# 3. Multivariate Normal Distribution

The MVN distribution is a generalization of the univariate normal distribution which has the density function (p.d.f.)

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} \qquad -\infty < x < \infty$$

where $\mu$ = mean of distribution, $\sigma^2$ = variance. In $p-$dimensions the density becomes

$$f(\boldsymbol{x}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right\} \tag{3.1}$$

Within the mean vector $\boldsymbol{\mu}$ there are $p$ (independent) parameters and within the symmetric co-variance matrix $\boldsymbol{\Sigma}$ there are $\frac{1}{2}p(p+1)$ independent parameters [ $\frac{1}{2}p(p+3)$ independent parameters in total]. We use the notation

$$\boldsymbol{x} \sim \boldsymbol{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \tag{3.2}$$

to denote a RV $\boldsymbol{x}$ having the $p-$variate MVN distribution with

$$\mathbb{E}(\boldsymbol{x}) = \boldsymbol{\mu}$$
$$Cov(\boldsymbol{x}) = \boldsymbol{\Sigma}$$

Note that MVN distributions are entirely characterized by the first and second moments of the distribution.

## 3.1 Basic properties

If $\boldsymbol{x}$ $(p \times 1)$is MVN with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$

- Any linear combination of $\boldsymbol{x}$ is MVN

  Let $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} + \boldsymbol{c}$ with $\boldsymbol{A}\,(q \times p)$ and $\boldsymbol{c}\,(q \times 1)$ then

  $$\boldsymbol{y} \sim \boldsymbol{N}_q\left(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y\right)$$

  where $\boldsymbol{\mu}_y = \boldsymbol{A}\boldsymbol{\mu} + \boldsymbol{c}$ and $\boldsymbol{\Sigma}_y = \boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^T$.

- Any subset of variables in $\boldsymbol{x}$ has a MVN distribution.

- If a set of variables is uncorrelated, then they are independently distributed. In particular

  i) if $\sigma_{ij} = 0$ then $x_i, x_j$ are independent.

ii) if $\boldsymbol{x}$ is MVN with covariance matrix $\boldsymbol{\Sigma}$, then $\boldsymbol{Ax}$ and $\boldsymbol{Bx}$ are independent if and only if

$$Cov\left(\boldsymbol{Ax}, \boldsymbol{Bx}\right) = \boldsymbol{A\Sigma B}^T \qquad (3.3)$$
$$= \boldsymbol{0}$$

- Conditional distributions are MVN.

### Result

For the MVN distribution, variable are uncorrelated $\Leftrightarrow$ variable are independent.

### Proof

Let $\boldsymbol{x}$ $(p \times 1)$ be partitioned as

$$\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \end{bmatrix} \begin{matrix} q \\ p-q \end{matrix}$$

with mean vector

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \begin{matrix} q \\ p-q \end{matrix}$$

and covariance matrix

$$\boldsymbol{\Sigma} = \begin{matrix} q & p-q \\ \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} & \begin{matrix} q \\ p-q \end{matrix} \end{matrix}$$

i) Independent $\Rightarrow$ uncorrelated (always holds).

Suppose $\boldsymbol{x}_1, \boldsymbol{x}_2$ are independent. Then $f\left(\boldsymbol{x}_1, \boldsymbol{x}_2\right) = h\left(\boldsymbol{x}_1\right) g\left(\boldsymbol{x}_2\right)$ is a factorization of the multivariate p.d.f.and $\boldsymbol{\Sigma}_{12} = Cov\left(\boldsymbol{x}_1, \boldsymbol{x}_2\right) = \mathbb{E}\left[\left(\boldsymbol{x}_1 - \boldsymbol{\mu}_1\right)\left(\boldsymbol{x}_2 - \boldsymbol{\mu}_2\right)^T\right]$ factorizes into the product of $\mathbb{E}\left[\left(\boldsymbol{x}_1 - \boldsymbol{\mu}_1\right)\right]$ and $\mathbb{E}\left[\left(\boldsymbol{x}_2 - \boldsymbol{\mu}_2\right)^T\right]$ which are both zero since $\mathbb{E}\left(\boldsymbol{x}_1\right) = \boldsymbol{\mu}_1$ and $\mathbb{E}\left(\boldsymbol{x}_2\right) = \boldsymbol{\mu}_2$. Hence $\boldsymbol{\Sigma}_{12} = 0$.

ii) Uncorrelated $\Rightarrow$ independent (for MVN)

This result depends on factorizing the p.d.f. (3.1) when $\boldsymbol{\Sigma}_{12} = 0$.

