longrightarrow$ design matrix has no 1st column of 1s.</span>

**Ex** $p=1$    model    <u>$y = \beta_1 x + \varepsilon$</u>