

Lecture 19: PCAR

PCA regression, instead of regressing Y onto $X_{n \times p}$

we can regress Y onto $Z_{n \times q}$

(if $q \ll p$, this is basically some "smart" variable selection)

Steps for PCAR

① mean center X

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

cols of X

② Do PCA: $X_c = UDV^T$

$$Z = X_c V_g$$

first q cols of V

③ regress Y onto Z

typically want to include intercept,

$$D = \begin{bmatrix} 1 & \bar{z} \\ 1 & 1 \end{bmatrix}$$

$$\text{reg. coef: } \hat{\beta}^{(PCR)} = (D^T D)^{-1} D^T Y. \in \mathbb{R}^{q+1}$$

What about prediction on new data?

$$X^{\text{test}} \in \mathbb{R}^{n \times p}$$

$$\text{For training, } \hat{y} = D \hat{\beta}^{\text{PCR}}$$

Careful and do the same processing steps to X^{test}

(0) Center test data

$$X_c^{\text{test}} = \begin{bmatrix} x_1^{\text{test}} - \text{mean}(x_1) & x_2^{\text{test}} - \text{mean}(x_2) & \dots \\ | & | & \\ | & | & \end{bmatrix}$$

(1) apply PCA

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

$$(2) \hat{y}^{\text{test}} = Z^{\text{test}} \underset{\substack{m \times q \quad q \times 1}}{\beta}^{\text{PCR}}$$

Comparison w/ Ridge

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

$$\begin{aligned} \hat{y}_{\text{train}}^{\text{ridge}} &= X \hat{\beta}^{\text{ridge}} = X (X^T X + \lambda I)^{-1} X^T Y \\ &= \sum_{j=1}^p \left(\frac{\sigma_j^2}{\sigma_j^2 + \lambda} \right) u_j u_j^T Y \end{aligned}$$

$u_j =$ left singular vecs. of X
 $\sigma_j =$ sing. vals of X

$$= \sum_{j=1}^p \Delta_j U_j U_j^T Y \quad \text{where } \Delta_j = \frac{\sigma_j^2}{\sigma_j^2 + \lambda}$$

↪ proj. Y onto PCs and
take weighted avg. w/
weights Δ_j

PCR: $\hat{Y}_{\text{train}} = Z \hat{\beta}^{\text{PCR}}$

$$= Z (Z^T Z)^{-1} Z^T Y$$

$Z = X V_g = U_g D_g$

$$= U_g D_g (D_g U_g^T U_g D_g)^{-1} D_g U_g^T Y$$

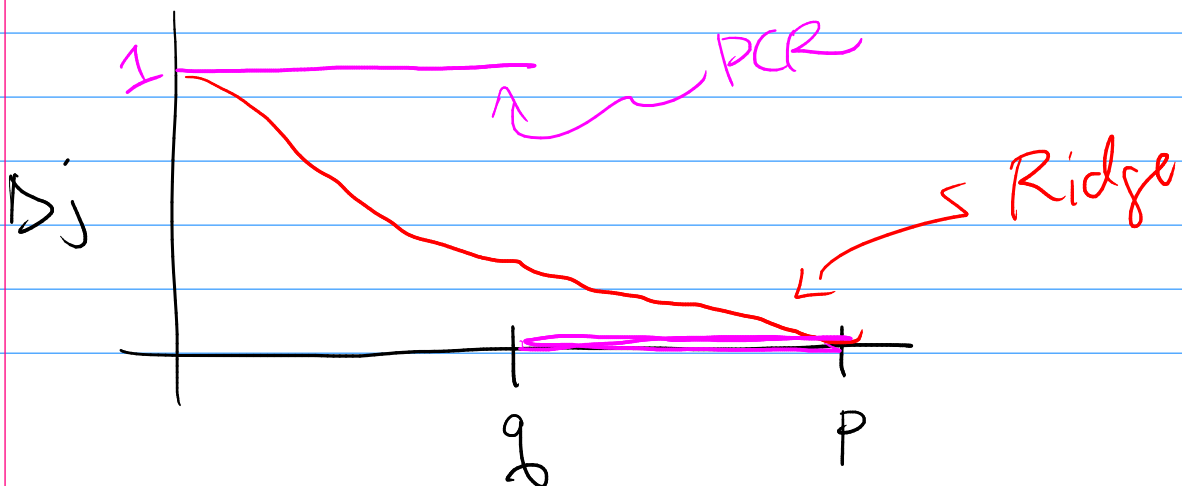
I

$$= U_g U_g^T Y$$

$$= \sum_{j=1}^g U_j U_j^T Y$$

$$= \sum_{j=1}^p \Delta_j^{\text{PCR}} U_j U_j^T Y$$

$$\Delta_j^{\text{PCR}} = \begin{cases} 1 & j \leq g \\ 0 & j > g \end{cases}$$



If X is full (ol rank ($\text{rank } X = P$))

