# Principle Components Analysis

(1) PCs of Population:

① Definition:

$X = (X_1 \cdots X_p)^T$, $E(X) = M$, $Var(X) = \Sigma_{p \times p}$.

Consider use $Z_i = a_i^T X = \sum_{j=1}^{p} a_{ji} X_j$ linear combination of $\{X_i\}$ to replace $\{Z_i\}$.

$Var(Z_i) = a_i^T \Sigma a_i$, $Cov(Z_i, Z_j) = a_i^T \Sigma a_j$

Note that: if $Var \uparrow$, then the data include more information.

Besides, we don't want to the information of $Z_1 \cdots Z_i$ overlaps. i.e. $Cov(Z_i, Z_j) = 0$, $i \neq j$.

Def: $Z_i = a_i^T X$ is $i^{th}$ component of $X$ if.

     i) $a_i^T a_i = 1$, $\forall 1 \leq i \leq p$.

     ii) $a_i^T \Sigma a_j = 0$, $\forall i \neq j$.

     iii) $Var(Z_i) = \max \{ Var(a^T x) | a^T a = 1, a^T \Sigma a_j = 0$

         for $\forall 1 \leq j \leq i-1 \}$

② Find Principle components:

   i) For $Z_1$:

Note that $\dfrac{a^T \Sigma a}{a^T a} \le \lambda_1$. $\lambda_1 \ge \lambda_2 \cdots \lambda_p$ is eigenvalues of $\Sigma$. Therefore, the $1^{st}$ component is $e_1$ eigenfunc. correspond $\lambda_1$. st. $\|e_1\|_2^2 = 1$.

ii) For $Z_k$:

$Var(\alpha^T X) \le \lambda_k$, if $\alpha^T \alpha = 1$. $\alpha^T \Sigma a_i = 0$. $1 \le i \le k-1$.

Then $Z_k = e_k$. correspond eigenvalue $\lambda_k$, $\|e_k\|_2^2 = 1$.

Thm. $Z = (Z_1 \cdots Z_p)^T$ is principle components of $X$.

     if i) $Z = A^T X$. $A$ is orthonormal. $A = (a_1 \cdots a_p)$

     ii) $Var(Z) = diag\{\lambda_1 \cdots \lambda_p\}$. $\lambda_1 \ge \lambda_2 \cdots \ge \lambda_p$.

     Rmk. To find $\{z_i\}_1^p$. Apply $\Sigma$ by a orthonormal

     diagonalization.

③ Properties:

i) $\sum_i^r \sigma_{ii} = tr(\Sigma) = \sum_i^r \lambda_i$.

     Rmk. If $\exists\, m \in \mathbb{Z}^+$. st. $\sum_i^r \sigma_{ii} \approx \sum_i^m \lambda_i$. Then

     replace $\{X_i\}_1^p$ by $\{Z_i\}_1^m$ to reduce data

ii) $\ell(Z_k. X_i) =: Cor(Z_k. X_i) = \sqrt{\lambda_k}\, a_{ik} / \sqrt{\sigma_{ii}}$.

     1f: $\ell(Z_k. X_i) = \dfrac{Cov(a_k^T X. e_i^T X)}{\sqrt{Var(Z_k) Var(X_i)}} = \dfrac{\lambda_k a_{ik}}{\sqrt{\lambda_k}\, \sigma_{ii}}$

     Rmk. We call $\ell(Z_k. X_i)$ is factor loading.

iii) $\sum\limits_{k=1}^{p} e^2(z_k, x_i) = \sum\limits_{k}^{p} \dfrac{\lambda_k a_{ik}}{\sigma_{ii}} = 1$

If: $\Sigma = A \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} A^T$  $\sigma_{ii} = a_i \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} a_i^T$

Rmk: Sum of squre of factor loading on

$\quad$ $x_i$ is 1. ( Full correlated )

iv) $\sum\limits_{i=1}^{p} \sigma_{ii} e^2(z_k, x_i) = \lambda_k$

If: $A^T \Sigma A = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix}$. $\lambda_k = a_k^T \Sigma a_k$.

Def: i) $\lambda_k / \sum\limits_{i}^{p} \lambda_i$ is rate of contribution of $z_k$.

