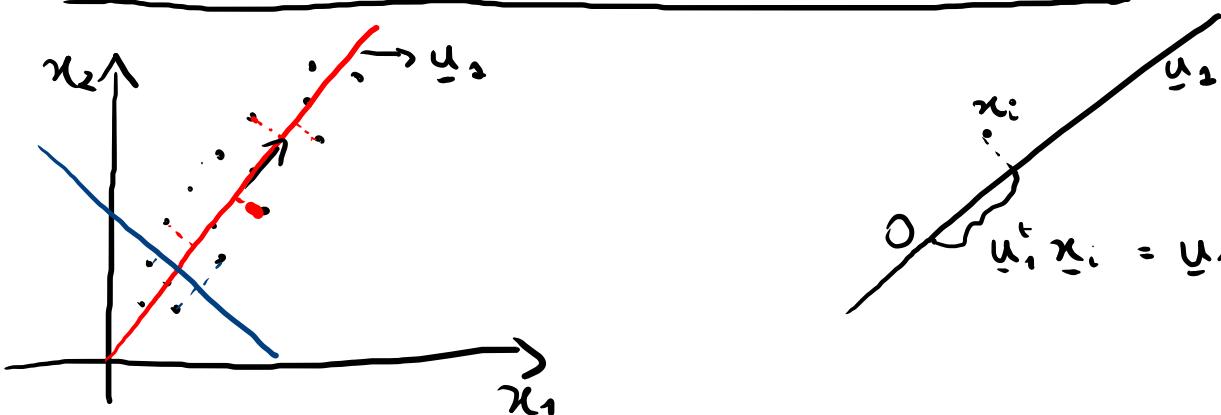


ALTERNATIVE DERIVATION OF PCA

10/05/21



$$u_i^T \underline{x}_i = \underline{u}_1 \cdot \underline{x}_i = u_{11}x_{i1} + u_{21}x_{i2} + \dots + u_{p1}x_{ip} = z_{i1}, \quad i=1, \dots, n$$

$$\underline{u}_1 = \begin{pmatrix} u_{11} \\ \vdots \\ u_{p1} \end{pmatrix} \rightarrow \text{loadings}$$

Assume that the data are centered ($\bar{x}_j = 0 \quad \forall j$) \Rightarrow z scores are also centered

Total variance of projections:

$$\underbrace{\frac{1}{n} \sum_{i=1}^n (z_{i1})^2}_{\text{OPTIMIZATION}} = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p u_{j1} x_{ij} \right)^2$$

$$\left(\frac{\sum_{i=1}^n z_{i1}}{n} = u_{11} \bar{x}_1 + \dots + u_{p1} \bar{x}_p = 0 \right)$$

PROBLEM

$$\underset{\underline{u}_1 \in \mathbb{R}^p, \|\underline{u}_1\|=1}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p u_{j1} x_{ij} \right)^2$$

$$\begin{aligned}
 & \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p u_{j1} x_{ij} \right)^2 = \frac{1}{n} \sum_{i=1}^n (u_1^t x_i)^2 \\
 & = \frac{1}{n} \sum_{i=1}^n (u_1^t x_i) \cdot (x_i^t u_1) \\
 & = \frac{1}{n} u_1^t \underbrace{\sum_{i=1}^n x_{i1} \cdot x_{i1}^t}_{\substack{px1 \ 1xp \\ pxp}} u_1 \\
 & = u_1^t \boxed{\frac{\sum_{i=1}^n x_{i1} \cdot x_{i1}^t}{n}} u_1
 \end{aligned}$$

SAMPLE COVARIANCE (S)

$$= u_1^t \frac{X^t X}{n} u_1 = \lambda_1$$

$$\text{Setting } u_1^t S u_1 = \lambda_1 \Rightarrow \underbrace{u_1 u_1^t}_I S u_1 = \lambda_1 u_1 \Rightarrow S u_1 = \lambda_1 u_1$$

So u_1 is the eigenvector of S corresponding to the largest eigenvalue.

