

**Multivariate Analysis**

Mater in Eng. and Data Science & Master in Mathematics and Applications

1st Test

Duration: 1.5 hours

1st Semester – 2020/2021

19/11/2020 – 20:00

Group I

1. Let X_1, X_2 , and X_3 be three independent, univariate Normal distributed random variables, with unitary mean and variance ($\mathcal{N}(1, 1)$). Let $\mathbf{Y} = (Y_1, Y_2, Y_3)^t$, where $Y_1 = X_1 - 3X_2 + 2$, $Y_2 = 2X_1 - X_2 - 1$, and $Y_3 = X_3 - 1$. Determine the distribution of \mathbf{Y} .

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} = \begin{bmatrix} X_1 - 3X_2 + 2 \\ 2X_1 - X_2 - 1 \\ X_3 - 1 \end{bmatrix} = \begin{bmatrix} 1 & -3 & 0 \\ 2 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} + \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix}$$

Thus, $E(\underline{Y}) = A\underline{\mu} + \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 & -3 & 0 \\ 2 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix}$

$$= \begin{bmatrix} -2+2 \\ 1-1 \\ 1-1 \end{bmatrix} = \underline{\Omega}$$

$$\text{Var}(\underline{Y}) = \underline{\text{Var}(\underline{X})} A^t = \begin{bmatrix} 1 & -3 & 0 \\ 2 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 0 \\ -3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} =$$

$$= \begin{bmatrix} 10 & 5 & 0 \\ 5 & 5 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \underline{\Sigma}_Y$$

Since $\underline{Y} = A\underline{X} + \underline{\Omega}$ and $\underline{X} \sim \mathcal{N}\left(\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \underline{\Sigma}_X\right)$ then $\underline{Y} \sim \mathcal{N}_3(\underline{\Omega}, \underline{\Sigma}_Y)$

2. Suppose $\mathbf{X} = (X_1, X_2, X_3)^t$ has a multivariate normal distribution. Show that if:

- (i) X_1 and $X_2 + X_3$ are independent,
 - (ii) X_2 and $X_1 + X_3$ are independent, and
 - (iii) X_3 and $X_1 + X_2$ are independent,
- then X_1, X_2 , and X_3 are independent random variables.

$$\text{diag}(\text{Cov}(\underline{\mathbf{X}}, \underbrace{\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \underline{\mathbf{X}}})) = \underline{\mathbf{0}}$$

$$(\Leftarrow) \text{diag}(\underline{\Sigma}_A) = \underline{\mathbf{0}} (\Leftarrow)$$

$$(\Leftarrow) \text{diag}\left(\begin{bmatrix} \tilde{\sigma}_{11} & \tilde{\sigma}_{12} & \tilde{\sigma}_{13} \\ \tilde{\sigma}_{12} & \tilde{\sigma}_{22} & \tilde{\sigma}_{23} \\ \tilde{\sigma}_{13} & \tilde{\sigma}_{23} & \tilde{\sigma}_{33} \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}\right) = \underline{\mathbf{0}}$$

(\Leftarrow)

$$\begin{cases} \tilde{\sigma}_{12} + \tilde{\sigma}_{13} = 0 \\ \tilde{\sigma}_{12} - \tilde{\sigma}_{23} = 0 \\ \tilde{\sigma}_{13} - \tilde{\sigma}_{23} = 0 \end{cases} \quad \begin{cases} -\tilde{\sigma}_{23} - \tilde{\sigma}_{23} = 0 \\ \tilde{\sigma}_{12} = -\tilde{\sigma}_{23} \\ \tilde{\sigma}_{13} = -\tilde{\sigma}_{23} \end{cases} \quad \begin{cases} \tilde{\sigma}_{23} = 0 \\ \tilde{\sigma}_{12} = 0 \\ \tilde{\sigma}_{13} = 0 \end{cases}$$

$$\text{thus, } \underline{\Sigma} = \text{diag}(\tilde{\sigma}_{11}, \tilde{\sigma}_{22}, \tilde{\sigma}_{33})$$

and given that $\underline{\mathbf{X}} \sim N_p(\underline{\mu}, \underline{\Sigma})$ then