$\quad$ $\sum\limits_{i}^{m} \lambda_k / \sum\limits_{i}^{p} \lambda_i$ is accumulation of rate

Consider two
$\Rightarrow$ factors when
select pcs

ii) Contribution of $[z_k]_1^m$ to $x_i : v_i^m$ def

$\quad$ by sum of squre of factor loading if $z_i$.

$\quad$ i.e. $v_i^m = \sum\limits_{k=1}^{r} \lambda_k a_{ik}^2 / \sigma_{ii}$

④ PcA of Standardization:

To eliminate the influence of units. We

can standardilization $X$. i.e. $X_i^* = \dfrac{X_i - E(X_i)}{\sqrt{\sigma_{ii}}}$

Then $Var(X^*) = R$. correlation matrix of $X$.

$\Rightarrow \sum\limits_{i}^{p} Var(z_i^*) = p = \sum\limits_{i}^{p} \lambda_i^*$. $z_i^* = a_i^T X^*$.

**Rmk:** If the Variance of $X_i$ differs a lot

Then the direction of PCA will differ

a lot from PCA on Standardilization data.

## (2) PCs of Samples:

When $\mu, \Sigma$ are unknown. We should infer from

the data matrix $X = (X_{ij})_{n \times p} = \begin{pmatrix} X_{(1)}^T \\ \vdots \\ X_{(n)}^T \end{pmatrix} = (X_1 \cdots X_p)$

① Common Case:

To find principle components:

i) replace population dist. by empirical dist.

ii) replace $\Sigma$ by $S = \frac{1}{n-1} \sum (X_{(t)} - \bar{X})^T (X_{(t)} - \bar{X})$

$\Rightarrow$ Find $(\hat{\lambda}_1, \hat{e}_1) \cdots (\hat{\lambda}_p, \hat{e}_p)$ eigenvalue - func pair

of $S$. $\hat{\lambda}_1 \geq \hat{\lambda}_2 \cdots \geq \hat{\lambda}_p$. Then $i^{th}$ pc is:

$\hat{\eta}_i = \hat{e}_i^T X = \sum_{i=1}^{i} \hat{e}_{ij} X_j$ for $1 \leq i \leq p$.

We obtain:

i) $\sum_i^p S_{ii} = tr(S) = \sum_i^p \hat{\lambda}_i$

ii) $\ell(\hat{\eta}_i, X_k) = \hat{e}_{ik} \sqrt{\hat{\lambda}_i} / \sqrt{S_{kk}}$.

② Standardilization:

If the data matrix is observed after standardilization

Then $R = \frac{1}{n-1} X^T X$ is correlation sample matrix.

Find $p$ pair eigenvalue-tuvre $(\lambda_i, a_i)$. $1 \leq i \leq p$, of $R$.

$\lambda_1 \geq \lambda_2 \cdots \geq \lambda_p$. $A = (a_1 \cdots a_p)$ is orthonormal.

<u>Def.</u> PC Score of $i^{th}$ PC $a_i^T X$ at $t^{th}$ sample
is $a_i^T X_{(t)}$. Denote by $z_{ti}$ (of $Z_i = A_i^T X$)

<u>properties:</u>

i) $\bar{z} = \frac{1}{n} \sum_i Z_{(i)} = 0$. $Z_i^T Z_j = (n-1)\lambda_i \delta_{ij}$

where $Z = XA = \begin{pmatrix} Z_{(1)}^T \\ \vdots \\ Z_{(n)}^T \end{pmatrix} = (z_1 \cdots z_p)$.

<u>Pf:</u> $\begin{pmatrix} Z_{(1)}^T \\ \vdots \\ Z_{(n)}^T \end{pmatrix} = \begin{pmatrix} X_{(1)}^T \\ \vdots \\ X_{(n)}^T \end{pmatrix} (a_1 \cdots a_p)$

$= (X_{(i)}^T a_j)_{n \times p} = (a_j^T X_{(i)})_{n \times p}$

ii) Principle components can minimize SSE.

Consider linear model:
(Exact $m$ PCs. Others as residue $\varepsilon_i$)
$$\begin{cases} X_1 = b_{11} z_1 + \cdots + b_{1m} z_m + \vec{\varepsilon_1} \\ \vdots \\ X_p = b_{p1} z_1 + \cdots + b_{pm} z_m + \vec{\varepsilon_p} \end{cases}$$

$B = (b_{ij})_{p \times m}$ $Z^* = (z_1 \cdots z_m)$. $X = Z^* B^T + E$.

