

## Kapitel 4

# Multivariat normalfördelning

### 4.1 ?

En multivariat normalfördelad stokastisk vektorvariabel  $X = (X_1, \dots, X_p)^T$  med populationsväntevärde  $\boldsymbol{\mu}$  och (icke-singulär) populationskovariansmatris  $\boldsymbol{\Sigma}$  (notation:  $\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ) har tätheten

$$= \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} \exp(-(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})/2)$$

#### Egenskaper

- *Linjära kombinationer av komponenterna är normalfördelade* (JW, Thm 4.2, Thm 4.3, Thm 4.8)
- *Alla delvektorer har en normalfördelning* (JW, Thm 4.4).
- En populationskovarians lika med 0 implicerar oberoende mellan motsvarande komponenter (JW, Thm 4.5).
- De betingade (villkorliga) fördelningarna hos komponenterna är normalfördelade (JW, Thm 4.6).
- Den stokastiska variabeln  $(\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu})$  har en chi-tvåfördelning med  $p$  frihetsgrader,  $\chi^2(p)$ , (JW Thm 4.7a). Härav följer att sannolikheten att  $X$  finns i området

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

är  $\alpha$  (JW, Thm 4.7b)

Området kallas en konfidensellipsoid. Axlarna i ellipsoiden bestäms av egenvektorer till  $\Sigma^{-1}$  som i sin tur bestäms av egenvektorer till  $\Sigma$ .

Om  $\lambda$  och  $e$  är egenvärde och egenvektor till  $\Sigma$  så är  $1/\lambda$  och  $e$  egenvärde och egenvektor till  $\Sigma^{-1}$ .

(Om  $\Sigma e = \lambda e$  så  $e = \lambda \Sigma^{-1} e$  och  $\Sigma^{-1} e = \frac{1}{\lambda} e$ )

### Maximum likelihoodskattningar

Om de stokastiska vektorerna  $X_1, \dots, X_n$  är ett slumpmässigt stickprov från en  $\mathcal{N}_p(\mu, \Sigma)$ -fördelning gäller att

$$\hat{\mu} = \bar{X} \quad \text{och} \quad \hat{\Sigma} = S_n$$

är maximum likelihoodskattningarna av  $\mu$  och  $\Sigma$ .

### Gemensam fördelning för $\bar{X}$ och $S$

$\bar{X}$  och  $S$  är oberoende.

$\bar{X}$  har en  $\mathcal{N}(\mu, \Sigma/n)$ -fördelning

$(n-1)S$  har en Wishartfördelning med  $n-1$  frihetsgrader

### Allmänt

Det finns multivariata versioner av stora talens lag (JW, Thm 4.12) och av centrala gränsvärdessatsen (JW, Thm 4.13).

## 4.2 4.2

Ex: 2-dim normalfördelning:

Om

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$$

gäller att

$$\begin{aligned} \Sigma^{-1} &= \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}^{-1} = \frac{1}{|\Sigma|} \begin{pmatrix} \sigma_2^2 & -\rho\sigma_1\sigma_2 \\ -\rho\sigma_1\sigma_2 & \sigma_1^2 \end{pmatrix} = \\ &= \frac{1}{1-\rho^2} \begin{pmatrix} \frac{1}{\sigma_1^2} & -\frac{\rho}{\sigma_1\sigma_2} \\ -\frac{\rho}{\sigma_1\sigma_2} & \frac{1}{\sigma_2^2} \end{pmatrix} \end{aligned}$$

Egenskaper för multivariat normalfördelning:

*Linjärkombinationer*

$\mathbf{X} \in \mathcal{N}_p(\mu, \Sigma) \Rightarrow$  Varje linjärkombination

$$Y = \mathbf{a}'\mathbf{X} = \sum_{i=1}^k a_i X_i$$

är  $\mathcal{N}(\mathbf{a}'\mu, \mathbf{a}'\Sigma\mathbf{a})$  (Result 4.2)

Speciellt  $a_i = 1/k$ : Ett medelvärde av komponenterna hos en multivariat normalfördelad variabel är (univariat) normalfördelad.

Flera linjärkombinationer: Om  $\mathbf{A}$  är en  $k \times p$  matris gäller att

$$\mathbf{A}\mathbf{X} \in N_k(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}')$$

(Result 4.3)

*Uppdelning*

Speciellt gäller alltså att alla komponenterna i  $\mathbf{X}$  är normalfördelade.

### 4.3 Stickprov från multivariat normalfördelning och maximum-likelihoodskattning

Om  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  är  $n$  oberoende observationer av  $\mathbf{X} \in \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  gäller att

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n) = \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right)$$

Baserad på dessa observationer är likelihooden för  $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ :

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})}$$

Maximum-likelihood skattningarna av  $\boldsymbol{\mu}$  och  $\boldsymbol{\Sigma}$  är

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}}, \quad \hat{\boldsymbol{\Sigma}} = \frac{n-1}{n} \mathbf{S} = \mathbf{S}_n$$

### 4.4 Fördelning för $\bar{\mathbf{X}}$ och $\mathbf{S}$

*Wishartfördelning*

$W_m(\cdot, \boldsymbol{\Sigma})$  = Wishartfördelning med  $m$  frihetsgrader = fördelningen för

$$\sum_{j=1}^m \mathbf{x}_j \mathbf{x}_j'$$

där  $\mathbf{x}_j \in \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$ .

