TMA4267 Linear Statistical Models V2017 [L5] Part 1: Multivariate RVs, and the multivariate normal distribution Properties of the multivariate normal distribution [H:2.6,4.4,5.1]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: January 24, 2017

Last lecture: derived the MGF and pdf of the multivariate normal distribution

1.
$$Z \sim N_1(0, 1)$$

 \blacktriangleright MGF: $M_Z(t) = E(e^{tz}) = e^{\frac{1}{2}t^2}$
2. Z_1, Z_2, \dots, Z_p iid $N_1(0, 1) \rightarrow Z_{p \times 1} \sim N_p(0, I)$
 \blacktriangleright MGF: $M_Z(t) = E(e^{t^T z}) = e^{\frac{1}{2}t^T t}$
3. $X = AZ + \mu$, $AA^T = \Sigma$ gives $X_{p \times 1} \sim N_p(\mu, \Sigma)$
 \blacktriangleright MGF: $M_X(t) = E(e^{t^T x}) = e^{t^T \mu + \frac{1}{2}t^T t}$
 \triangleright pdf (Σ invertible):

$$f(\boldsymbol{x}) = \frac{1}{(2\pi)^{\frac{\rho}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} exp\{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\}$$

Why is the mulitivariate normal distribution so important in statistics?

- Many natural phenomena may be modelled using this distribution (just as in the univariate case).
- Multivariate version of the central limit theorem- the sample mean will be approximately multivariate normal for large samples.
- Good interpretability of the covariance.
- Mathematically tractable.
- Building block in many models and methods.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $A\Sigma B^{T} = \mathbf{0} \Leftrightarrow AX$ and BX are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $X_2 \mid (X_1 = x_1) \sim N_{p2}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 \mu_1), \Sigma_{22} \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}).$

Diabetes data

We will study a data set on diabetes in Part 2. The data set has measurements on n = 442 diabetes patients, and p = 11 different measurements are taken for each patients. These measurements are:

- age
- sex
- body mass index (bmi)
- mean arterial blood pressure (map)
- six blood serum measurements: total cholesterol (tc), ldl cholesterol (ldl), hdl cholesterol (hdl), tch, ltg, glu.
- a quantitative measurement of disease progression one year after baseline (prog)

We will look at the four variables bmi, map, tc and ldl. Can we assume that these follow a multivariate normal distribution?

Contours of multivariate normal distribution

Contours of constant density for the *p*-dimensional normal distribution are ellipsoids defined by *x* such that

$$(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}) = b$$

where b > 0 is a constant.

These ellipsoids are centered at $\boldsymbol{\mu}$ and have axes $\pm \sqrt{b\lambda_i} \boldsymbol{e}_i$, where $\boldsymbol{\Sigma} \boldsymbol{e}_i = \lambda_i \boldsymbol{e}_i$, for i = 1, ..., p.

- $(\mathbf{x} \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} \boldsymbol{\mu})$ is distributed as χ^2_{ρ} .
- ▶ The volume inside the ellipsoid of *x* values satisfying

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \le \chi_p^2(\alpha)$$

has probability $1 - \alpha$.

Example: Slightly modified version of Exam K2014 1b

Let
$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$$
 be a bivariate normal random vector with mean
 $\boldsymbol{\mu} = \mathrm{E}(\mathbf{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix
 $\mathbf{\Sigma} = \mathrm{Cov}(\mathbf{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$.

You find the eigenvalues and eigenvectors of the covariance matrix $\pmb{\Sigma}$ on the next slide.

Describe the graph of the equation $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = b$ where b > 0 is a constant. Make a drawing of the graph, for b = 1 found above. What is the probability that a random sample from this distribution will be inside this graph?

Example: Exam K2014 1b

```
> sigma <- matrix(c(1,0.5,0.5,2),ncol=2)
> eigen(sigma)
$values
[1] 2.2071068 0.7928932
```

\$vectors

[,1] [,2] [1,] 0.3826834 -0.9238795 [2,] 0.9238795 0.3826834



bmi and map

bmi and







9 / 20

Multivariate distributions - in 3D: task for the intermission! Let $\mathbf{\Sigma} = \begin{bmatrix} \sigma_x^2 & \rho \sigma_x \sigma_y \\ \rho \sigma_x \sigma_y & \sigma_y^2 \end{bmatrix}$.