LSE is $\hat{B}^T = ((Z^*)^T Z^*)^{-1} (Z^*)^T X$, $A^* = (a_1 \cdots a_m)$

$= (A^{*T} X^T X A^*)^{-1} A^{*T} X^T X$

$= \text{diag} [\lambda_1 \cdots \lambda_m]^{-1} (RA^*)^T$

$= (a_1 \cdots a_m)^T = A^{*T}$

$\Rightarrow Q(A^*) = \min Q(B)$

<u>Rmk:</u> See geometry of PC line in (3) ☺.

② Large Sample:

For $(\hat{\lambda}_k, \hat{e}_k, \hat{\eta}_k)_{p=1,2,\cdots,p}$ in ①.

Assume: i) $X_{(k)}$ is random sample from normal dist.

    ii) eigenvalues of $\Sigma$ satisfies: $\lambda_1 > \lambda_2 \cdots > \lambda_p > 0$

Then: $\hat{\lambda} = (\hat{\lambda}_1 \cdots \hat{\lambda}_k)$. $\sqrt{n}(\hat{\lambda} - \vec{\lambda}) \sim AN_p(0, 2\Lambda^2)$.

where $\Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix}$. $\vec{\lambda} = (\lambda_1 \cdots \lambda_p)$.

Denote: $E_i = \lambda_i \sum_{k \neq i}^{p} \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} e_k e_k^T$. $e_k$ is eigen func.

correspond to $\lambda_k$ w.r.t. $\Sigma$.

Then $\sqrt{n}(\hat{e}_i - e_i) \sim AN_p(0, E_i)$

Rmk: i) $\hat{\lambda}_i$ is indept with $\hat{e}_i$.

    ii) Obtain confidence interval with $\alpha$ level:

$$\lambda_i \in \left[ \hat{\lambda}_i / (1 \pm z(\tfrac{\alpha}{2}) \sqrt{2/n}) \right].$$

(3) Application:

① Connection with SVD:

· Apply SVD on $X(I - P_n) = X - J\bar{X}^T \underset{n \times p}{=} U_{n \times p} L_{p \times p} V^T_{p \times p}$

where $U^T U = I_p$. $V^T V = I_p$ $L = diag\{l_1 \cdots l_p\}$

$\Rightarrow S = \frac{1}{n-1}(X - J\bar{X}^T)^T(X - J\bar{X}^T) = V(\frac{L^2}{n-1})V^T$

It means:

    i) $(V)$ is eigenfunction of $S$.

ii) $l_i^2/(n-1)$ : eigenvalues of $S$.

iii) $Y = (X - J_n \bar{X}^T) V = UL$ : PC score.
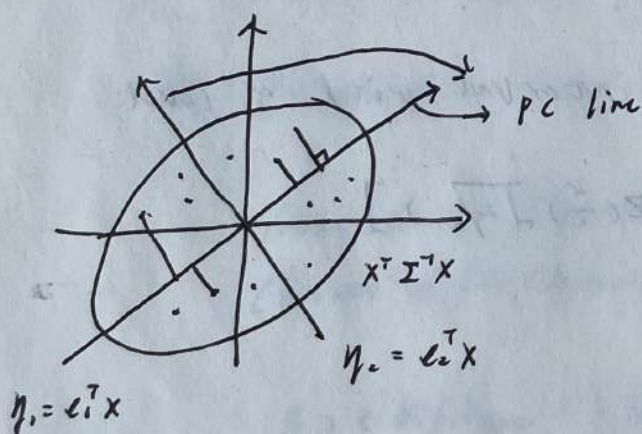
$$\tilde{B} = U \begin{pmatrix} l_1 & & \\ & \ddots & l_k \\ & & O_{p-k} \end{pmatrix} V^T = \arg\min \{ \text{tr}[(X - J_n \bar{X} - B)^T$$