In this case $\left(\boldsymbol{x} - \boldsymbol{\mu}\right)^T \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{x} - \boldsymbol{\mu}\right)$ has the partitioned form

$$\begin{bmatrix} \boldsymbol{x}_1^T - \boldsymbol{\mu}_1^T, & \boldsymbol{x}_2^T - \boldsymbol{\mu}_2^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Sigma}_{22} \end{bmatrix}^{-1} \begin{bmatrix} \boldsymbol{x}_1 - \boldsymbol{\mu}_1 \\ \boldsymbol{x}_2 - \boldsymbol{\mu}_2 \end{bmatrix}$$

$$= \begin{bmatrix} \boldsymbol{x}_1^T - \boldsymbol{\mu}_1^T, & \boldsymbol{x}_2^T - \boldsymbol{\mu}_2^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_{11}^{-1} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Sigma}_{22}^{-1} \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_1 - \boldsymbol{\mu}_1 \\ \boldsymbol{x}_2 - \boldsymbol{\mu}_2 \end{bmatrix}$$

$$= \left(\boldsymbol{x}_1 - \boldsymbol{\mu}_1\right)^T \boldsymbol{\Sigma}_{11}^{-1} \left(\boldsymbol{x}_1 - \boldsymbol{\mu}_1\right) + \left(\boldsymbol{x}_2 - \boldsymbol{\mu}_2\right)^T \boldsymbol{\Sigma}_{22}^{-1} \left(\boldsymbol{x}_2 - \boldsymbol{\mu}_2\right)$$

so that $\exp\{(x - \mu)^T \Sigma^{-1}(x - \mu)\}$ factorizes into the product of $\exp\left\{(x_1 - \mu_1)^T \Sigma_{11}^{-1}(x_1 - \mu_1)\right\}$ and $\exp\left\{(x_2 - \mu_2)^T \Sigma_{22}^{-1}(x_2 - \mu_2)\right\}$.

Therefore the p.d.f. can be written as

$$f(x) = g(x_1) h(x_2)$$

proving that $x_1$ and $x_2$ are independent. ∎

## 3.2 Conditional distribution

Let $X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \begin{matrix} q \\ p - q \end{matrix}$ be a partitioned MVN random $p-$vector,

with mean $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and covariance matrix

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

The conditional distribution of $X_2$ given $X_1 = x_1$ is MVN with

$$\mathbb{E}(X_2|X_1 = x_1) = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 - \mu_1) \tag{3.4a}$$

$$Cov(X_2|X_1 = x_1) = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} \tag{3.4b}$$

**Note**: the notation $X_1$ to denote the $r.v.$ and $x_1$ to denote a specific constant value (realization of $X_1$) will be very useful here.

**Proof of 3.4a**

Define a transformation from $(X_1, X_2)$ to new variables $X_1$ and $X_2' = X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1$. This is achieved by the linear transformation

$$\begin{bmatrix} X_1 \\ X_2' \end{bmatrix} = \begin{bmatrix} I & 0 \\ -\Sigma_{21}\Sigma_{11}^{-1} & I \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \tag{3.5a}$$

$$= AX \quad \text{say.} \tag{3.5b}$$

This linear relationship shows that $X_1, X_2'$ are jointly MVN (by first property of MVN stated above.)

We now show that $X_2'$ and $X_1$ are *independent* by proving that $X_1$ and $X_2'$ are uncorrelated. *Approach 1:*

$$\begin{aligned}
Cov\left(\boldsymbol{X}_1, \boldsymbol{X}_2'\right) &= Cov\left(\boldsymbol{X}_1, \boldsymbol{X}_2 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{X}_1\right) \\
&= Cov(\boldsymbol{X}_1, \boldsymbol{X}_2) - Cov\left(\boldsymbol{X}_1, \boldsymbol{X}_1\right)\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \\
&= \boldsymbol{\Sigma}_{12} - \boldsymbol{\Sigma}_{11}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \\
&= \boldsymbol{0}
\end{aligned}$$

