

Multivariate Analysis

1st Sem 2019/2020

1st Test - 5th Nov 2019 -

Group I

1. $(X, Y)^t \sim N_2(\underline{\mu}, \Sigma) : \text{Var}(X) = \text{Var}(Y)$

$$\underline{Z} = \begin{bmatrix} X+Y \\ X-Y \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} = \underline{A} \begin{bmatrix} X \\ Y \end{bmatrix}, \text{ where } \underline{A} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$\text{Var}(\underline{Z}) = \underline{A} \text{Var}(\underline{X}) \underline{A}^t = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12} & \Sigma_{22} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} \Sigma_{11} + \Sigma_{12} & \Sigma_{12} + \Sigma_{11} \\ \Sigma_{11} - \Sigma_{12} & \Sigma_{12} - \Sigma_{11} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} 2(\Sigma_{11} + \Sigma_{12}) & 0 \\ 0 & 2(\Sigma_{11} - \Sigma_{12}) \end{bmatrix} = \Sigma_Z$$

Since (i) $\text{Cov}(X+Y, X-Y) = 0$ and

(ii) $\Sigma = \underline{A} \Sigma \underline{A}^t$ and $\underline{X} \sim N_2(\underline{\mu}, \Sigma)$

then $\underline{Z} \sim N_2(\underline{\mu}_Z, \Sigma_Z)$ and

$Z_1 \perp\!\!\!\perp Z_2$ i.e. $X+Y \perp\!\!\!\perp X-Y$ QED

2. $\text{Var}(\underline{X}_1 + \underline{X}_2) \neq \Sigma_{11} + \Sigma_{22} + 2 \Sigma_{12}$

$$\mathbb{E}(\underline{X}_1 + \underline{X}_2) = \mathbb{E}(\underline{X}_1) + \mathbb{E}(\underline{X}_2) = \underline{\mu}_1 + \underline{\mu}_2$$

$$\begin{aligned}
\text{Var}(X_1 + X_2) &= E((X_1 + X_2)(X_1 + X_2)^t) - (\mu_1 + \mu_2)(\mu_1 + \mu_2)^t \\
&= E\left\{ X_1 X_1^t + X_1 X_2^t + X_2 X_1^t + X_2 X_2^t \right\} - \mu_1 \mu_1^t - \mu_2 \mu_2^t \\
&\quad - \mu_2 \mu_1^t - \mu_2 \mu_2^t = \\
&= [E(X_1 X_1^t) - \mu_1 \mu_1^t] + [E(X_1 X_2^t) - \mu_1 \mu_2^t] \\
&\quad [E(X_2 X_1^t) - \mu_2 \mu_1^t] + [E(X_2 X_2^t) - \mu_2 \mu_2^t] \\
&= \text{Var}(X_1) + \text{Cov}(X_1, X_2) + \text{Cov}(X_2, X_1) \\
&\quad + \text{Var}(X_2) = \Sigma_{11} + \Sigma_{12} + \Sigma_{21} + \Sigma_{22}
\end{aligned}$$

In general, $\text{Cov}(X_{1i}, X_{2j}) \neq \text{Cov}(X_{1j}, X_{2i})$
and being so $\Sigma_{12} \neq \Sigma_{21}$

3. $\underline{X} \sim N_3(\underline{\mu}, \underline{\Sigma})$, $\underline{\mu} = (-3, 1, 4)^t$, $\underline{\Sigma} = \underline{\Gamma} \underline{\Lambda} \underline{\Gamma}^t$

$$\underline{\Sigma} = \begin{bmatrix} 1 & 1 & 0 \\ 2 & 0 & 3 \end{bmatrix}, \underline{\Gamma} = \begin{bmatrix} 0 & 0.53 & 0.85 \\ 0 & 0.85 & -0.53 \\ 1 & 0.00 & 0.00 \end{bmatrix}, \underline{\Lambda} = \begin{bmatrix} 3 & 2.62 & 0.38 \end{bmatrix}$$

(a) $P(\text{RD}(\underline{X}, \underline{\mu}) > 3.0574) = 1 - P(\text{RD}^2(\underline{X}, \underline{\mu}) \leq 9.348)$

$$= 1 - F_{\chi_{(3)}^2}(9.348) = 1 - 0.9707 = 0.0293$$

$$0.95 = F_{\chi_{(3)}^2}(7.815) \leq F_{\chi_{(3)}^2}(9) \leq F_{\chi_{(3)}^2}(9.348) = 0.975$$

$$0.025 \leq 1 - F_{\chi_{(3)}^2}(9) \leq 0.05$$

(b) If $\underline{Y} = \underline{\Sigma}^{-1/2}(\underline{X} - \underline{\mu})$ then