What about the second component?

choose \underline{u}_2 so that $\underbrace{\underline{u}_2^T S \underline{u}_2}_{\text{variance of projections on } \underline{u}_2}$ is maximized and $\underline{u}_2 \perp \underline{u}_1$

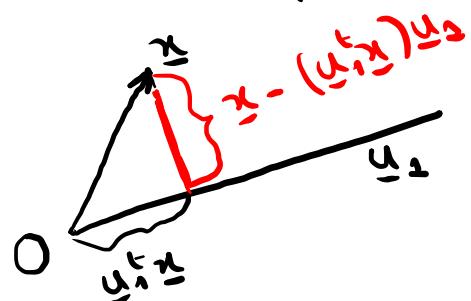
$\Rightarrow \underline{u}_2$ is the eigenvector corresponding to the second largest eigenvalue.

And so forth:

$\underline{u}_1, \underline{u}_2, \underline{u}_3, \dots, \underline{u}_p$

We've reconstructed the same basis as last time!

There is a final justification for this particular choice of vectors.



We can state the problem as finding \underline{u}_2 that minimizes the average squared distances of the points to the chosen direction.

OPTIMIZATION
PROBLEM

$$\underset{\substack{\underline{u}_1 \in \mathbb{R}^p \\ \|\underline{u}_1\|=1}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \|x_i - (\underline{u}_1^t x_i) \underline{u}_1\|^2$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|x_i - (\underline{u}_1^t x_i) \underline{u}_1\|^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - (\underline{u}_1^t x_i) \underline{u}_1)^t (x_i - (\underline{u}_1^t x_i) \underline{u}_1) \\ &= \frac{1}{n} \sum_{i=1}^n [x_i^t x_i - x_i^t (\underline{u}_1^t x_i) \underline{u}_1 - (\underline{u}_1^t x_i) \underline{u}_1^t x_i + \underline{u}_1^t (\underline{u}_1^t x_i) (\underline{u}_1^t x_i) \underline{u}_1] \\ &= \frac{1}{n} \sum_{i=1}^n [x_i^t x_i - \underbrace{(\underline{u}_1^t x_i)(\underline{u}_1^t \underline{u}_1)}_{\underline{u}_1^t \underline{u}_1} - (\underline{u}_1^t x_i)^2 + (\underline{u}_1^t x_i)^2 \underbrace{\underline{u}_1^t \underline{u}_1}_{\|\underline{u}_1\|^2 = 1}] \\ &= \frac{1}{n} \sum_{i=1}^n [x_i^t x_i - (\underline{u}_1^t x_i)^2] \\ &= \underbrace{\frac{1}{n} \sum_{i=1}^n x_i^t x_i}_{\text{does not depend on } \underline{u}_1} - \underbrace{\frac{1}{n} \sum_{i=1}^n \underline{u}_1^t x_i x_i^t \underline{u}_1}_{\underline{u}_1^t S \underline{u}_1} \Rightarrow \text{So } \underline{u}_1 \text{ should} \\ &\quad \text{maximize } \underline{u}_1^t S \underline{u}_1 \\ &\quad (\text{as before!}) \end{aligned}$$

With two components, we would look at the plane that is closest to the data points:

$$\frac{1}{n} \sum_{i=1}^n \| \underbrace{\underline{x}_i - \underline{z}_{i1} \underline{u}_1 - \underline{z}_{i2} \underline{u}_2}_{\underline{u}_1^\top \underline{x}_i} \|^2 \quad \leftarrow \text{We want to find } \underline{u}_1 \text{ and } \underline{u}_2 \text{ that minimize this quantity with } \underline{u}_1 \perp \underline{u}_2$$