*Approach 2:*

In (3.3), write $\boldsymbol{A} = \begin{bmatrix} \boldsymbol{B} \\ \boldsymbol{C} \end{bmatrix}$ where $\boldsymbol{B} = \begin{bmatrix} \boldsymbol{I} & \boldsymbol{0} \end{bmatrix}$ and $\boldsymbol{C} = \begin{bmatrix} -\boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1} & \boldsymbol{I} \end{bmatrix}$

$$\begin{aligned}
Cov\left(\boldsymbol{X}_1, \boldsymbol{X}_2'\right) &= Cov\left(\boldsymbol{B}\boldsymbol{X}, \boldsymbol{C}\boldsymbol{X}\right) \\
&= \boldsymbol{B}\boldsymbol{\Sigma}\boldsymbol{C}^T \\
&= \begin{bmatrix} \boldsymbol{I} & \boldsymbol{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \begin{bmatrix} -\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \\ \boldsymbol{I} \end{bmatrix} \\
&= \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \end{bmatrix} \begin{bmatrix} -\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \\ \boldsymbol{I} \end{bmatrix} \\
&= \boldsymbol{0}
\end{aligned}$$

Since $\boldsymbol{X}_2'$ and $\boldsymbol{X}_1$ are MVN and uncorrelated they are independent. Thus

$$\begin{aligned}
\mathbb{E}\left(\boldsymbol{X}_2'|\boldsymbol{X}_1 = \boldsymbol{x}_1\right) &= \mathbb{E}\left(\boldsymbol{X}_2'\right) \\
&= \mathbb{E}\left(\boldsymbol{X}_2 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{X}_1\right) \\
&= \boldsymbol{\mu}_2 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1
\end{aligned}$$

Now, as $\boldsymbol{X}_2' = \boldsymbol{X}_2 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{X}_1$ and $\boldsymbol{X}_1 = \boldsymbol{x}_1$ is given, we have

$$\begin{aligned}
\mathbb{E}\left(\boldsymbol{X}_2|\boldsymbol{X}_1 = \boldsymbol{x}_1\right) &= \mathbb{E}\left(\boldsymbol{X}_2'|\boldsymbol{X}_1 = \boldsymbol{x}_1\right) + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{x}_1 \\
&= \boldsymbol{\mu}_2 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{x}_1 \\
&= \boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\left(\boldsymbol{x}_1 - \boldsymbol{\mu}_1\right)
\end{aligned}$$

as required.

**Proof of 3.4b**

Because $\boldsymbol{X}_2'$ is independent of $\boldsymbol{X}_1$

$$Cov\left(\boldsymbol{X}_2'|\boldsymbol{X}_1 = \boldsymbol{x}_1\right) = Cov\left(\boldsymbol{X}_2'\right)$$

The left hand side is

$$
\begin{aligned}
LHS &= Cov\left(\mathbf{X}_2'|\mathbf{X}_1 = \mathbf{x}_1\right) \\
&= Cov\left(\mathbf{X}_2 - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{x}_1|\mathbf{X}_1 = \mathbf{x}_1\right) \\
&= Cov\left(\mathbf{X}_2|\mathbf{X}_1 = \mathbf{x}_1\right)
\end{aligned}
$$

The right hand side is

$$
\begin{aligned}
RHS &= Cov\left(\mathbf{X}_2'\right) \\
&= Cov\left(\mathbf{X}_2 - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{X}_1\right) \\
&= \mathbf{\Sigma}_{22} - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12}
\end{aligned}
$$

following from the general expansion

$$
\begin{aligned}
Cov\left(\mathbf{X}_2 - \mathbf{D}\mathbf{X}_1\right) = \quad & Cov\left(\mathbf{X}_2, \mathbf{X}_2\right) - \mathbf{D}Cov\left(\mathbf{X}_1, \mathbf{X}_2\right) \\
& -Cov\left(\mathbf{X}_2, \mathbf{X}_1\right)\mathbf{D}^T + \mathbf{D}Cov\left(\mathbf{X}_1, \mathbf{X}_1\right)\mathbf{D}^T
\end{aligned}
$$

with $\mathbf{D} = \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}$. Therefore

$$
Cov\left(\mathbf{X}_2|\mathbf{X}_1 = \mathbf{x}_1\right) = \mathbf{\Sigma}_{22} - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12}
$$

as required.

**Example**

Let $\mathbf{x}$ have a MVN distribution with covariance matrix

$$
\mathbf{\Sigma} = \begin{bmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & 0 \\ \rho^2 & 0 & 1 \end{bmatrix}
$$

Show that the conditional distribution of $(X_1, X_2)$ given $X_3 = x_3$ is also MVN with mean

$$
\boldsymbol{\mu} = \begin{bmatrix} \mu_1 + \rho^2\left(x_3 - \mu_3\right) \\ \mu_2 \end{bmatrix}
$$

and covariance matrix

$$
\begin{bmatrix} 1 - \rho^4 & \rho \\ \rho & 1 \end{bmatrix}
$$

**Solution**

Let $\boldsymbol{Y}_1 = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$ and $\boldsymbol{Y}_2 = (X_3)$ then

$$
\begin{aligned}
\mathbb{E}\,\boldsymbol{Y}_1 &= \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \\
\mathbb{E}\,\boldsymbol{Y}_2 &= (\mu_3).
\end{aligned}
$$