$$(X - J_n \bar{X} - B)] \mid r(B) \leq k \}. \quad i.e.$$

$$\min_{B} \text{tr}[(X - J_n \bar{X} - B)^T (X - J_n \bar{X} - B)] = \sum_{k=1}^{p} \hat{\lambda}_i$$
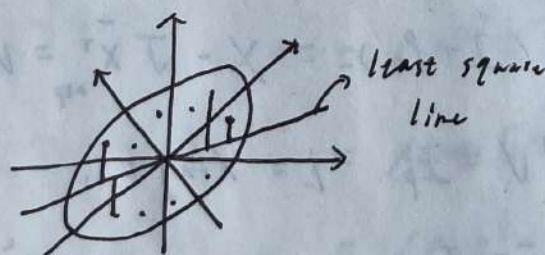
Rmk: $\text{tr}(X^T X) = \sum_{ij} |X_{ij}|^2$. $\text{tr}((A-B)^T(A-B))$

$$= \sum_{ij} (a_{ij} - b_{ij})^2. \quad \text{means} \quad SS \text{ of difference.}$$

② Geometry of PC lines:



$X^T \Sigma^{-1} X$

$\eta_2 = l_2^T X$

$\eta_1 = l_1^T X$

PC lines minimize the sum of squared orthogonal distances from each data to PC plane

Rmk: Compare to regression line: (least square line)



least square line

It minimizes the sum of vertical distance from data points to this line

③ Classification:

If $X$ is standardilized. $R = \frac{1}{n-1} X^T X$, has eigenvalues

$\lambda_1 \geq \lambda_2 \cdots \geq \lambda_p \geq 0$. PCs: $Z = (z_1 \cdots z_p)$

If we take the first $m$ components $Z^* = (z_1 \cdots z_m)$

$X^* = (X_1^* \cdots X_p^*) = (Z^* \ O_{n \times (p-m)}) A. = Z^* A^{*T}. \quad A = \begin{pmatrix} A^* \\ \tilde{A} \end{pmatrix}$

prop. $\text{tr}((X-X^*)^T (X-X^*)) = (n-1) \sum\limits_{m+1}^{p} \lambda_k$.

$\quad$ Pf: $\quad X - X^* = X(I - A^* A^{*T}). \quad A^* = \begin{pmatrix} a_1^T \\ \vdots \\ a_m^T \end{pmatrix}$

$\quad\quad \Rightarrow (X-X^*)^T (X-X^*) = (n-1)(I - A^* A^{*T}) R (I - A^* A^{*T})$

$\quad\quad\quad\quad = (n-1) (\sum\limits_{m+1}^{p} \lambda_i a_i a_i^T)$

i) If $r_{ij} \simeq 1$. Then $X_i, X_j$ can be classified as one class.

$\| X_i - X_j \|_2^2 = 2(n-1)(1 - r_{ij}) \simeq \| X_i^* - X_j^* \|_2^2$. For $1 - r_{ij}$:

$1 - r_{ij} = \frac{1}{2(n-1)} \| \sum\limits_{1}^{m} (a_{ik} - a_{jk}) z_k \|_2^2 = \frac{1}{2(n-1)} \sum\limits_{k=1}^{m} \lambda_k (a_{jk} - a_{ik})^2$

$\quad = \frac{1}{2(n-1)} \sum\limits_{1}^{m} (\ell_{ik} - \ell_{jk})^2. \quad \ell_{jk} = \ell(z_k, X_i) = \sqrt{\lambda_k} a_{ik}$

prop. If $\sum\limits_{1}^{m} (\ell_{ik} - \ell_{jk})^2 \simeq 0$. Then classify $X_i, X_j$

$\quad\quad$ as the same class

ii) Analogously. $\| X_{(i)} - X_{(j)} \|_2^2 \simeq 0 \Rightarrow X_{(i)}, X_{(j)}$ can be

$\quad$ recognized as samples from same class

$\| X_{(i)} - X_{(j)} \|_2^2 \simeq \| X_{(i)}^* - X_{(j)}^* \|_2^2 \quad (X_{(i)}^* = A^* z_i^{*T})$

$\quad\quad\quad\quad = \sum\limits_{k=1}^{m} (z_{ik} - z_{jk})^2$

④ Drawback:

PCA doesn't use information of $k^{th}$ moments $(k \geq 2)$

It may miss nonlinear structure or distort by outlier.