

Example Multivariate Gaussian:  $N_p(\mu, \Sigma)$

If  $\Sigma > 0$  it has a density

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

$$x = (x_1, \dots, x_p)^T \in \mathbb{R}^p$$

with parameters  $\mu \in \mathbb{R}^p$ ,  $\Sigma \in \mathbb{R}^{p \times p}$ ,  $\Sigma > 0$ .

### 3.2. Expectation and Covariance

Given some r.v.  $X = (X_1, \dots, X_p)^T$ .

D. 3.1. a)  $E(X) = (E(X_1), \dots, E(X_p))^T$  is called expectation (vector).

b)  $\text{Cov}(X) = E[(X - E(X))(X - E(X))^T]$  is called covariance matrix. 】

o Covariance matrix has as its  $(i,j)$ -th component

$$\begin{aligned} (\text{Cov}(X))_{i,j} &= E[(X_i - E(X_i))(X_j - E(X_j))] \\ &= \text{Cov}(X_i, X_j) \end{aligned}$$

T. 3.2.  $X = (X_1, \dots, X_p)^T, Y = (Y_1, \dots, Y_p)^T$

- a)  $E(AX+b) = A E(X) + b$
- b)  $E(X+Y) = E(X) + E(Y)$
- c)  $\text{Cov}(AX+b) = A \text{Cov}(X) A^T$
- d)  $\text{Cov}(X+Y) = \text{Cov}(X) + \text{Cov}(Y)$ , if  $X$  and  $Y$  are stoch. independent.
- e)  $\text{Cov}(X) \geq 0$  (u.u.d.) [

Proof. Ex a)-d)

$$\begin{aligned} e) \quad a^T \text{Cov}(X) a &= \text{Cov}(a^T X) && \forall a \in \mathbb{R}^p \\ &= \text{Var}(a^T X) \geq 0 \end{aligned}$$

- o Show the following If  $X \sim N_p(\mu, \Sigma)$ , then  $E(X) = \mu$ ,  $\text{Cov}(X) = \Sigma$ .

T. 3.3. (Steiner's rule)

Given a r.v.  $X = (X_1, \dots, X_p)^T$ . It holds

$$E((X-b)(X-b)^T) = \underbrace{\text{Cov}(X)}_{\text{covariance}} + \underbrace{(b-E(X))(b-E(X))^T}_{\text{bias}} \quad \forall b \in \mathbb{R}^p$$
[

Proof. Let  $\mu = EX$

$$\begin{aligned} & E[(X-\mu+b)(X-\mu+b)^T] \\ &= E((X-\mu)(X-\mu)^T) + E((\mu-b)(\mu-b)^T) \\ &= \text{Cov}(X) + (EX-b)(EX-b)^T \end{aligned}$$

since  $E((X-\mu)(\mu-b)^T) = 0$  and  $E(a) = a \quad \forall a \in \mathbb{R}^p$ .

■

Theorem 3.4. Let  $X$  be a random variate vector with  $EX = \mu$  and  $\text{Cov}(X) = V$  then

$$P(X \in \text{Im}(V) + \mu) = 1 \quad \left[ \begin{array}{cc} 1 & 1 \\ 1 & 1 \end{array} \right] \left[ \begin{array}{c} -1 \\ -1 \end{array} \right] = \left[ \begin{array}{c} 0 \\ 0 \end{array} \right]$$

Proof: Let  $\text{Ker}(V) = \{x \in \mathbb{R}^p \mid Vx = 0\}$ :

Assume a basis for  $\text{Ker}(V) = \langle a_1, \dots, a_r \rangle$

$$Va_i = 0 \quad i=1, \dots, r$$

$$a_i^T X = \underbrace{\text{Var}(a_i^T X)}_{\text{Cov}(X)} = a_i^T \underbrace{\text{Var}(X)}_{\text{Cov}(X)} a_i = a_i^T Va_i = 0$$

almost surely  $a_i^T X$  is equal to its expectation.