We have $Cov \begin{bmatrix} \boldsymbol{Y}_1 \\ \boldsymbol{Y}_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$ where

$$
\begin{aligned}
\boldsymbol{\Sigma}_{11} &= \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \\
\boldsymbol{\Sigma}_{12} &= \begin{bmatrix} \rho^2 \\ 0 \end{bmatrix} = \boldsymbol{\Sigma}_{21}^T \\
\boldsymbol{\Sigma}_{22} &= [1]
\end{aligned}
$$

Hence

$$
\begin{aligned}
\mathbb{E}\left[\boldsymbol{Y}_1 | \boldsymbol{Y}_2 = x_3\right] &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\left(x_3 - \mu_3\right) \\
&= \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} \rho^2 \\ 0 \end{bmatrix}\left(\boldsymbol{x}_3 - \boldsymbol{\mu}_3\right) \\
&= \begin{bmatrix} \mu_1 + \rho^2\left(x_3 - \mu_3\right) \\ \mu_2 \end{bmatrix}
\end{aligned}
$$

and .

$$
\begin{aligned}
Cov\left[\boldsymbol{Y}_1 | \boldsymbol{Y}_2 = x_3\right] &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21} \\
&= \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} - \begin{bmatrix} \rho^2 \\ 0 \end{bmatrix}\begin{bmatrix} \rho^2 & 0 \end{bmatrix} \\
&= \begin{bmatrix} 1 - \rho^4 & \rho \\ \rho & 1 \end{bmatrix}
\end{aligned}
$$

## 3.3 Maximum-likelihood estimation

Let $\boldsymbol{X}^T = (\boldsymbol{x}_1, ..., \boldsymbol{x}_n)$ contain an independent random sample of size $n$ from $N_p\left(\boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$.

The maximum likelihood estimates (MLE 's) of $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ are the sample mean and covariance matrix (with divisor $n$)

$$
\begin{aligned}
\hat{\boldsymbol{\mu}} &= \bar{\boldsymbol{x}} & \text{(3.6a)} \\
\hat{\boldsymbol{\Sigma}} &= \boldsymbol{S} & \text{(3.6b)}
\end{aligned}
$$

The likelihood function is a function of the parameters $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ given the data $\boldsymbol{X}$

$$L\left(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{X}\right) = \prod_{r=1}^{n} f\left(\boldsymbol{x}_r | \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) \tag{3.7}$$

The RHS is evaluated by substituting the individual data vectors $\{\boldsymbol{x}_1, ..., \boldsymbol{x}_n\}$ in turn into the p.d.f. of $N_p\left(\boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$ and taking the product.

$$\prod_{r=1}^{n} f\left(\boldsymbol{x}_r | \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) = (2\pi)^{-\frac{np}{2}} |\boldsymbol{\Sigma}|^{-n/2}$$

$$\exp\left\{-\frac{1}{2} \sum_{r=1}^{n} (\boldsymbol{x}_r - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_r - \boldsymbol{\mu})\right\}$$

Maximizing $L$ is equivalent to *minimizing* the "log likelihood" function

$$\begin{aligned}
l\left(\boldsymbol{\mu}, \boldsymbol{\Sigma}\right) &= -2\log L \\
&= -2\sum_{r=1}^{n} \log f\left(\boldsymbol{x}_r | \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) \\
&= K + n\log |\boldsymbol{\Sigma}| + \sum_{r=1}^{n} (\boldsymbol{x}_r - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_r - \boldsymbol{\mu}) \tag{3.8}
\end{aligned}$$

where $K$ is a constant independent of $\boldsymbol{\mu}, \boldsymbol{\Sigma}$.

**Result 3.3**

$$l\left(\boldsymbol{\mu}, \boldsymbol{\Sigma}\right) = n\left\{\log |\boldsymbol{\Sigma}| + tr\left[\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{S} + \boldsymbol{d}\boldsymbol{d}^T\right)\right]\right\} \tag{3.9}$$

up to an additive constant, where $\boldsymbol{d} = \bar{\boldsymbol{x}} - \boldsymbol{\mu}$.