$\underline{\mathbf{X}} = (X_1, X_2, X_3)^t$ are indep. r.v

3. Let $\mathbf{X} = (X_1, X_2, X_3, X_4)^t$ be a random vector with multivariate normal distribution with parameters:

$$\boldsymbol{\mu} = \begin{pmatrix} 1 \\ 2 \\ 1 \\ 3 \end{pmatrix} \quad \text{and} \quad \boldsymbol{\Sigma} = \begin{pmatrix} 2 & 1 & 2 & 3 \\ 1 & 2 & 1 & 3 \\ 2 & 1 & 3 & 3 \\ 1 & 3 & 3 & 7 \end{pmatrix}$$

Determine the distribution of $(X_3, X_4)^t | (X_1, X_2)^t = (x_1, x_2)^t$.

Suggestion: If $\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ then $E(\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2) = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$ and $Var(\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2) = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$.

Let $\underline{Y} = (X_1, X_2)^t, \underline{Z} = (X_3, X_4)^t$

then $\underline{Y} | \underline{Z} = \underline{z} \sim \mathcal{N}_2(\boldsymbol{\mu}_{Y|Z}, \boldsymbol{\Sigma}_{Y|Z})$

$$\begin{aligned} E(\underline{Y} | \underline{Z} = \underline{z}) &= E(\underline{Y}) + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\underline{z} - E(\underline{Z})) \\ &= \begin{bmatrix} 1 \\ 2 \\ 1 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} (\underline{z} - \begin{bmatrix} 9 \\ 2 \end{bmatrix}) \end{aligned}$$

$$Var(\underline{Y} | \underline{Z} = \underline{z}) = \begin{bmatrix} 3 & 3 \\ 3 & 7 \end{bmatrix} - \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 2 & 3 \\ 1 & 3 \end{bmatrix} =$$

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} = \frac{1}{4-1} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

$$Var(\underline{Y} | \underline{Z} = \underline{z}) = \begin{bmatrix} 3 & 3 \\ 3 & 7 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 1 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} 3 & 3 \\ 3 & 7 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 4-1 & 6-3 \\ -2+2 & -3+6 \end{bmatrix}$$

$$= \begin{bmatrix} 3 & 3 \\ 3 & 7 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 3 & 3 \\ 0 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} 3 & 3 \\ 3 & 7 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 6 & 6+3 \\ 9 & 9+9 \end{bmatrix} = \begin{bmatrix} 3 & 3 \\ 3 & 7 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 6 & 9 \\ 9 & 18 \end{bmatrix}$$

$$= \begin{bmatrix} 3-2 & 3-3 \\ 3-3 & 7-6 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{aligned}
E(\underline{y} | \underline{z} = \underline{x}) &= E(\underline{y}) + \sum_{i2} \sum_{22}^{-1} (\underline{x} - E(\underline{z})) \\
&= \begin{bmatrix} 1 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} (\underline{x} - \begin{bmatrix} 1 \\ 2 \end{bmatrix}) \\
&= \begin{bmatrix} 1 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix} \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -7 & 2 \end{bmatrix} (\underline{x} - \begin{bmatrix} 1 \\ 2 \end{bmatrix}) \\
&= \begin{bmatrix} 1 \\ 3 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 4-1 & -2+2 \\ 6-3 & -3+6 \end{bmatrix} (\underline{x} - \begin{bmatrix} 1 \\ 2 \end{bmatrix}) \\
&= \begin{bmatrix} 1 \\ 3 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 3 & 0 \\ 3 & 3 \end{bmatrix} (\underline{x} - \begin{bmatrix} 1 \\ 2 \end{bmatrix}) \\
&= \begin{bmatrix} 1 \\ 3 \end{bmatrix} + \begin{bmatrix} x_1-1 \\ (x_1-1) + (x_2-2) \end{bmatrix} \\
&= \begin{bmatrix} x_1 \\ x_1+x_2 \end{bmatrix}
\end{aligned}$$

thus, $\begin{bmatrix} x_3 \\ x_4 \end{bmatrix} \mid \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim N_2 \left(\begin{bmatrix} x_1 \\ x_1+x_2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$

