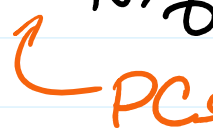# Principal Components Regression (PCR)

Instead of regressing $Y$ onto $X_{N \times P}$,

regress $Y$ onto $Z_{N \times q}$

↳ PCs

# Steps for PCR

⓪ wean center $X$

$$X_c = \begin{bmatrix} X_1 - \hat{\mu}_1 & X_2 - \hat{\mu}_2 & \cdots & X_P - \hat{\mu}_P \end{bmatrix}$$

$$\hat{\mu}_i = \text{mean}(X_i)$$

① do PCA on $X_c = UDV^T$

$$Z = X_c V_{1:q}$$

(2) regress $Y$ onto $Z$
(typically include intercept)

$$D = \begin{bmatrix} 1 & | \\ 1 & Z \\ 1 & | \end{bmatrix}$$

$$\hat{\beta}^{(PCR)} = (D^TD)^{-1}D^TY \in \mathbb{R}^{g+1}$$

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

What about new data?

Let $X^{test}$ be $M \times P$

For training: $\hat{Y} = D\hat{\beta}^{(PCR)}$

Need to calc. $D^{test}$ by applyg some procedure to $X^{test}$

(0) mean center $X^{test}$

$$X_c^{test} = \begin{bmatrix} X_1^{test} - \hat{\mu}_1 & \cdots & X_p^{test} - \hat{\mu}_p \end{bmatrix}$$

$$\Lambda_c = \begin{bmatrix} \Lambda_1 - \mu_1 & & & \mu_p & \mu_p \end{bmatrix}$$

from traing data

① apply PCA:

$$Z^{test} = X_c^{test} \, V_{1:q}$$

from training

② $D^{test} = \begin{bmatrix} 1 & Z^{test} \\ 1 & 1 \end{bmatrix}$

③ $\hat{y}^{test} = D^{test} \hat{\beta}^{(PCR)}$  from train.

---

## Compare PCR w/ ridge regression

$$\hat{Y}^{(ridge)} = X \hat{\beta}^{(ridge)}$$

up scaling cols of $Z$

$\vdots$

$$= \sum_{j=1}^{P} \left( \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \right) u_j u_j^T \gamma$$

$$\underbrace{\qquad}_{\Delta_j}$$

Can also show:

$$\hat{\gamma}^{(PCR)} = Z \hat{\beta}^{(PCR)} = \cdots = \sum_{j=1}^{P} \Delta_j u_j u_j^T \gamma$$

$$\Delta_j = \begin{cases} 1, & j \leq q \\ 0 & j > q \end{cases}$$

——  —  —  —  —  —  ——

OLS: $\Delta_j = 1$

ridge: $\Delta_j = \dfrac{\sigma_j^2}{\sigma_j^2 + \lambda}$

PCR: $\Delta_j = \mathbb{1}(j \leq q)$

——  —  —  —  ——

Consider $X$ to be full rank

Then as $\lambda \to 0$ we have

Then as $\lambda \to 0$ we have

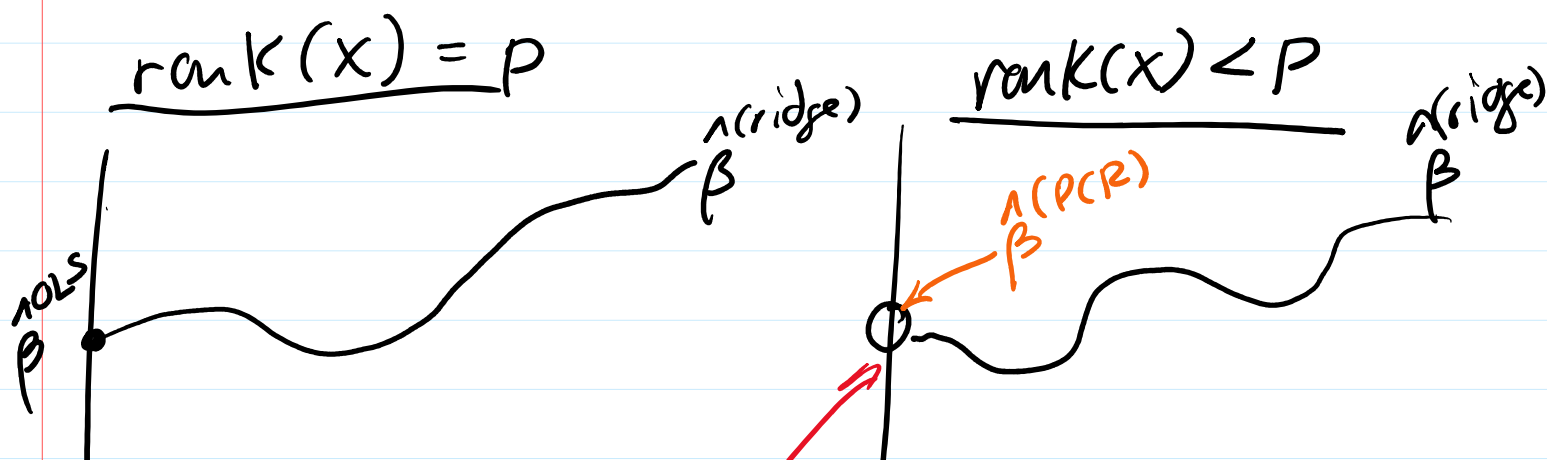$$\hat{\beta}^{(ridge)} \longrightarrow \hat{\beta}^{(OLS)}$$