Proof

Noting that $\boldsymbol{x}_r - \boldsymbol{\mu} = (\boldsymbol{x}_r - \bar{\boldsymbol{x}}) + \boldsymbol{d}$ the final term in the likelihood expression (3.8) becomes

$$\begin{aligned}
&\sum_{r=1}^{n} (\boldsymbol{x}_r - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_r - \boldsymbol{\mu}) \\
&= \sum_{r=1}^{n} (\boldsymbol{x}_r - \bar{\boldsymbol{x}})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_r - \bar{\boldsymbol{x}}) + n\boldsymbol{d}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{d} \\
&= ntr\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{S}\right) + n\boldsymbol{d}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{d} \\
&= ntr\left[\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{S} + \boldsymbol{d}\boldsymbol{d}^T\right)\right]
\end{aligned}$$

proving the expression (3.9). Note that the cross-product terms have vanished because $\sum_{r=1}^{n} \boldsymbol{x}_r =$

$n\bar{x}$ and therefore

$$
\begin{aligned}
\sum_{r=1}^{n} \boldsymbol{d}^{T}\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{x}_{r}-\overline{\boldsymbol{x}}\right) &= \boldsymbol{d}^{T}\boldsymbol{\Sigma}^{-1}\sum_{r=1}^{n}\left(\boldsymbol{x}_{r}-\overline{\boldsymbol{x}}\right) \\
&= \sum_{r=1}^{n}\left(\boldsymbol{x}_{r}-\overline{\boldsymbol{x}}\right)^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{d} \\
&= 0
\end{aligned}
$$

In (3.9) the dependence on $\boldsymbol{\mu}$ is entirely through $\boldsymbol{d}$. Now assume that is positive definite (p.d.), then so is $\boldsymbol{\Sigma}^{-1}$ as

$$
\boldsymbol{\Sigma}^{-1} = \boldsymbol{V}\boldsymbol{\Lambda}^{-1}\ \boldsymbol{V}^{T}
$$

where $\boldsymbol{\Sigma} = \boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{V}^{T}$ is the eigenanalysis of $\boldsymbol{\Sigma}$. Thus $\forall \boldsymbol{d} \neq \boldsymbol{0}$ we have $\boldsymbol{d}^{T}\boldsymbol{\Sigma}^{-1}\ \boldsymbol{d} > 0$. Hence $l\left(\boldsymbol{\mu},\boldsymbol{\Sigma}\right)$ is minimized with respect to $\boldsymbol{\mu}$ for fixed $\boldsymbol{\Sigma}$ when $\boldsymbol{d} = 0$ i.e.

$$
\hat{\boldsymbol{\mu}} = \bar{\boldsymbol{x}}
$$

*Final part of proof:* to minimize the log-likelihood $l\left(\hat{\boldsymbol{\mu}},\boldsymbol{\Sigma}\right)$ w.r.t. $\boldsymbol{\Sigma}$ let

$$
\begin{aligned}
l\left(\hat{\boldsymbol{\mu}},\boldsymbol{\Sigma}\right) &= n\left\{\log\left|\boldsymbol{\Sigma}\right| + tr\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{S}\right)\right\} \\
&= \Phi\left(\boldsymbol{\Sigma}\right)
\end{aligned}
\tag{3.10}
$$

We show that

$$
\begin{aligned}
\Phi\left(\boldsymbol{\Sigma}\right) - \Phi\left(\boldsymbol{S}\right) &= n\left\{\log\left|\boldsymbol{\Sigma}\right| - \log\left|\boldsymbol{S}\right| + tr\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{S}\right) - p\right\} \\
&= n\left\{tr\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{S}\right) - \log\left|\boldsymbol{\Sigma}^{-1}\boldsymbol{S}\right| - p\right\} \\
&\geq 0
\end{aligned}
\tag{3.11}
$$

**Lemma 1**

$\boldsymbol{\Sigma}^{-1}\boldsymbol{S}$ is positive semi-definite (proved elsewhere). Therefore the eigenvalues of $\boldsymbol{\Sigma}^{-1}\boldsymbol{S}$ are positive.

**Lemma 2**

For any set of positive numbers

$$
A \geq \log G + 1
$$

where $A$ and $G$ are the arithmetic, geometric means respectively.

**Proof**

For all $x$ we have $e^x \geq 1 + x$ (simple exercise). Consider a set of $n$ strictly positive numbers $\{y_i\}$

$$
\begin{aligned}
y_i &\geq 1 + \log y_i \\
\sum y_i &\geq n + \sum \log y_i \\
A &\geq 1 + \log \left( \prod y_i \right)^{\frac{1}{n}} \\
&= 1 + \log G
\end{aligned}
$$

as required.