$$= \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \underline{z}_{i1} \underline{u}_1 - \underline{z}_{i2} \underline{u}_2)^\top (\underline{x}_i - \underline{z}_{i1} \underline{u}_1 - \underline{z}_{i2} \underline{u}_2)$$

$$= \frac{1}{n} \sum_{i=1}^n (\underline{x}_i^\top \underline{x}_i - \underbrace{\underline{z}_{i1} \underline{x}_i^\top \underline{u}_1}_{=0} - \underbrace{\underline{z}_{i2} \underline{x}_i^\top \underline{u}_2}_{=0} - \underline{z}_{i1} \underline{u}_1^\top \underline{x}_i + \underbrace{\underline{z}_{i1}^2 \underline{u}_1^\top \underline{u}_1}_{\|\underline{u}_1\|^2=1} + \underbrace{\underline{z}_{i1} \underline{z}_{i2} \underline{u}_1^\top \underline{u}_2}_{=0 \text{ since } \underline{u}_1 \perp \underline{u}_2} + \underline{z}_{i2} \underline{u}_2^\top \underline{u}_2)$$

$$= \frac{1}{n} \sum_{i=1}^n [\cancel{\underline{x}_i^\top \underline{x}_i} - \cancel{\underline{z}_{i1}^2} - \cancel{\underline{z}_{i2}^2} - \cancel{\underline{z}_{i1}^2} + \cancel{\underline{z}_{i1}^2} - \cancel{\underline{z}_{i2}^2} + \cancel{\underline{z}_{i2}^2}]$$

$$= \underbrace{\frac{1}{n} \sum_{i=1}^n \underline{x}_i^\top \underline{x}_i}_{\text{does not depend on } \underline{u}_1 \text{ nor } \underline{u}_2} - \underbrace{\frac{1}{n} \sum_{i=1}^n \underline{z}_{i1}^2}_{\substack{\text{maximized by } \underline{u}_1 \\ \text{eigenvector corresponding to largest eigenvalue of } S}} - \underbrace{\frac{1}{n} \sum_{i=1}^n \underline{z}_{i2}^2}_{\substack{\text{second largest eigenvalue} \\ \text{of } S}}$$

CHOOSING THE NUMBER OF COMPONENTS

TOTAL VARIANCE

Trace ($\frac{X^t X}{n}$)

Sum of diagonal terms of sample covariance

$$= \text{Trace} \left(\frac{U D^t V^t V D U^t}{n} \right)$$

$$X = V D U^t \quad = \frac{\text{Trace}(U D^t D U^t)}{n} = \frac{\text{Trace}(D^t D U^t U)}{n}$$

$$= \frac{\text{Trace}(D^t D)}{n} =$$

$$= \frac{\lambda_1 + \lambda_2 + \dots + \lambda_p}{n}$$

$$D_{n \times p} = \begin{bmatrix} d_{11} & 0 & \dots & 0 \\ 0 & d_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & d_{pp} \end{bmatrix}, d_{ii} = \sqrt{\lambda_i}$$

$$D^t D_{p \times p} = \begin{bmatrix} d_{11}^2 & 0 & \dots & 0 \\ 0 & d_{22}^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & d_{pp}^2 \end{bmatrix}$$