Om  $\mathbf{A}_1 \in W_{m_1}(\cdot, \boldsymbol{\Sigma})$  och  $\mathbf{A}_2 \in W_{m_2}(\cdot, \boldsymbol{\Sigma})$  så gäller att  $\mathbf{A}_1 + \mathbf{A}_2 \in W_{m_1+m_2}(\cdot, \boldsymbol{\Sigma})$ .

Om  $\mathbf{A} \in W_m(\cdot, \boldsymbol{\Sigma})$  så gäller att  $\mathbf{C}\mathbf{A}\mathbf{C}' \in W_m(\cdot, \mathbf{C}\boldsymbol{\Sigma}\mathbf{C}')$ .

Det gäller att

1.  $\bar{\mathbf{X}} \in \mathcal{N}_p(\boldsymbol{\mu}, \frac{1}{n}\boldsymbol{\Sigma})$
2.  $(n-1)\mathbf{S} \in W_p(n-1, \boldsymbol{\Sigma})$
3.  $\bar{\mathbf{X}}$  och  $\mathbf{S}$  är oberoende.

## 4.5 Asymptotiska resultat för $\bar{\mathbf{X}}$ och $\mathbf{S}$

$\bar{\mathbf{X}}$  och  $\mathbf{S}$  konvergerar i sannolikhet mot  $\boldsymbol{\mu}$  och  $\boldsymbol{\Sigma}$ .

Centrala gränsvärdessatsen:

$$\sqrt{n}(\bar{\mathbf{X}} - \boldsymbol{\mu}) \sim \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$$

$$\sqrt{n}\boldsymbol{\Sigma}^{-1/2}(\bar{\mathbf{X}} - \boldsymbol{\mu}) \sim \mathcal{N}_p(\mathbf{0}, \mathbf{I})$$

$$n(\bar{\mathbf{X}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}) \sim \chi^2(p)$$

## 4.6 Fastställande av normalfördelningsantagandet

Vi vet att om  $\mathbf{X}$  är multivariat normalfördelad så är alla marginalfördelningar (multivariat) normalfördelade. Spec:

- 1-dimensionella fördelningar är en-dimensionell normalfördelade,
- 2-dimensionella fördelningar är 2-dimensionellt normalfördelade.

### 1. Histogram

ungefär 95 % av observationerna mellan  $\boldsymbol{\mu} \pm 1.96\sigma$

Q-Q-plot: Plot av stickprovskvantiler mot förväntade kvantiler om variabeln normalfördelad: Ordna observationerna i växande ordning

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

Plotta  $x_{(j)}$  mot  $q_j$  där

$$P(Z \leq q_j) = \int_{-\infty}^{q_j} \varphi(z) dz = \Phi(q_j) = \frac{j - \frac{1}{2}}{n}$$

Om variabeln normalfördelad bör  $(x_{(j)}, q_j)$  ligga på ungefär rät linje.

### 2. Bivariat normalitet?

Sannolikheten för  $(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq \chi_{0.5}^2(2)$  är 0.5. Skatta  $\boldsymbol{\mu}$  med  $\bar{\mathbf{x}}$  och  $\boldsymbol{\Sigma}$  med  $\mathbf{S}$ .

Enklare att bestämma "generaliserade kvadratiske avstånd

$$d_j^2 = (\mathbf{x}_j - \bar{\mathbf{x}})' \mathbf{S}^{-1} (\mathbf{x}_j - \bar{\mathbf{x}})$$

för  $j = 1, 2, \dots, n$  samt konstruera ett chi-två plot:

- (a) Ordna avstånden i växande ordning

$$d_{(1)}^2 \leq d_{(2)}^2 \leq \dots \leq d_{(n)}^2$$

- (b) Plotta  $d_{(j)}^2$  mot  $\chi_{\frac{j-1}{n}}^2(p)$ , dvs  $100(j-\frac{1}{2})/n$ -percentilen i  $\chi^2$ -fördelningen med  $p$  frihetsgrader.

$$(q_{c,p}(\alpha) = \chi_{\alpha}^2(p).)$$

## 4.7 Transformation till normalitet

Några vanliga transformationer

1. Antal  $y$ :  $\sqrt{y}$  (rot-transformation)
2. Proportion  $p$ :  $\text{logit}(\hat{p}) = \frac{1}{2} \ln\left(\frac{\hat{p}}{1-\hat{p}}\right)$
3. Korrelation  $r$ : Fishers  $z = \arctan(r) = \frac{1}{2} \ln\left(\frac{1+r}{1-r}\right)$

Box-Cox:

$$x^{(\lambda)} = \begin{cases} \frac{x^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln(x) & \lambda = 0 \end{cases}$$

Kan generaliseras till multivariata observationer.