Recall that for any $(n \times n)$ matrix $\boldsymbol{A}$, we have $tr\,(\boldsymbol{A}) = \sum_{i=1}^{n} \lambda_i$ the sum of the eigenvalues, and $|\,\boldsymbol{A}| = \prod \lambda_i$ the product of the eigenvalues. Let $\lambda_i$ $(i = 1, ..., p)$ be the positive eigenvalues of $\boldsymbol{\Sigma}^{-1}\boldsymbol{S}$ and substitute in $(3.11)$

$$
\begin{aligned}
\log |\boldsymbol{\Sigma}^{-1}\boldsymbol{S}| &= \log \left( \prod \lambda_i \right) \\
\\
&= p \log G
\end{aligned}
$$

$$
\begin{aligned}
tr\left( \boldsymbol{\Sigma}^{-1}\boldsymbol{S} \right) &= \sum \lambda_i \\
\\
&= pA
\end{aligned}
$$

Hence

$$
\begin{aligned}
\Phi\left( \boldsymbol{\Sigma} \right) - \Phi\left( \boldsymbol{S} \right) &= np\left\{ A - \log G - 1 \right\} \\
&\geq 0
\end{aligned}
$$

This proves that the MLE's are as stated in $(3.6)$.

## 3.3 Sampling distribution of $\bar{\mathrm{x}}$ and S

### The Wishart distribution (Definition)

*If* $\boldsymbol{M}$ $(p \times p)$ *can be written* $\boldsymbol{M} = \boldsymbol{X}^T\boldsymbol{X}$ *where* $\boldsymbol{X}$ $(m \times p)$ *is a data matrix from* $N_p\left( \boldsymbol{0}, \boldsymbol{\Sigma} \right)$ *then* $\boldsymbol{M}$ *is said to have a Wishart distribution with scale matrix* $\boldsymbol{\Sigma}$ *and degrees of freedom* $m$. *We write*

$$
\boldsymbol{M} \sim W_p\left( \boldsymbol{\Sigma}, m \right) \tag{3.12}
$$

*When* $\boldsymbol{\Sigma} = \boldsymbol{I}_p$ *the distribution is said to be in standard form.*

**Note:**

The Wishart distribution is the multivariate generalization of the chi-square $\chi^2$ distribution

**Additive property of matrices with a Wishart distribution**

Let $\boldsymbol{M}_1$, $\boldsymbol{M}_2$ be matrices having the Wishart distribution

$$\begin{aligned} \boldsymbol{M}_1 &\sim W_p\left(\boldsymbol{\Sigma}, m_1\right) \\ \boldsymbol{M}_2 &\sim W_p\left(\boldsymbol{\Sigma}, m_2\right) \end{aligned}$$

independently, then

$$\boldsymbol{M}_1 + \boldsymbol{M}_2 \sim W_p\left(\boldsymbol{\Sigma}, m_1 + m_2\right)$$

This property follows from the definition of the Wishart distribution because data matrices are additive in the sense that if

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{X}_1 \\ \boldsymbol{X}_2 \end{bmatrix}$$

is a combined data matrix consisting of $m_1 + m_2$ rows then

$$\boldsymbol{X}^T\boldsymbol{X} = \boldsymbol{X}_1^T\boldsymbol{X}_1 + \boldsymbol{X}_2^T\boldsymbol{X}_2$$

is matrix (known as the "Gram matrix") formed from the combined data matrix $\boldsymbol{X}$.

**Case of $p = 1$**

When $p = 1$ we know from the definition of $\chi_r^2$ as the distribution of the sum of squares of $r$ independent $N(0,1)$ variates that

$$\boldsymbol{M} = \sum_{i=1}^{m} x_i^2 \sim \sigma^2\chi_m^2$$

so that

$$W_1\left(\sigma^2, m\right) \equiv \sigma^2\chi_m^2$$

**Sampling distributions**

Let $\boldsymbol{x}_1, \boldsymbol{x}_2, ..., \boldsymbol{x}_n$ be a random sample of size $n$ from $N_p\left(\boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$. Then

1. The sample mean $\bar{\boldsymbol{x}}$ has the normal distribution

$$\bar{\boldsymbol{x}} \sim N_p\left(\boldsymbol{\mu}, \frac{1}{n}\boldsymbol{\Sigma}\right)$$

2. The (scaled) sample covariance matrix has the Wishart distribution:

$$(n-1)\boldsymbol{S}_u \sim W_p\left(\boldsymbol{\Sigma}, n-1\right)$$

3. The distributions of $\bar{\boldsymbol{x}}$ and $\boldsymbol{S}_u$ are independent.

## 3.4 Estimators for special circumstances

### 3.4.1 $\mu$ proportional to a given vector

Sometimes $\boldsymbol{\mu}$ is known to be proportional to a given vector, so $\boldsymbol{\mu} = k\boldsymbol{\mu}_0$ with $\boldsymbol{\mu}_0$ being a known vector.