VARIANCE EXPLAINED BY k^{th} COMPONENT

$\frac{1}{n}$

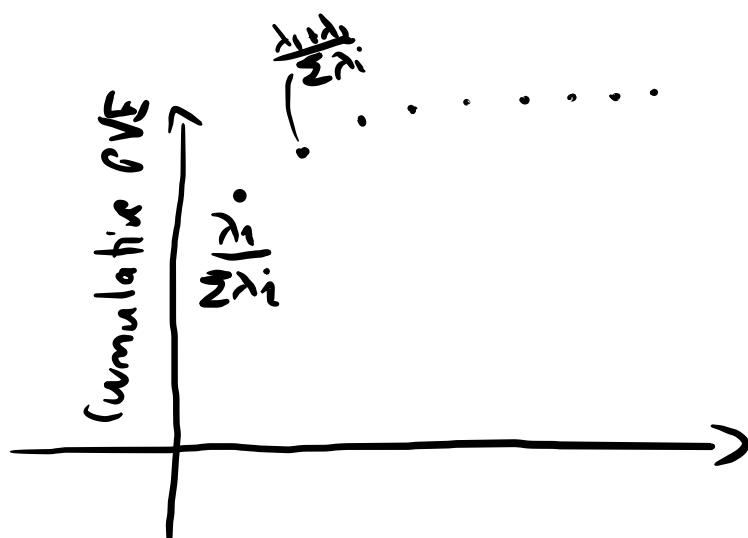
$$\sum_{i=1}^n z_{ik}^2 = \frac{\lambda_k}{n} \text{ where } \lambda_k \text{ is the } k^{th} \text{ eigenvalue of } X^t X$$

PROPORTION OF VARIANCE EXPLAINED (PVE)

$$PVE_k = \frac{\frac{\lambda_k}{n}}{\frac{\lambda_1 + \dots + \lambda_p}{n}} = \frac{\lambda_k}{\lambda_1 + \dots + \lambda_p}$$

$$\Rightarrow \sum_{k=1}^p PVE_k = 1$$

TYPICAL PLOTS :

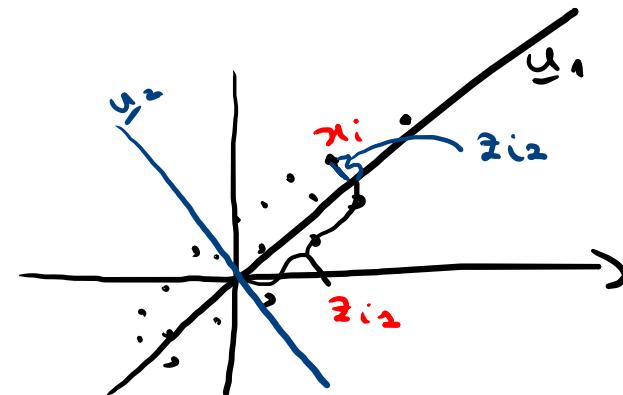


PCA IN PRACTICE

① Get familiar with the terminology!

$$z_k = u_{1k} x_1 + u_{2k} x_2 + \dots + u_{pk} x_p$$

$$z_{ik} = \sum_{j=1}^p u_{jk} x_{ij}, \quad i=1, \dots, n$$



u_1, \dots, u_p

u_{1k}, \dots, u_{pk}

K PRINCIPAL COMPONENT

SCORES OF Kth PRINCIPAL COMPONENT
(projections)

PRINCIPAL COMPONENT
LOADING VECTORS → directions

LOADINGS OF Kth COMPONENT

→ how much each variable x_1, \dots, x_p contributed to component K.

→ scores and loadings are both "visible" in a biplot

② SCALING OF DATA

→ Always advised

③ IMPLEMENTATION

Two functions:

- stats : prcomp (SVD)
- stats : princomp (ED of sample covariance)

Results unique up to sign flip !.

PCA IN SUPERVISED LEARNING (Section 6.3)

Consider these two cases:

- ① The number of variables p is very large, e.g. in genomics
→ Dimensionality reduction to reduce p
- ② Predictors are highly correlated
→ Summarising information via principal components
→ One component would represent a number of correlated variables

Two methods:

- PRINCIPAL COMPONENT REGRESSION (PCR)
- PARTIAL LEAST SQUARES (PLS)

PCK

$x_1, \dots, x_p \rightarrow \text{extract } M (< p) \text{ components}$

z_1, \dots, z_M

$\rightarrow \text{use } z_1, \dots, z_M \text{ as my new predictors}$

$$y = z\theta + \epsilon$$

$$= z_0\theta_0 + z_1\theta_1 + \dots + z_M\theta_M + \epsilon$$

$$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_M \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad z = \begin{pmatrix} z_{11} & \dots & z_{1M} \\ \vdots & \ddots & \vdots \\ z_{n1} & \dots & z_{nM} \end{pmatrix}$$

(scores)

The vector θ is related to the original regression coefficients β ($y = X\beta + \epsilon$).

