PCA (Principal Component Analysis) → Dimension Reduction Technique.
(unsupervised)

(D)

| $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ |
|---|---|---|---|---|
| | | | | |

$$\left. \begin{array}{l} x_1 = 2, 2, 2, 2, 1, 2, 1 \\ \\ x_2 = 2, 12, 121, 31, 46, 89 \end{array} \right]-$$

$$var(x_2) > var(x_1)$$

$$D^{n\times 5} \rightsquigarrow D^{*\ n\times p}, \quad p < 5$$

variance of a data:

$$\text{variance}(D) = var(x_1) + var(x_2) + \dots + var(x_5)$$

$$= \sum_{i=1}^{5} var(x_i)$$

$var(X_1), var(X_4)$ are v. small

| $X_2$ | $X_3$ | $X_5$ |
|---|---|---|
| | | |

$$D^*_{n\times 3}$$
$$=$$

(?)

(1)

# Objective of PCA

Reduce the dimension of the data without reducing the variance.

Reduce the dimension of the data with a very small reduction in the variance.

$$D_{n \times p} \longrightarrow \overset{*}{D}_{n \times q} \ , \quad \boxed{q \ll p}$$

$$Var\left(\overset{*}{D}\right) \underset{\leq}{\approx} Var\left(D\right)$$

$p = 1000$ variables

$\downarrow$

$\boxed{30 - 40}$

Linear regression.

Case? Compulsion
.... ?

②

① $[X_1, X_2]$ ⟶ D* $[P_1, P_2]$



$X_2$

$(100,120)$ → $P_1$

$P_2$

sp.$(x_2)$

$O$ → $X_1$

spread$(x_1)$

$P_2$

$O$ ⟶ $P_1$

acc. false

$$H_0: b_i = 0$$
$$H_1: b_i \neq 0$$

p value increas.

$P_i$ | Y

$X_i$ | Y

$X$

Type 2 error.

③

$$p_{11} = a_{11} x_{11} + a_{21} x_{21}$$

$$p_{21} = a_{12} x_{11} + a_{22} x_{21}$$

$$\begin{pmatrix} p_{11} \\ p_{21} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{pmatrix} \begin{pmatrix} x_{11} \\ x_{21} \end{pmatrix}$$

$$\Rightarrow \boxed{\underset{\sim}{p_1} = A^T \underset{\sim}{x_1}}$$

$$p_{12} = a_{11} x_{12} + a_{21} x_{22}$$

$$p_{22} = a_{12} x_{12} + a_{22} x_{22}$$

$$\begin{pmatrix} p_{12} \\ p_{22} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{21} \\ a_{12} & a_{23} \end{pmatrix} \begin{pmatrix} x_{12} \\ x_{22} \end{pmatrix}$$

$$\Rightarrow \boxed{\underset{\sim}{p_2} = A^T \underset{\sim}{x_2}}$$

$$\dots \quad \boxed{\underset{\sim}{p_j} = A^T \underset{\sim}{x_j}}$$

$a_{ij}$ = loading corresponding to $i^{th}$ variable for $j^{th}$ PC

④

$$p_1 = A^T \underset{\sim}{x}_1$$

$$p_2 = A^T \underset{\sim}{x}_2$$

$$\vdots \qquad \vdots$$

$$p_n = A^T \underset{\sim}{x}_n$$

$\boxed{\#}$ $\bigcirc{P}$ $n \times 2$

|  | $P_1$ | $P_2$ | $[\cdots P_d]$ |
|---|---|---|---|
| $\underset{\sim}{x}_1$ | $p_{11}$ | $p_{21}$ |  |
| $\underset{\sim}{x}_2$ | $p_{12}$ | $p_{22}$ |  |
| $\vdots$ | | | |

$\#$ $\bigcirc{X}$ $n \times 2$

|  | $X_1$ | $X_2$ | $\cdots$ $X_d$ |
|---|---|---|---|
| $\underset{\sim}{x}_1$ | $x_{11}$ | $x_{21}$ |  |
| $\underset{\sim}{x}_2$ | $x_{12}$ | $x_{22}$ |  |
| $\vdots$ | | | |

$$\begin{pmatrix} p_1 & p_2 & \cdots & p_n \end{pmatrix} = \begin{pmatrix} A^T \underset{\sim}{x}_1 & A^T \underset{\sim}{x}_2 & \cdots & A^T \underset{\sim}{x}_n \end{pmatrix}$$

$$= A^T \begin{pmatrix} \underset{\sim}{x}_1 & \underset{\sim}{x}_2 & \cdots & \underset{\sim}{x}_n \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} p_{11} & p_{12} & \cdots & p_{1n} \\ p_{21} & p_{22} & & p_{2n} \end{pmatrix} = A^T \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & & x_{2n} \end{pmatrix} \Rightarrow \boxed{P^T = A^T X^T}$$

$$2 \times n \qquad\qquad\qquad 2 \times n$$

$$\boxed{P^T_{(2 \times n)} = A^T_{(2 \times 2)} X^T_{(2 \times n)}}$$

$\bigcirc{5}$

$$\boxed{X}$$

$$\begin{array}{c} \underset{\sim}{z_1} \\ \underset{\sim}{z_2} \\ \vdots \\ \underset{\sim}{z_n} \end{array} \quad \begin{array}{cccc} x_1 & x_2 & \cdots & x_p \\ \hline \left( x_{11} \right. & x_{21} & \cdots & \left. x_{p1} \right) \\ \left( x_{12} \right. & x_{22} & \cdots & \left. x_{p2} \right) \\ \vdots & \vdots & & \vdots \\ \left( x_{1n} \right. & x_{2n} & \cdots & \left. x_{pn} \right) \end{array}$$