For example if $\boldsymbol{x}$ represents a sample of repeated measurements then $\boldsymbol{\mu} = k\mathbf{1}$ where $\mathbf{1} = (1, 1, ..., 1)^T$ is the $p-$vector of $1's$.

We find the MLE of $k$ for this situation. Suppose $\boldsymbol{\Sigma}$ is known and $\boldsymbol{\mu} = k\boldsymbol{\mu}_0$. Let $\boldsymbol{d}_0 = \bar{\boldsymbol{x}} - k\boldsymbol{\mu}_0$. The log likelihood is

$$
\begin{aligned}
l(k) &= -2\log L \\
&= n\left[\log|\boldsymbol{\Sigma}| + tr\ \left\{\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{S} + \boldsymbol{d}_0\boldsymbol{d}_0^T\right)\right\}\right] \\
&= n\left[\log|\boldsymbol{\Sigma}| + tr\ \left(\boldsymbol{\Sigma}^{-1}\boldsymbol{S}\right) + (\bar{\boldsymbol{x}} - k\boldsymbol{\mu}_0)^T\boldsymbol{\Sigma}^{-1}(\bar{\boldsymbol{x}} - k\boldsymbol{\mu}_0)\right] \\
&= n\left[\bar{\boldsymbol{x}}^T\boldsymbol{\Sigma}^{-1}\bar{\boldsymbol{x}} - 2k\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\bar{\boldsymbol{x}} + k^2\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0\right] \\
&\quad + \text{constant terms indept of } k
\end{aligned}
$$

Set $\dfrac{dl}{dk} = 0$ to minimize $l(k)$ w.r.t. $k$

$$
-2\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\bar{\boldsymbol{x}} + 2\left(\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0\right)k = 0
$$

from which

$$
\hat{k} = \frac{\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\bar{\boldsymbol{x}}}{\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0} \tag{3.13}
$$

Properties

We now show that $\hat{k}$ is an unbiased estimator of $k$ and determine the variance of $\hat{k}$

In (3.13) $\hat{k}$ takes the form $\dfrac{1}{\alpha}\boldsymbol{c}^T\bar{\boldsymbol{x}}$ with $\boldsymbol{c}^T = \boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}$ and $\alpha = \boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0$ so

$$
\begin{aligned}
\mathbb{E}\left[\hat{k}\right] &= \frac{\boldsymbol{c}^T\mathbb{E}\left[\bar{\boldsymbol{x}}\right]}{\alpha} \\
&= \frac{k\boldsymbol{c}^T\boldsymbol{\mu}_0}{\alpha}. \\
&= \frac{k\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0}{\alpha}
\end{aligned}
$$

since $\mathbb{E}\left[\bar{\boldsymbol{x}}\right] = k\boldsymbol{\mu}_0$. Hence

$$
\mathbb{E}\left[\hat{k}\right] = k \tag{3.14}
$$

showing that $\hat{k}$ is an unbiased estimator.

Note that $Var\left[\bar{\boldsymbol{x}}\right] = \dfrac{1}{n}\boldsymbol{\Sigma}$ and therefore that $Var\left[\boldsymbol{c}^T\bar{\boldsymbol{x}}\right] = \dfrac{1}{n}\boldsymbol{c}^T\boldsymbol{\Sigma}\boldsymbol{c}$ we have

$$
\begin{aligned}
Var\left(\hat{k}\right) &= \frac{1}{n\alpha^2}\boldsymbol{c}^T\boldsymbol{\Sigma}\boldsymbol{c} \\
&= \frac{1}{n}\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0\left(\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0\right)^{-2} \\
&= \frac{1}{n\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_0}
\end{aligned}
\tag{3.15}
$$

### 3.4.2 Linear restriction on $\mu$

We determine an estimator for $\boldsymbol{\mu}$ to satisfy a linear restriction

$$\boldsymbol{A}\boldsymbol{\mu} = \boldsymbol{b}$$

where $\boldsymbol{A}$ $(m \times p)$ and $\boldsymbol{b}$ $(m \times 1)$ are given constants and $\boldsymbol{\Sigma}$ is assumed to be known.