Similarly, as $q \to P$ then
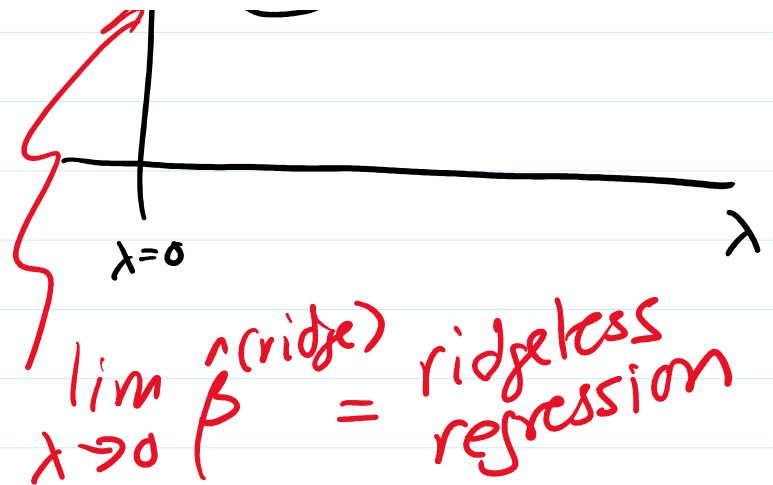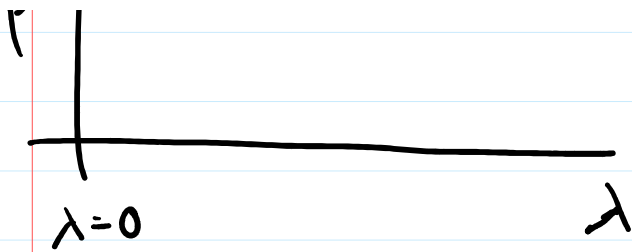
$$\hat{\beta}^{(PCR)} \longrightarrow \hat{\beta}^{(OLS)}$$

## If $rank(X) < P$

Then OLS doesn't exist.
(Equiv. ridge w/ $\lambda = 0$ doesn't exist)

However

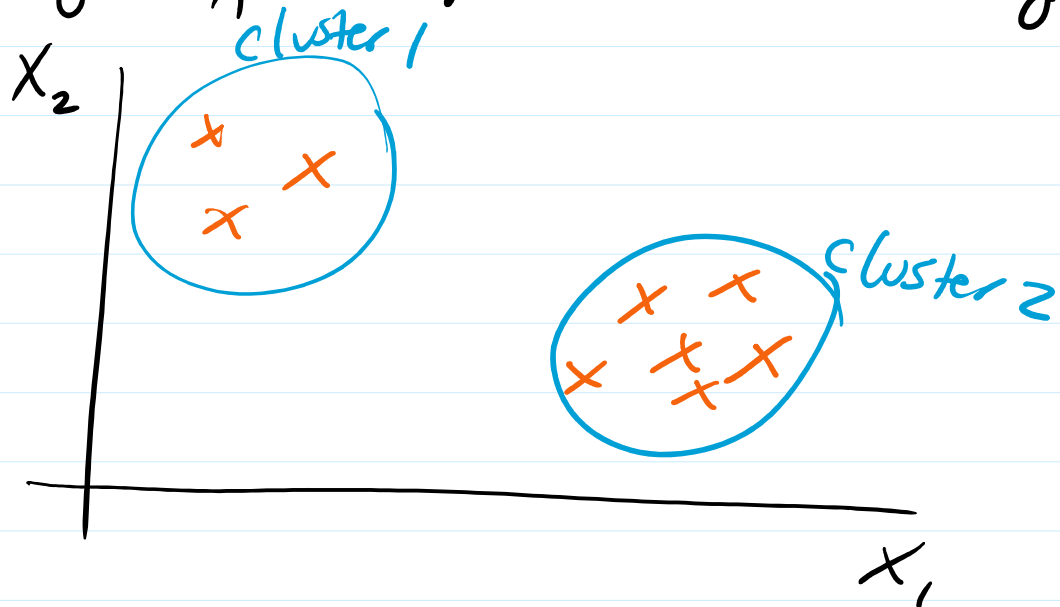$$\lim_{\lambda \to 0} \hat{\beta}^{(ridge)} = \hat{\beta}^{(PCR)}$$
$$w/ \ q = rank(X)$$