4. Let \mathbf{X}_1 and \mathbf{X}_2 be two independent random vectors, where $\mathbf{X}_i \sim \mathcal{N}_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$, $i = 1, 2$. Consider two independent random samples, with sizes n_1 and n_2 from each population. Prove that

$$\frac{n_1 n_2}{n_1 + n_2} (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2))^t \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)) \sim \chi^2_{(p)}.$$

Given that $x_i \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ then

\bar{x}_i is also $N_p(E(\bar{x}_i), \text{Var}(\bar{x}_i))$

$$\text{where } E(\bar{x}_i) = \sum_{h=1}^{n_i} \frac{E(x_{hi})}{n_i} = \frac{n_i E(\mu_i)}{n_i} = \mu_i$$

$$\text{Var}(\bar{x}_i) = \text{Var}\left(\sum_{h=1}^{n_i} \frac{x_{hi}}{n_i}\right) = \frac{1}{n_i^2} \sum_{h=1}^{n_i} \text{Var}(x_{hi}) = \frac{n_i}{n_i^2} \sum = \frac{\sum}{n_i}$$

since $x_{hi} \perp x_{hj}$ if $h \neq j$

i.e. $\bar{x}_i \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$.

Given that the two random samples are indep, so

$$\bar{x}_1 \perp \bar{x}_2 \quad \text{and} \quad \bar{x}_1 - \bar{x}_2 \sim N_p(E(\bar{x}_1 - \bar{x}_2), \text{Var}(\bar{x}_1 - \bar{x}_2))$$

thus,

$$E(\bar{x}_1 - \bar{x}_2) = E(\bar{x}_1) - E(\bar{x}_2) = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$$

$$\text{Var}(\bar{x}_1 - \bar{x}_2) = \text{Var}(\bar{x}_1) + \text{Var}(\bar{x}_2) = \frac{1}{n_1} \sum + \frac{1}{n_2} \sum =$$

$$\begin{aligned} & \bar{x}_1 \perp \bar{x}_2 \\ & = \frac{(n_1 + n_2)}{n_1 n_2} \sum \end{aligned}$$

Being so $\bar{z} = \bar{x}_1 - \bar{x}_2 \sim N_p(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2, \frac{n_1 + n_2}{n_1 n_2} \sum)$

Thus $(\bar{z} - E(\bar{z}))^T \sum_{zz}^{-1} (\bar{z} - E(\bar{z})) \sim \chi^2_{(p)}$

(=)

$$(\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2))^T \frac{n_1 n_2}{n_1 + n_2} \sum^{-1} (\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)) \sim \chi^2_{(p)}$$

QED

Group II**10.0 points**

The U.S. crime data set consists of the reported number of crimes in the 50 U.S. states in 1985. The crimes were classified according to 7 categories:

- X_1 : murder,
- X_2 : rape,
- X_3 : robbery,
- X_4 : assault
- X_5 : burglary,
- X_6 : larceny,
- X_7 : auto theft.

