# Principal Components (Analysis) Regression

## Unsupervised Techniques:

① Dimensionality Reduction

$$X_1, \ldots, X_p \rightsquigarrow Z_1, \ldots, Z_q \quad \text{when } q \ll p$$

thereby we go from $p$ dims to $q \ll p$ dims.

② PCA summarize the $X$s into $Z$s that are Linear Combs of $X$s:

$$Z_j = \sum_{i=1}^{p} W_{ij} X_i = X W_j$$

i.e. $\quad Z = X W$

$$\begin{bmatrix} | & & | \\ Z_1 & \cdots & Z_q \\ | & & | \end{bmatrix} = \begin{bmatrix} | & & | \\ X_1 & \cdots & X_p \\ | & & | \end{bmatrix} \begin{bmatrix} | & & | \\ W_1 & \cdots & W_q \\ | & & | \end{bmatrix}$$

$$\qquad N \times q \qquad\qquad N \times P \qquad\qquad P \times q$$

③ Capture as much "info" in the orig. $X$s in the new $Z$s

$$\text{"info"} = \text{Variance}$$

So $\quad W_1 = \max \text{Var}(Z_1) \qquad \boxed{Z_1 = X W_1}$

$$\|W\| = 1$$

othervise

$$\text{Var}(1000 \ X \ w) = 1000^2 \ \text{Var}(Xw)$$

$$W_2 = \max_{\substack{\|w\|=1 \\ \text{Cor}(z_1, z_2) = 0}} \text{Var}(z_2)$$

don't capture
same info
$\left( \begin{array}{c} \text{parsimoniasly} \\ \text{reducig dimm} \end{array} \right)$

$$W_3 = \max_{\substack{\|w\|=1 \\ \text{Cor}(z_1, z_3) = \text{Cor}(z_2, z_3) = 0}} \text{Var}(z_3)$$

ete. up to $r = \text{rank}(x)$

Soln: $W = V_{1:q}$ where $X = UDV^T$

od $Z = X V_{1:q}$.

Typically, we ① mean center $X_j$s

Sometimes ② re-scal $X_j$s by their S.d.

One more interpretation of PCA:

$X V_{1:q} = Z$

$\wedge V_{1:q} = \tau$

2-D

proj. data onto this subspace

Zs are coords wrt. the basis of this space

Instead look at the projected data in the original basis of the space:

$$X_q = X P_w = \overbrace{X V_{1:q}}^{\text{coords}} V_{1:q}^T$$

proj. onto subspace

alt. if $X = U \Sigma V^T$

$$X_q = U_{1:q} D_q V_{1:q}^T = \text{truncated SVD}$$

$q \times q$ upper sub mtx

notice ⚡ have up to $r = \text{rank}(X)$ possible PCs ($Z_s$)

If $q = r$ then $X_q = X$

minimize info lost:

## Eckhart-Yang Theorem:

$$X_q = \underset{B:\, rank(B)=q}{\arg\min} \|X - B\|_F$$

$$\|A\|_F = \sqrt{\sum_i \sum_j a_{ij}^2}$$

where $q < r = rank(X)$

use $X_q$ over $X$ how much I lose

**PCR:** Instead of regressing $Y$ onto $X_{N\times P}$ lets regress $Y$ onto $Z_{N\times q}$ $q \ll P$.

(basically variable selection (?))

## Steps for PCR:

① Assume that $X$s have been <u>mean centered</u>
$$mean(X_j) = 0$$
(optionally: re-scale by S.dis of $X_j$)

② Let $X = UDV^T = $ PCA basically
then
$$Z = XV_{1:q}$$

since $V_{1:q} = W$

$$XV = UD$$

$U_{1:q}$ basically $Zs$

regress $Y$ onto $Z$.

(want an intercept.)

$$\left[ \text{often} \quad Z = \begin{bmatrix} 1 & XV_{1:q} \\ \vdots & \end{bmatrix} \quad \text{``want an intercept.''} \right]$$

i.e. $\hat{\gamma}^{(PCR)} = (\underline{Z^T Z})^{-1} Z^T Y \in \mathbb{R}^q \quad \left[ \text{or } \mathbb{R}^{q+1} \right]$

③ <u>Prediction:</u>

(traing data) $\hat{Y}_{train} = Z \hat{\gamma}$.