(i) $\underline{Y} \sim N_3(0, I)$ since $\underline{Y} = A\underline{X}$ and
 \underline{X} has $N_3(\underline{\mu}, \underline{\Sigma})$

$$(ii) E(\underline{Y}) = \underline{\Sigma}^{-1/2} E(\underline{X} - \underline{\mu}) = 0$$

$$\begin{aligned} \text{Var}(\underline{Y}) &= \text{Var}(\underline{\Sigma}^{-1/2}(\underline{X} - \underline{\mu})) = \\ &= \underline{\Sigma}^{-1/2} \sum (\underline{\Sigma}^{-1/2})^t = \\ &= (\underline{\Sigma}^{-1/2} \underline{\Gamma}^t) \underline{\Gamma} \underline{\Gamma}^t (\underline{\Sigma}^{-1/2} \underline{\Gamma}^t) \\ &= \underline{\Gamma} \underline{\Sigma}^{-1/2} \underline{\Sigma}^{-1/2} \underline{\Gamma}^t = \underline{\Gamma} \underline{\Gamma}^t = I \end{aligned}$$

where

$$\begin{aligned} \underline{\Sigma}^{-1/2} &= \underline{\Gamma} \underline{\Sigma}^{-1/2} \underline{\Gamma}^t = \\ &= \begin{bmatrix} 0 & 0.53 & 0.85 \\ 0 & 0.85 & -0.53 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{3} & & \\ & 1/\sqrt{2.62} & \\ & & 1/\sqrt{0.38} \end{bmatrix} \end{aligned}$$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 0.53 & 0.85 & 0 \\ 0.85 & -0.53 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0.53 & 0.85 \\ 0 & 0.85 & -0.53 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0.577 \\ 0.325 & 0.526 & 0 \\ 1.376 & -0.851 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1.342 & -0.447 & 0 \\ -0.447 & 0.894 & 0 \\ 0 & 0 & 0.577 \end{bmatrix}$$

Group II

1. Give the classical estimates of the eigen-value we conclude that the 1st three PCs explain 86.46% of the observed total variability

$$\frac{\sum_{i=1}^3 \hat{d}_i}{\sum_{i=1}^7 \hat{d}_i} = \frac{\sum_{i=1}^3 \hat{d}_i}{7} = 0.8646$$

Note that the estimates were obtained from the standardized raw-values so:

$$\text{tr}(\hat{R}) = 7 = \sum_{i=1}^7 \hat{d}_i$$

If we use the criteria to select the number of PCs which $\hat{d}_i \geq \bar{d}$ ($\Leftrightarrow \hat{d}_i \geq 1$) then only one PC should be selected since

$$\hat{d}_1 \geq 1 \quad (\Leftrightarrow \sqrt{\hat{d}_1} \geq 1)$$

and the R output returns the value $\sqrt{\hat{d}_i}, i=1, \dots, 7$.

2. $\hat{y}_1 = +0.45 \left(\frac{x_1 - 2.69}{0.51} \right) + 0.31 \left(\frac{x_2 - 1.79}{0.05} \right) + 0.40 \left(\frac{x_3 - 13.17}{1.50} \right)$

$$+ 0.43 \left(\frac{x_4 - 2.02}{0.94} \right) + 0.45 \left(\frac{x_5 - 6.21}{0.40} \right)$$

$$+ 0.24 \left(\frac{x_6 - 41.28}{3.47} \right) + 0.30 \left(\frac{x_7 - 28.52}{6.15} \right)$$

3. The first PC is a weighted mean of all the events, and since large scores on the events means a better performance, competitors with the best marks have very high values on the 1st PC and competitors with bad performances have the lowest possible values on PC1.

The second PC, PC₂, is essentially a contrast between high jump performance and the 800m run (all other weights are considered approximately zero). Being so, a competitor with very good result on high jump ($\hat{\beta}_{22} = -0.65$) and very poorly on 800m ($\hat{\beta}_{27} = 0.66$) will have a very negative score. Contrarily, competitors with good results on 800m run and bad on high jump will have very high scores on the PC₂.

$$4. \quad \underline{x}_1 = (3.73, 1.86, 15.8, 4.05, 7.27, 45.66, 34.92)^t$$

while the maximum values (best scores on all events) are :

Max: 3.73, 1.86, 16.23, 4.05, 7.27, 47.50, 39.23

So, the USA competitor won the 100m hurdles, high jump, 200m, and the long jump.

So

$$\begin{aligned} \hat{y}_1 &= +0.45 \left(\frac{3.73 - 2.69}{0.51} \right) + 0.31 \left(\frac{1.86 - 1.79}{0.05} \right) + 0.40 \left(\frac{15.8 - 13.17}{1.50} \right) \\ &\quad + 0.43 \left(\frac{4.05 - 2.02}{0.94} \right) + 0.45 \left(\frac{7.27 - 6.21}{0.40} \right) \\ &\quad + 0.24 \left(\frac{45.66 - 41.28}{3.47} \right) + 0.30 \left(\frac{34.92 - 28.52}{6.15} \right) = \dots \end{aligned}$$

so, her score is expected to be very high, which is confirmed by the plot PC₁ vs Total score.

5. The robust estimates of PC lead to :

(i) higher % of variability explained by the first two PC :

	class	RCD
$\frac{\sum \hat{d}_i}{P}$	0.7461	0.7685

(ii) the reinterpretation of the 1st PC stays the same - overall measure of performance in both cases.

The 2nd PC reinterpretation changes:

(i) PC₂(RCD) is mainly the javelin results
 PC₂(class) is a contrast between high jump and run 800m

So far the point of view of building a overall index of the quality/performance of the competitors (and their order) there are no major differences.

6. Yes, plot PC1 vs Total Score clear indicates a high positive correlation between these two indexes.