The following four 3D-printed figures have been made:

• D:
$$\sigma_x = 1$$
, $\sigma_y = 2$, $\rho = 0$

The figures have the following colours:

- white
- purple
- red
- black

Task: match letter and colour by writing the correct letter after the name of the colour on the available sheets and take the sheet with you. We report on the solution after the intermission.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $A\Sigma B^{T} = \mathbf{0} \Leftrightarrow AX$ and BX are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $X_2 \mid (X_1 = x_1) \sim N_{p2}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 \mu_1), \Sigma_{22} \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}).$

Example: Exam K2014 1a

Let $\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ be a bivariate normal random vector with mean $\boldsymbol{\mu} = \mathrm{E}(\mathbf{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix $\mathbf{\Sigma} = \mathrm{Cov}(\mathbf{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$.

Let $\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$, where $Y_1 = 3X_1 - 2X_2$ and $Y_2 = X_1 + X_2$. What is the distribution of \mathbf{Y} ?

What is the distribution of Y_1 ?

Example: Exam K2014 1a (slightly modified)

Let $\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ be a bivariate normal random vector with mean $\boldsymbol{\mu} = \mathrm{E}(\mathbf{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix $\mathbf{\Sigma} = \mathrm{Cov}(\mathbf{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$. Let $\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$, where $Y_1 = 3X_1 - 2X_2$ and $Y_2 = X_1 + X_2$. What is the distribution of \mathbf{Y} ?

What is the distribution of Y_1 ?

Example: Exam K2014 1a (slightly modified)

Let $\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ be a bivariate normal random vector with mean $\mu = E(\mathbf{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix $\mathbf{\Sigma} = Cov(\mathbf{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$. Let $\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$, where $Y_1 = 3X_1 - 2X_2$ and $Y_2 = X_1 + X_2$. What is the distribution of \mathbf{Y} ?

What is the distribution of Y_1 ?

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $A\Sigma B^{T} = \mathbf{0} \Leftrightarrow AX$ and BX are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $X_2 \mid (X_1 = x_1) \sim N_{p2}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 \mu_1), \Sigma_{22} \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}).$

Independent variables?

Let $oldsymbol{X}_{p imes 1} \sim N_p(oldsymbol{\mu}, oldsymbol{\Sigma})$, with

$$\boldsymbol{\Sigma} = \left[\begin{array}{rrrr} 2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 1 \\ 0 & 0 & 2 & 1 \\ 0 & 1 & 1 & 2 \end{array} \right]$$

List the pairs of variables that are independent.

Example: Exam K2014 1a - cont.

Let
$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$$
 be a bivariate normal random vector with mean
 $\boldsymbol{\mu} = \mathrm{E}(\mathbf{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix
 $\mathbf{\Sigma} = \mathrm{Cov}(\mathbf{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$.
Let $\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$, where $Y_1 = 3X_1 - 2X_2$ and $Y_2 = X_1 + X_2$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $A\Sigma B^{T} = \mathbf{0} \Leftrightarrow AX$ and BX are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $X_2 \mid (X_1 = x_1) \sim N_{p2}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 \mu_1), \Sigma_{22} \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}).$

Example: Exam V2010, Problem 1

Let
$$\boldsymbol{X} = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} \sim N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
 where $\boldsymbol{\mu} = \begin{pmatrix} 4 \\ -3 \\ 1 \end{pmatrix}$ and $\boldsymbol{\Sigma} = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 1 & -1.5 \\ 0 & -1.5 & 5 \end{pmatrix}$.

a) Find the distribution of $X_1 + X_2 + X_3$ and of X_2 given $X_1 = x_1$ and $X_3 = x_3$.

Help: for
$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} \sim N\left(\begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \mathbf{\Sigma}_{11} & \mathbf{\Sigma}_{12} \\ \mathbf{\Sigma}_{21} & \mathbf{\Sigma}_{22} \end{pmatrix}\right)$$
 we have $\mathbf{X}_2 \mid (\mathbf{X}_1 = \mathbf{x}_1) \sim N(\boldsymbol{\mu}_2 + \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1), \mathbf{\Sigma}_{22} - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12})$

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition). [CompEx1.1b]
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF). [CompEx1.1a]
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF). [CompEx1.1a]
- 5. $A\Sigma B^{T} = 0 \Leftrightarrow AX$ and BX are independent (will be very important in Part 2). [CompEx1.2b]
- 6. The conditional distributions of the components are (multivariate) normal. $X_2 \mid (X_1 = x_1) \sim N_{p2}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 \mu_1), \Sigma_{22} \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}).$