We write the restriction in vector form $\boldsymbol{g}\left(\boldsymbol{\mu}\right) = \boldsymbol{0}$ and form the Lagrangean

$$\mathcal{L}\left(\boldsymbol{\mu},\boldsymbol{\lambda}\right) = l\left(\boldsymbol{\mu}\right) + 2\boldsymbol{\lambda}^T\boldsymbol{g}\left(\boldsymbol{\mu}\right)$$

where $\boldsymbol{\lambda}^T = \left(\lambda_1,...,\lambda_m\right)$ is a **vector** of Lagrange multipliers (the factor 2 is inserted just for convenience).

$$
\begin{aligned}
\mathcal{L}\left(\boldsymbol{\mu},\boldsymbol{\lambda}\right) &= l\left(\boldsymbol{\mu}\right) + 2\boldsymbol{\lambda}^T\left(\boldsymbol{A}\boldsymbol{\mu} - \boldsymbol{b}\right) \\
&= n\left\{\left(\bar{\boldsymbol{x}} - \boldsymbol{\mu}\right)^T\boldsymbol{\Sigma}^{-1}\left(\bar{\boldsymbol{x}} - \boldsymbol{\mu}\right) + 2\boldsymbol{\lambda}^T\left(\boldsymbol{A}\boldsymbol{\mu} - \boldsymbol{b}\right)\right\}
\end{aligned}
$$
ignore constant terms involving $\boldsymbol{\Sigma}$

Set $\dfrac{d}{d\boldsymbol{\mu}}\mathcal{L}\left(\boldsymbol{\mu},\boldsymbol{\lambda}\right) = \boldsymbol{0}$ using results from Example Sheet 2:

$$
\begin{aligned}
-2\boldsymbol{\Sigma}^{-1}\left(\bar{\boldsymbol{x}} - \boldsymbol{\mu}\right) + 2\boldsymbol{A}^T\boldsymbol{\lambda} &= \boldsymbol{0} \\
\bar{\boldsymbol{x}} - \boldsymbol{\mu} &= \boldsymbol{\Sigma}\boldsymbol{A}^T\boldsymbol{\lambda}
\end{aligned}
\tag{3.16}
$$

We use the constraint $\boldsymbol{A}\boldsymbol{\mu} = \boldsymbol{b}$ to evaluate the Lagrange multipliers $\boldsymbol{\lambda}$. Premultiply by $\boldsymbol{A}$

$$
\begin{aligned}
\boldsymbol{A}\bar{\boldsymbol{x}} - \boldsymbol{b} &= \boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^T\boldsymbol{\lambda} \\
\boldsymbol{\lambda} &= \left(\boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^T\right)^{-1}\left(\boldsymbol{A}\bar{\boldsymbol{x}} - \boldsymbol{b}\right)
\end{aligned}
$$

Substitute into (3.16)

$$\boxed{\hat{\boldsymbol{\mu}} = \bar{\boldsymbol{x}} - \boldsymbol{\Sigma}\boldsymbol{A}^T\left(\boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^T\right)^{-1}\left(\boldsymbol{A}\bar{\boldsymbol{x}} - \boldsymbol{b}\right)} \tag{3.17}$$

### 3.4.3 Covariance matrix $\Sigma$ proportional to a given matrix

We consider estimating $k$ when $\boldsymbol{\Sigma} = k\boldsymbol{\Sigma}_0$, where $\boldsymbol{\Sigma}_0$ is a given.constant matrix. The likelihood (3.8) takes the form when $\boldsymbol{d} = \boldsymbol{0}$ ($\hat{\boldsymbol{\mu}} = \bar{\boldsymbol{x}}$)

$$l(k) = n\left\{ \log|k\boldsymbol{\Sigma}_0| + tr\left(\frac{1}{k}\boldsymbol{\Sigma}_0^{-1}\boldsymbol{S}\right) \right\}$$

plus constant terms (not involving $k$).

$$
\begin{aligned}
l(k) &= \left\{ p\log k + \frac{1}{k}tr\left(\boldsymbol{\Sigma}_0^{-1}\boldsymbol{S}\right) \right\} \\
&\quad + \text{ constant terms} \\
\frac{dl}{dk} &= 0 \Rightarrow \frac{p}{k} - \frac{1}{k^2}tr\left(\boldsymbol{\Sigma}_0^{-1}\boldsymbol{S}\right) = 0
\end{aligned}
$$

Hence

$$\boxed{\hat{k} = \frac{tr\left(\boldsymbol{\Sigma}_0^{-1}\boldsymbol{S}\right)}{p}} \tag{3.18}$$