# Lecture 14: PCR

Instead of regressing $Y$ onto $X_{N \times p}$
we can regress $Y$ onto $Z_{N \times q}$, $q \ll p$

— — — — — — — — — — — — —

## Steps for PCR:

**(0)** mean center $X$

$$X_c = \begin{bmatrix} | & | & \\ X_1 - \text{mean}(X_1) & X_2 - \text{mean}(X_2) & \cdots \\ | & | & \end{bmatrix}$$

**(1)** do PCA: $X_c = UDV^T$

$$Z = X_c V_q \quad \longleftarrow \text{first } q \text{ cols of } V$$

**(2)** regress $Y$ onto $Z$

typically want to include intercept

so let

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

then $\hat{\beta}^{(PCR)} = (D^T D)^{-1} D^T Y \in \mathbb{R}^{q+1}$

What about predicting on new data?

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

for traing $\hat{Y} = D \hat{\beta}^{(PCR)}$

Need to form $D_{test}$ by applying same steps to $X^{test}$.

⓪ center the data

means of traing data

$$X_c^{test} = \begin{bmatrix} | & & | & \\ X_1^{test} - mean(X_1) & X_2^{test} - mean(X_2) & \cdots \\ | & & | & \end{bmatrix}$$

① apply PCA

from traing

$$Z^{test} = X_c^{test} V_q$$

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

then $\hat{y}^{test} = D_{test} \hat{\beta}^{(PCR)}$

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

## Compare w/ ridge regression

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

saw that

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

<span style="color:red">$\propto$ PCs</span>

<span style="color:red">Shrink contribution of small $\sigma_i$ more then large $\sigma_i$</span>

$$= \sum_{j=1}^{P} \Delta_j \, U_j U_j^T Y$$

$$\Delta_j = \frac{\sigma_j^2}{\sigma_j^2 + \lambda}$$

## Same analysis for PCR:

$$\hat{Y} = Z\hat{\beta}^{(PCR)} = Z(Z^TZ)^{-1}Z^T Y$$

<span style="color:magenta">$Z = XV_q = U_q D_q$</span>

$$= U_q D_q \left( (U_q D_q)^T U_q D_q \right)^{-1} (U_q D_q)^T Y$$

$$= U_q D_q (D_q U_q^T U_q D_q)^{-1} D_q^T U_q^T Y$$

<span style="color:magenta">$I$</span>

$$= U_q D_q (D_q^T D_q)^{-1} D_q U_q^T Y$$

$$\underbrace{\phantom{D_q (D_q^T D_q)^{-1} D_q}}_{I}$$

$$= U_q U_q^T Y$$

$$= \sum_{j=1}^{q} U_j U_j^T Y \qquad j\text{th col of } U$$

$$= \sum_{j=1}^{P} \Delta_j U_j U_j^T Y \qquad \Delta_j = \begin{cases} 1 & j \leq q \\ 0 & j > q \end{cases}$$



Consider $X$ to be full rank

then as $\lambda \longrightarrow 0$ we have

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

Similarly, as $q \to P$

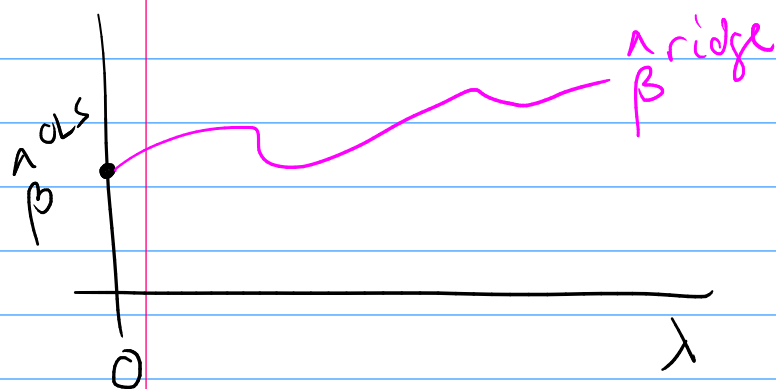$$\hat{\beta}^{PCR} \longrightarrow \hat{\beta}^{OLS}$$