(new data) $\hat{Y}_{test} = \dots ?$

⓪ $A_{N \times P} = \begin{bmatrix} a_1 & \dots & a_P \\ | & & | \end{bmatrix} \longleftarrow$ new data

① $\tilde{A} = \begin{bmatrix} 1 & a_1 - \bar{X}_1 & \dots & a_P - \bar{X}_P \\ 1 & | & & | \end{bmatrix} \longleftarrow$ mean center [opt: re-scale]

② $\tilde{A} V_{1:q} \longleftarrow$ proj. into subspace

③ $\boxed{\hat{Y}_{new} = \underbrace{\tilde{A} V_{1:q} \hat{\gamma}}}$

$\boxed{\begin{array}{l} X = UDV^T \\ N\times P \quad N\times N \quad P\times P \end{array}}$

<u>Notice:</u> $\hat{\beta}^{(PCR)} \in \mathbb{R}^P$ since $V_{1:q}$ is $P \times q$

So $\hat{\gamma} \in \mathbb{R}^q$

so

$$\hat{Y} = \tilde{A}\hat{\beta}^{(PCR)} \qquad \hat{\gamma} \in \mathbb{R}^{q}$$

$N \times P$     $\mathbb{R}^{P}$ ($\sigma$ $\mathbb{R}^{P+1}$)

---

# Formalize comparison to Ridge Regression

## Ridge:

$$\hat{\beta}^{(ridge)} = (X^TX + \lambda I)^{-1} X^T Y \qquad D = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \sigma_r \\ & & 0 \end{bmatrix}$$

New data $A \in \mathbb{R}^{N \times P}$

$$X = UDV^T \quad \leftarrow \text{training data}$$

$$\hat{Y} = A\hat{\beta}^{(ridge)} = A \underbrace{VD^{-1}}_{\dot{A}} \overbrace{D V^T \hat{\beta}^{(ridge)}}^{I}$$

$\dot{A} = $ proj. $A$ onto $V$-basis and
$\phantom{\dot{A} = }$ re-scale by SVs
$\phantom{\dot{A} = } = A$ under different basis

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

$$= DV^T (VD^2V^T + \lambda I)^{-1} VDU^T Y$$

$$= D \underbrace{(D^2 + \lambda I)^{-1} D}_{\text{Diag}\left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right) = \Delta} U^T Y$$

$$\hat{Y} = \dot{A}\Delta U^T \hat{\beta}^{(ridge)} = \dot{A}\Delta U^T Y \qquad \xleftarrow{\text{weight } \Delta_i}$$

$$= \sum_{j=1}^{P} \dot{A}_j \left(\frac{\sigma_j^2}{\sigma_j^2 + \lambda}\right) U_j^T Y$$

          btwn

## Same genre for PCR

$$A\hat{\beta}^{(PCR)} = AUD^{-1}DV^{T}\hat{\beta}^{(PCR)}$$

$$= \dot{A}DV^{T}V_{1:q}\hat{\gamma} \qquad Z = XV_{1:q} \approx U_{1:q}D_q$$

$$= \dot{A}DV^{T}V_{1:q}(Z^{T}Z)^{-1}Z^{T}Y \qquad X^{T}X = VD^{2}V^{T}$$

$$\approx \dot{A}DV^{T}V_{1:q}(V_{1:q}^{T}X^{T}XV_{1:q})^{-1}V_{1:q}^{T}X^{T}Y$$

$$= \dot{A}DV^{T}\underbrace{V_{1:q}}_{I}(D_q^{2})^{-1}D_q U_{1:q}^{T}Y$$

$$= \dot{A}D(D_q)^{-2}D_q U_{1:q}^{T}Y$$

$$= \sum_{j=1}^{q}\dot{A}_j U_j^{T}Y \left[ = \sum_{j=1}^{P}\dot{A}_j \Delta_j U_j^{T}Y \right.$$

cord of A wrt PCs

$$\Delta_j = \begin{cases} 1 & \text{if } j \leq q \\ 0 & \text{else.} \end{cases}$$

## PCAR:

Scaling

$\uparrow$ PCAR

1

Scaling
$\triangleright \hat{j}$



$q$         $p$ ridge

If $X$ is full rank ($X^TX$ invertible)

as $\lambda \to 0$

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

Similarly as $q \to p = \text{rank}(X)$

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

Hoowever: if $X$ isn't full rank, however

as

$$\lambda \to 0 \quad \hat{\beta}^{(ridge)} \longrightarrow \hat{\beta}^{(PCR)}$$

w/ $q = \text{rank}(X)$.