$$P(a_i^T X = a_i^T \mu) = 1 \Rightarrow P(a_i^T (X - \mu) = 0) = 1$$

$$\forall i \in \{1, \dots, r\}$$

$$X - \mu \perp a_i \quad a.s.$$

- FBDA 3.11.2017

$$\mathbb{P}(X - \mu \in a_i^\perp) = 1 \quad \forall i \in \{1, \dots, r\}$$

$$\mathbb{P}(X - \mu \in \underbrace{a_1^\perp \cap a_2^\perp \dots \cap a_r^\perp}_{\text{Ker}(\bar{V})^\perp}) = 1$$

$$\begin{aligned} \text{Im}(V) &= \text{Ker}(V)^\perp = \langle a_1, \dots, a_r \rangle^\perp \\ &= a_1^\perp \cap a_2^\perp \dots \cap a_r^\perp \end{aligned}$$

$$\Rightarrow \mathbb{P}(X - \mu \in \text{Im}(V)) = 1.$$

$$\begin{array}{ll} P_x = 0 & P_y \neq 0 \Rightarrow P_y \in \text{Im}(V) \\ y \notin \text{Ker}(V) & x^T P y = 0 \Rightarrow (x^T P)y = (P^T x)^T y \\ & = (Px)^T y \\ & = 0 \end{array}$$

### 3.3. Conditional distribution

Let  $X = (X_1, \dots, X_p)^T$  be a random vector s.t.

$X = (Y_1, Y_2)^T$  where  $Y_1 = (X_1, \dots, X_K)$  and

$$Y_2 = (X_{K+1}, \dots, X_p)$$

Suppose that  $X$  is absolutely continuous with density  $f_X = f_{Y_1, Y_2}$ . Then the conditional density of  $Y_1$  given  $Y_2 = y_2$  is defined by

$$f_{Y_1|Y_2}(y_1|y_2) = \frac{f_{Y_1, Y_2}(y_1, y_2)}{f_{Y_2}(y_2)}$$

where  $y_1 \in \mathbb{R}^K$  is a realization of  $Y_1$ .

$$\mathbb{P}(Y_1 \in B | Y_2 = y_2) = \int_B f_{Y_1|Y_2}(y_1|y_2) dy_1, \forall B \in \mathbb{R}^K$$

Theorem 3.5. Suppose that  $(Y_1, Y_2) \sim N_p(\mu, \Sigma)$

where

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \quad \Lambda = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ -\Lambda_{21} & \Lambda_{22} \end{bmatrix}$$

$$\Lambda = \Sigma^{-1}.$$

(a) The distribution of  $Y_1$  and  $Y_2$  are given by

$$Y_1 \sim N_K(\mu_1, \Sigma_{11}) \quad \text{and} \quad Y_2 \sim N_{p-K}(\mu_2, \Sigma_{22})$$

(b) The conditional density  $f_{Y_1|Y_2}(y_1|y_2)$  is given by

multivariate normal distribution

$$f_{Y_1|Y_2}(y_1|y_2) \sim N_K(\mu_{1|2}, \Sigma_{1|2}).$$

$$\begin{aligned}
 \mu_{1|2} &= \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (y_2 - \mu_2) \\
 &= \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (y_2 - \mu_2) \\
 &= \Sigma_{1|2} (\Lambda_{11} \mu_1 - \Lambda_{12} (y_2 - \mu_2))
 \end{aligned}$$

and

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} = \Lambda_{11}^{-1}$$

the Schur complement.

Ex. Bivariate normal distribution

### 3.4. Maximum Likelihood Estimation

Suppose  $\mathbf{x} = (x_1, \dots, x_n)$  is a random variable sampled from a p.d.f  $f(\mathbf{x}; \boldsymbol{\theta})$  where  $\boldsymbol{\theta}$  is a ~~param~~ parameter.

The function  $L(\mathbf{x}; \boldsymbol{\theta})$  is referred to as the likelihood function and defined as

$$L(\mathbf{x}; \boldsymbol{\theta}) = \prod_{i=1}^n f(x_i; \boldsymbol{\theta})$$

Furthermore, the log-likelihood function  $l(x; \vartheta)$  is defined as

$$l(x; \vartheta) = \log L(x; \vartheta) = \sum_{i=1}^n \log f(x_i; \vartheta)$$

For a given sample  $(x_1, \dots, x_n)$ , one can derive these functions, dependent on  $\vartheta$ , and select  $\vartheta$  such that they are maximized.

$$\hat{\vartheta} = \arg \max_{\vartheta} l(x; \vartheta)$$

$\hat{\vartheta}$  is called maximum likelihood estimation (MLE) of  $\vartheta$ .