If $rank(X) < P$ then $\hat{\beta}^{OLS}$ doesn't exist
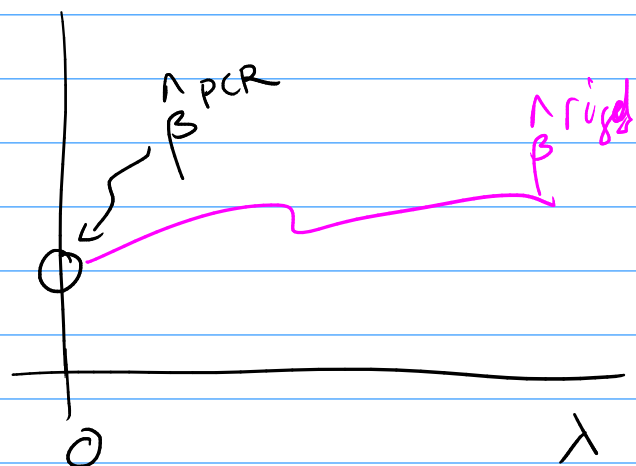
equivalently $\hat{\beta}^{ridge}$ w/ $\lambda = 0$ doesn't exist.

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

$rank(X) = P$
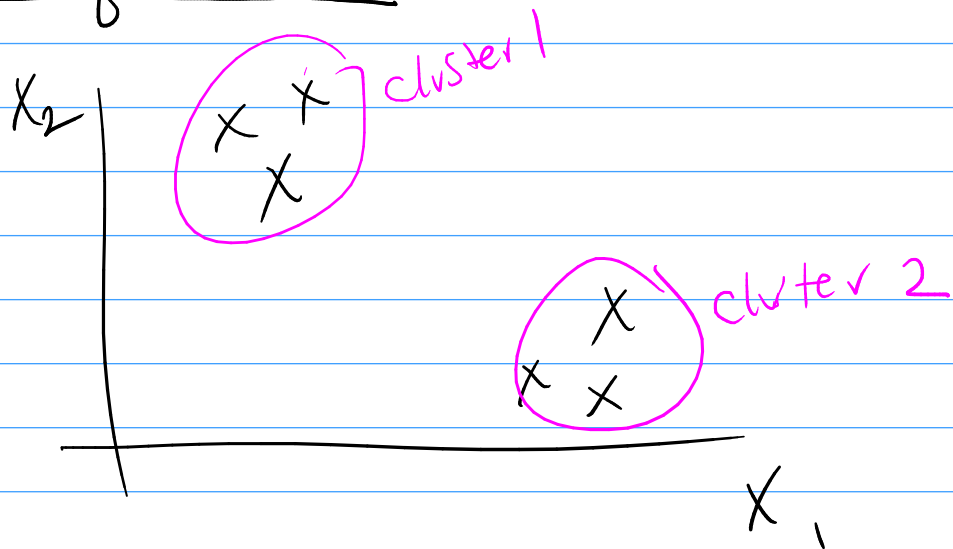


$rank(X) < P$

<u>Back to unsupervised methods</u>

Another class of unsupervised methods is
<u>clustering methods</u>



To find clusters I need some way of
defining how pts are similar/dissimilar

<u>Need</u>: Some dissimilarity measure

If I have $N$ observations then I
need to define some matrix

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

Where $D_{ii'}$ = dissimilarty measure between
obs. $i$ and $i'$.

For many methods don't explicitly need $X$
only need $D$

Need some properties to be true of D:

$$\text{(1)} \quad D_{ii} = 0$$

$$\text{(2)} \quad D_{ii'} \geq 0$$

$$\text{(3)} \quad D = D^T \quad \text{(symmetric)}$$

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

## K-means clustering

Assume each data point belongs to one of K clusters (or groups)

$$G_1, G_2, \ldots, G_K$$

<u>want</u>: assign each point $i$ to some cluster $G_k$

<u>how</u>: want to make assignments so that I minimize some measure of not being well clustered ("loss")

classic loss for clustering is

$$W = \begin{array}{c} \text{total w/in} \\ \text{cluster} \\ \text{dissimilarity} \end{array} = \sum_{k=1}^{K} \sum_{i, i' \in G_k} D_{ii'}$$

— W should be large if clustering is bad

— W " small if clustering is good

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

$$B = \begin{array}{l}\text{total between} \\ \text{cluster} \\ \text{dissim}\end{array} = \sum_{k,k'} \sum_{i \in G_k} \sum_{i' \in G_{k'}} D_{ii'}$$

One can show that

$$T = W + B$$

So to find $G_1, ..., G_K$ we should

either  ① minimize W

or  ② maximize B

Ideally: try all possible cluster assignments

practically: not comp. tractible for reasonably large N or K