then as $\lambda \rightarrow 0$

$$\hat{\beta}^{\text{ridge}} \rightarrow \hat{\beta}^{\text{OLS}}$$

similarly, as $q \rightarrow P$

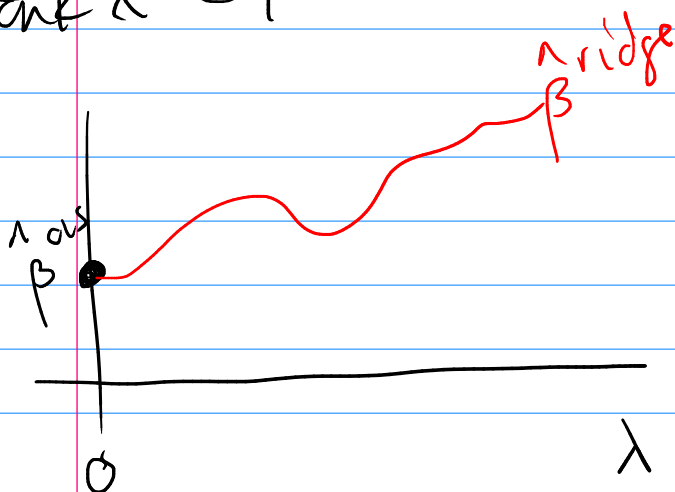
$$\hat{\beta}^{\text{PCR}} \rightarrow \hat{\beta}^{\text{OLS}}$$

However, if $\text{rank } X < P$. then $\hat{\beta}^{\text{OLS}}$ doesn't exist

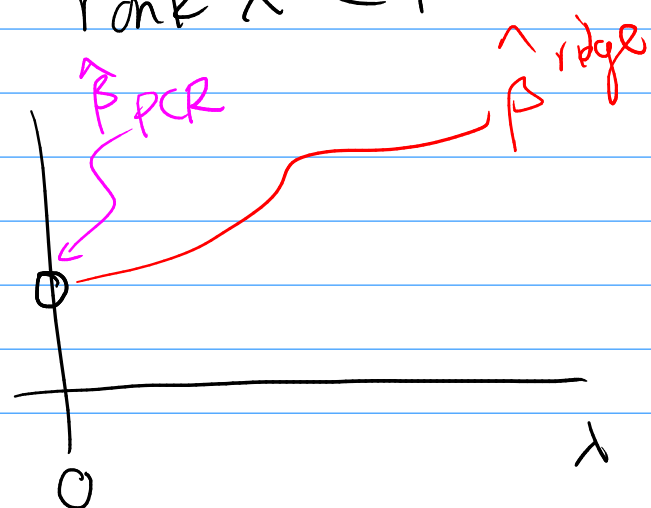
So $\hat{\beta}^{\text{ridge}}$ doesn't exist at $\lambda = 0$,

however, $\lim_{\lambda \rightarrow 0} \hat{\beta}^{\text{ridge}} = \hat{\beta}^{\text{PCR}}$ w/
 $q = \text{rank } X$.

$\text{rank } X = P$



$\text{rank } X < P$



Back to unsupervised learning

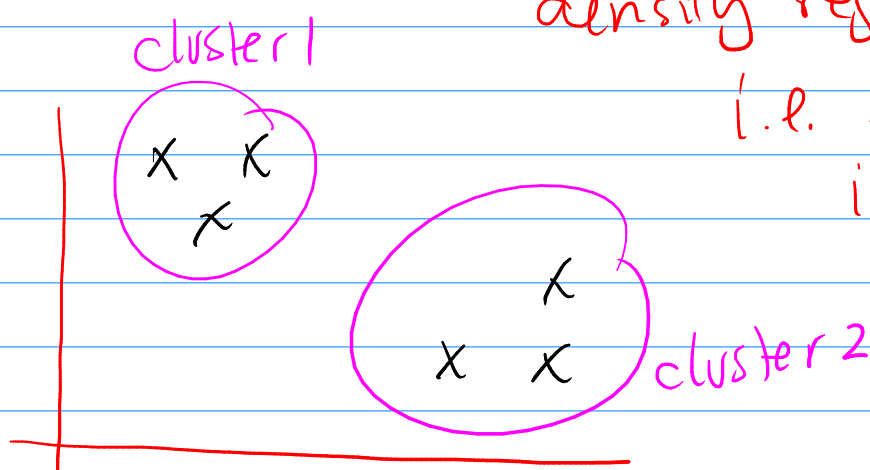
- Supervised problems: interested in $p(x, y)$

unsupervised problems: interested in $p(x)$

PCA: find a subspace where $p(x)$ is concentrated

clustering: want to find high density regions of $p(x)$

i.e. find "clusters" in our data



Goal: automatically find clusters

To do clustering we need some measure of either similarity or dissimilarity among observations

If I have N observations I need a matrix
to $(N \times N)$

Called the dissim. mtx where

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

[Can create dissim w/ dec. trans. of a sim. metric]

Most/all clustering algs only need D , not a data mtx X

Properties of D

(1) $D_{ii} = 0$ (diag. elements are zero)

(2) $D_{ii'} \geq 0$ (non-neg)

(3) $D = D^T$ (symmetric)

Ex. If I have a data mtx $X_{N \times P}$

I can calc. "attribute-based" dissims

$$D_{ii'} = d(x_i, x_{i'}) = \sum_{j=1}^P d_j(x_i, x_{i'})$$

↑ meas. for var j

Numeric

careat;
careful w/ scale

$$\text{Ex, } d_j(x_i, x_{i'}) = (x_{ij} - x_{i'j})^2$$

↑ euclidean

Ex, Categorical

$$d_j(x_i, x_{i'}) = \mathbb{1}(x_{ij} \neq x_{i'j})$$