1. Obtain the standard deviation of the second sample principal component, α , missing in the previous R output.

We know that $\text{Var}(\hat{\gamma}_i) = \hat{d}_i$

and $\sum_{i=1}^7 \hat{d}_i = 7$ since the loadings are estimated based on the sample correlation matrix. Being so, $\sum_{i=1}^7 \hat{d}_i = 7$

$$\begin{aligned}\text{and } \hat{d}_2 &= 7 - \sum_{\substack{i=1 \\ i \neq 2}}^7 \hat{d}_i \\ &= 7 - 1.9776^2 - 0.8224^2 - \dots - 0.3405^2 \\ &= 7 - 5.536711 = 1.463289\end{aligned}$$

$$\text{and } \alpha = \sqrt{\hat{d}_2} = \sqrt{1.463289} = 1.209665$$

2. Decide how many principal components to retain.

Suggestion: If you have not answer the previous question use $\alpha = 1.2$.

$$\begin{aligned}\hat{\lambda}_i &: 3.9107 & 1.4633 & 0.6767 & 0.3715 \\ \hat{\lambda}_1/7 & 0.5587 & 0.2090 & 0.0966 & 0.0531 \\ \sum_{i=1}^k \hat{\lambda}_i/7 & 0.5587 & 0.7677 & 0.8643 & 0.9174\end{aligned}$$

thus, if the criteria is:

$$(i) \hat{\lambda}_i > \bar{\lambda} \quad (e) \hat{\lambda}_i > 1 \Rightarrow k=2$$

$$(ii) \text{ if \% Variability explained is } > 0.80 \Rightarrow k=3$$

For simplicity reasons we choose $k=2$.

3. Interpret the first two principal components.

The 1st PC is a weighted mean of all the murder types, so it is a global measure of crime in US.

High (low) values on the 1st PC represents states with high (low) crime levels.

	Comp.1	Comp.2	Comp.3
murd	0.272	0.653	
rape	0.431	0.117	
robb	0.376	-0.051	
assa	0.397	0.455	
burg	0.425	-0.309	
larc	0.362	-0.370	
auto	0.360	-0.344	

the 2nd PC is a contrast between murder and assault against burglary, larceny, and auto theft. A state with high (low) value on the 2nd PC is a state with high (low) levels of crimes injured to people (murder and assault) and low (high) levels of crimes injured to property (burglary)

4. To identify atypical cities related with the crime an outlier detection based on Robust PCA, using the first three robust principal components and a false alarm rate of 0.001, was estimated leading to the following distance-distance, where GA - Georgia, MD - Maryland, MI - Michigan, and NY - New York.

- (a) Obtained the critical value for the score distance.
 (b) Interpret the results.

(a) $k = 3$, $\alpha = 0.001$, Assuming $\underline{x} \sim N_p(\mu, \Sigma)$
 Then $\hat{\underline{y}} = (\hat{y}_1, \hat{y}_2, \hat{y}_3)^t \sim N_3(\mu_y, \Sigma_y)$
 where $\hat{\mu}_y = \bar{E}(\hat{\underline{y}}) = \bar{E}(\underline{R}^t \Sigma) = \underline{R}^t \bar{E}(\Sigma) = \underline{0}$

$$\Sigma_y = \text{Var}(\hat{\underline{y}}) = \text{diag}(\hat{d}_1, \hat{d}_2, \hat{d}_3)$$

$$\begin{aligned} \text{Thus, } SDC(\Sigma) &= \hat{\underline{y}}^t \Sigma_y^{-1} \hat{\underline{y}}^t = n D^2(\underline{y}, \underline{0}) \\ &= \sum_{i=1}^3 \frac{\hat{y}_i^2}{\hat{d}_i} \sim \chi^2_{(3)} \end{aligned}$$

Being so,

$$c_{00} = P(D^2(\underline{x}) > c_{00}^2) = \alpha = 0.001$$

$$c_{00}^2 = F_{\chi^2_{(3)}}^{-1}(1-0.001) = 16.26624$$

$$\text{and } c_{00} = \sqrt{16.26624} = 4.03162$$

(b) According with the classification rule there are 4 states classified as atypical:

New York - with high SD and OD score and

Georgia, Maryland, and Michigan all target by having a high OD.