## $rank(X) = P$



$\hat{\beta}^{OLS}$     $\hat{\beta}^{(ridge)}$

## $rank(X) < P$



$\hat{\beta}^{(PCR)}$     $\hat{\beta}^{(ridge)}$

$$\lambda = 0 \qquad \lambda$$

$$\lambda = 0 \qquad \lambda$$

$$\lim_{\lambda \to 0} \hat{\beta}^{(ridge)} = \text{ridgeless regression}$$

## <u>Back to unsupervised</u>

Second type of unsupervised: clustering



cluster 1

$X_2$

cluster 2

$X_1$

To find clusters I need to define how pts are similar/dissim to each other

<u>Dissimilarity Measure</u>

## Dissimilarity Measure

If I have N obs then I can define a dissim mtx

$$D \qquad (N \times N)$$

where

$$D_{ii'} = \text{dissim meas betwn obs. } i \text{ and } i'.$$

For many clustering methods — don't need X but only D.

Properties of D:

① $D_{ii} = 0$

② $D_{ii'} \geqslant 0$

③ $D = D^T$

## K-means clustering

Assume that each data pt belongs to one of K clusters (or groups):

Assume that each data pt belongs to one of $K$ clusters (or groups):

$$G_1, \ldots, G_K$$

**want:** assign each point $i$ to some cluster $G_k$

**how:** want to make assignments so that we minimize some measure of not being well clustered (loss)

classic way to measure this is

$$W = \text{total within-cluster dissimilarity}$$

$$= \sum_{k=1}^{K} \sum_{i, i' \in G_k} D_{ii'}$$

small if clustering is good
large  "  "  bad

$$T = \text{total dissim} = \sum_{i,i'} D_{ii'}$$

$$B = \begin{array}{l}\text{total between}\\\text{cluster dissim}\end{array}$$

$$= \sum_{k \neq k'} \sum_{i \in G_k} \sum_{i \in G_{k'}} D_{ii'}.$$

Can show: $T = B + W$

So, to find $G_1, \ldots, G_K$

    ① minimize $W$

or ② maximize $B$

---

Ideally: try all possible cluster assignments

    practically: not possible — too computationally demanding

---

Assume all data is numeric: $x \in \mathbb{R}^p$

Assume all data is numeric: $x_n \in \mathbb{R}^p$

define
$$D_{ii'} = \|x_i - x_{i'}\|_2^2$$

If I do this then

$$W = \sum_{k=1}^{K} 2 N_k \sum_{i \in G_k} \|x_i - \bar{x}_k\|^2$$

dist of pt to mean of cluster

$N_k = |G_k|$

mean of pts in cluster $k$

Goal: try to minimize this $W$

Lloyd's Algorithm

Ⓞ Initialization Step

make some random initializations of cluster means:

$$\mu_1^{(0)}, \mu_2^{(0)}, \ldots, \mu_K^{(0)}$$

For $t = 1, 2, 3, \ldots$

   at the $t^{th}$ iteration:

     ① <u>Assignment Step:</u>
     assign $i$ (or $x_i$) to cluster $G_k$
     w/ the closest mean to $x_i$

$$G_k^{(t)} = \{ x_i : \| x_i - \mu_k^{(t)} \| \leq \| x_i - \hat{\mu}_{k'}^{(t)} \| \; \forall \, k' \}$$

     ② <u>Update Step:</u> re-calculate $\mu$s
     given current assignments

$$\hat{\mu}_k^{(t+1)} = \operatorname*{mean}_{i \in G_k^{(t)}} x_i$$