$(n \times p)$

$\longrightarrow$

$$\boxed{P}$$

$$\begin{array}{c} \underset{\sim}{p_1} \\ \underset{\sim}{p_2} \\ \vdots \\ \underset{\sim}{p_n} \end{array} \quad \begin{array}{cccc} P_1 & P_2 & \cdots & P_p \\ \hline \underline{P_{11} \quad P_{21} \quad \cdots \quad P_{p1}} \\ P_{12} & P_{22} & \cdots & P_{p2} \\ \vdots & \vdots & & \vdots \\ P_{1n} & P_{2n} & \cdots & P_{pn} \end{array}$$

$(n \times p)$

$\boxed{a_{ij}}$ loading corresponding to ith variable & jth PC.

$$P_{11} = a_{11}\, x_{11} + a_{12}\, x_{21} + a_{13}\, x_{31} + \cdots + a_{p1}\, x_{p1} \Rightarrow P_{11} =$$

$$P_{12} = a_{12}\, x_{11} + a_{22}\, x_{21} + a_{32}\, x_{31} + \cdots + a_{p2}\, x_{p1}$$

$$\begin{array}{c|l} & x_1 \qquad\qquad x_2 \qquad\qquad x_3 \qquad\qquad\qquad x_q \\ \hline p_1 & P_{11} = a_{11}\, x_{11} + a_{21}\, x_{21} + a_{31}\, x_{31} + \cdots + a_{p1}\, x_{p1} \\[2mm] p_2 & P_{21} = a_{12}\, x_{11} + a_{22}\, x_{21} + a_{32}\, x_{31} + \cdots + a_{p2}\, x_{p1} \\ \vdots & \\ p_p & P_{p1} = a_{1p}\, x_{11} + a_{2p}\, x_{21} + a_{3p}\, x_{31} + \cdots + a_{pp}\, x_{p1} \end{array}$$

PC's are obtained by rotation of the original data $X = [X_1, X_2 \ldots, X_p]$ s.t.

1st PC is along the max. variance of the data X.

2nd PC is uncorrelated to 1st PC and is along the 2nd max. variance of X.

3rd PC is uncorrelated to 1st PC and 2nd PC and is along the 3rd max. variance of the data X

:
:
:

pth PC is uncorrelated to all the other PC's and is along the smallest variance of the data X.

# Result 1

Let $\Sigma$ be the __variance-covariance__ matrix of $(X)_{n \times p}$

__cars2.__ $(406 \times 3)$

$$\Sigma = \begin{array}{c} \\ X_1 \\ X_2 \\ \vdots \\ X_p \end{array} \begin{array}{cccc} X_1 & X_2 & \cdots & X_p \\ \left[ \begin{matrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ & \sigma_{22} & \cdots & \sigma_{2p} \\ & & \ddots & \vdots \\ & & & \sigma_{pp} \end{matrix} \right]_{p \times p} \end{array}$$

$(3 \times 3)$

$$\sigma_{11} = \frac{1}{n} \sum (x_{1i} - \bar{x}_1)^2 = Var(x_1)$$

$$\sigma_{12} = \frac{1}{n} \sum (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)$$

$$= Cov(x_1, x_2)$$

$\Sigma$ is a square matrix of dimension $p \times p$. It will have $p$ eigenvalue-eigenvector pairs, say, $(\lambda_1, \underset{\sim}{e}_1)$, $(\lambda_2, \underset{\sim}{e}_2)$, ...., $(\lambda_p, \underset{\sim}{e}_p)$.

If $\lambda_1 \geqslant \lambda_2 \geqslant \cdots \geqslant \lambda_p$, it can be shown that

$\underset{\sim}{e}_1$ is the 1st PC loading

$\underset{\sim}{e}_2$ is the 2nd PC loading, and so on...

$$\begin{pmatrix} (\lambda_1, \underset{\sim}{e}_1) \\ (\lambda_2, \underset{\sim}{e}_2) \\ (\lambda_3, \underset{\sim}{e}_3) \end{pmatrix}$$

$$A = \begin{pmatrix} \underset{\sim}{e}_1 & \underset{\sim}{e}_2 \cdots & \underset{\sim}{e}_p \end{pmatrix}$$

$$A = \begin{pmatrix} \underset{\sim}{e}_1 & \underset{\sim}{e}_2 & \underset{\sim}{e}_3 \end{pmatrix}$$