In order to see this:

$$z_{ik} = u_{ik} x_{i1} + \dots + u_{ip} x_{ip}, \quad i=1, \dots, n, \quad k=1, \dots, M$$

In matrix form:

$$z = XU$$

$n \times M$ $n \times p$ $p \times M$

$$\text{with } U = \begin{pmatrix} u_{11} & \dots & u_{1M} \\ \vdots & \ddots & \vdots \\ u_{p1} & \dots & u_{pM} \end{pmatrix}$$

So

$$\begin{aligned}y &= \mathbf{Z}\boldsymbol{\theta} + \varepsilon \\&= \mathbf{X}\underbrace{\mathbf{U}\boldsymbol{\theta}}_{\mathbf{B}} + \varepsilon \\&= \mathbf{X}(\underbrace{\mathbf{U}\boldsymbol{\theta}}_{\mathbf{B}}) + \varepsilon\end{aligned}$$

So $\underline{\beta} = \sum_{p \times 1}^{p \times M} \mathbf{U} \underline{\theta}_{M \times 1}$

$$\beta_j = \sum_{k=1}^n u_{jk} \theta_k$$

- $\hat{\theta} \rightarrow \hat{\beta}$ interpretation
- Or check variables most associated to one component

choice of $M \rightarrow$
• Either like before
• Or, better, via CV or test error (supervised learning!)

PLS

This method uses \mathbf{y} (the response) also in the selection of the directions. In particular, we look for directions that:

- ① Explain/summarise variation in \mathbf{x}
- ② Are most related with \mathbf{y}

FIRST PLS DIRECTION

Find transformation $\underline{z}_1 = \mathbf{x} \underline{w}_1$ with $\|\underline{w}_1\|=1$ and such that $\underline{z}_1^t \underline{y}$ is maximised (with maximal correlation with \underline{y}).

This is found by projecting \mathbf{x} on \underline{y} , i.e.

$$\underline{w}_1 = \frac{\mathbf{x}^t \underline{y}}{\|\mathbf{x}^t \underline{y}\|}, \text{ i.e. } w_{j1} \propto \mathbf{x}_j^t \underline{y} \quad \text{proportional to regression coefficient of } \mathbf{x}_j \text{ on } \underline{y}$$

So $\underline{z}_{11} = \sum_{j=1}^p w_{j1} x_{ij} = w_{11} x_{i1} + \dots + w_{p1} x_{ip}$

maximal contribution by predictors that are most correlated with \mathbf{y}

SECOND PLS DIRECTION

$$x_1 = \gamma_1 z_1 + \varepsilon$$

⋮

$$x_p = \gamma_p z_1 + \varepsilon$$

Calculate residuals:

$$\hat{\varepsilon}_j = x_j - \hat{\gamma}_j z_1 \quad (\text{orthogonal to } z_1)$$

Create a new data matrix

$$E = \begin{pmatrix} \hat{\varepsilon}_{11} & \dots & \hat{\varepsilon}_{1p} \\ \vdots & & \vdots \\ \hat{\varepsilon}_{n1} & \dots & \hat{\varepsilon}_{np} \\ \hat{\varepsilon}_1 & & \hat{\varepsilon}_p \end{pmatrix}$$

Remaining information not accounted for by 1 PLS

Find a transformation $\underline{z}_2 = E \underline{w}_2$ with $\|\underline{w}_2\|=1$, such that $\underline{z}_2^T \underline{y}$ is maximised.

So $\underline{w}_2 \propto E^T \underline{y}$... And continue this way ...

The number of PLS components can be decided by CV.