Theorem 3.6. Let  $x_1, \dots, x_n$  be i.i.d. samples from the distribution  $X \sim N_p(\mu, \Sigma)$ . The MLE of  $\mu$  and  $\Sigma$  are given by

~~$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$~~ and

$$\hat{\Sigma} = S_n = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T.$$

Proof.

$$l(x_1, \dots, x_n; \mu, \Sigma) = \prod_{i=1}^n \log f(x_i; \mu, \Sigma). \quad (*)$$

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

$$\begin{aligned} \log f(x; \mu, \Sigma) &= \underbrace{\log \frac{1}{(2\pi)^{p/2}}}_{\text{---}} - \frac{1}{2} \log |\Sigma| \\ &\quad + \underbrace{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}_{\text{---}} \end{aligned}$$

$$(\hat{\mu}, \hat{\Sigma}) = \arg \max_{\mu, \Sigma} l(x_1, \dots, x_n; \mu, \Sigma)$$

$$= \arg \max_{\mu, \Sigma} \sum_{i=1}^n \left( -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)$$

$$= \arg \min_{\mu, \Sigma} \frac{n}{2} \log |\Sigma| + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$\sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \sum_{i=1}^n \text{Tr}((x_i - \mu)^T \Sigma^{-1} (x_i - \mu))$$

$$= \sum_{i=1}^n \text{Tr}(\Sigma^{-1} (x_i - \mu)(x_i - \mu)^T)$$

$$= \text{Tr}(\Sigma^{-1} \left[ \underbrace{\sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T}_{\text{---}} \right])$$

$$(\hat{\mu}, \hat{\Sigma}) = \arg \min_{\mu, \Sigma} \frac{n}{2} \log |\Sigma| + \underbrace{\text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T \right])}_{(*)}$$

→ First maximize (\*) w.r.t.  $\mu$ .

$$\begin{aligned} & \text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T \right]) = \\ &= \text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \mu)(\mathbf{x}_i - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \mu)^T \right]) \\ &= \text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \right]) \\ &+ \text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}} - \mu) + \sum_{i=1}^n (\bar{\mathbf{x}} - \mu)(\mathbf{x}_i - \bar{\mathbf{x}})^T \right]) \quad (1) \\ &+ \text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\bar{\mathbf{x}} - \mu)(\bar{\mathbf{x}} - \mu)^T \right]) \quad (2) \end{aligned}$$

$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \rightarrow$

$$(1) \quad \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}} - \mu) = (\bar{\mathbf{x}} - \mu) \left( \sum_{i=1}^n \mathbf{x}_i - n\bar{\mathbf{x}} \right)$$

$$= (\bar{\mathbf{x}} - \mu) \left( \sum_{i=1}^n \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right) = 0$$

$$(2) \quad \text{Tr}(\Sigma^{-1} \cdot n \cdot (\bar{\mathbf{x}} - \mu)(\bar{\mathbf{x}} - \mu)^T) =$$

$$n \text{Tr}((\bar{\mathbf{x}} - \mu)^T \Sigma^{-1} (\bar{\mathbf{x}} - \mu)) \geq 0$$

$$\Sigma^{-1} \succ 0$$

(a)

$$\text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T \right])$$

$$> \text{Tr}(\Sigma^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \right]) \quad (*)$$

eq. if  $(\bar{\mathbf{x}} - \boldsymbol{\mu})^T \Sigma^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) = 0$   $nS_n$

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

$$\boldsymbol{\mu}^* = \bar{\mathbf{x}}$$

$$\arg \min_{\Sigma} \frac{n}{2} \log |\Sigma| + \frac{1}{2} \text{Tr}(\Sigma^{-1} \cdot nS_n)$$

$$= \arg \min_{\Sigma} \frac{1}{2} \log |\Sigma| + \frac{1}{2} \text{Tr}(\Sigma^{-1} S_n)$$

$$\frac{\partial}{\partial \Sigma} \ell(\Sigma) = 0$$

$$\boxed{\begin{aligned} \frac{\partial}{\partial A} (\text{Tr}(AX)) &= X^T \\ \frac{\partial}{\partial A} \log |A| &= A^{-1} \end{aligned}} \quad (*)$$

$$\frac{\partial}{\partial \Sigma^{-1}} \left( \frac{1}{2} \text{Tr}(\Sigma^{-1} S_n) \right) = \frac{1}{2} S_n^T = \frac{1}{2} S_n \quad (*)$$

$$\begin{aligned} \frac{\partial}{\partial \Sigma^{-1}} \log |\Sigma| &= \frac{\partial}{\partial \Sigma^{-1}} (-\log |\Sigma^{-1}|) \oplus \\ &= -\Sigma \end{aligned}$$

$$-\frac{\Sigma}{2} + \frac{1}{2} S_n = 0 \Rightarrow \underline{\Sigma = S_n}$$

(10)