

## Chapter 3

# Multivariate Limit Theorems

This chapter discusses multivariate limit theorems, and follows Olive (2014, § 8.6, 8.7) closely.

### 3.1 Multivariate Limit Theorems

Many of the univariate results from Chapter 2 can be extended to random vectors. For the limit theorems, the vector  $\mathbf{X}$  is typically a  $k \times 1$  column vector and  $\mathbf{X}^T$  is a row vector. Let  $\|\mathbf{x}\| = \sqrt{x_1^2 + \cdots + x_k^2}$  be the Euclidean norm of  $\mathbf{x}$ .

**Definition 3.1.** Let  $\mathbf{X}_n$  be a sequence of random vectors with joint cdfs  $F_n(\mathbf{x})$  and let  $\mathbf{X}$  be a random vector with joint cdf  $F(\mathbf{x})$ .

a)  $\mathbf{X}_n$  **converges in distribution** to  $\mathbf{X}$ , written  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ , if  $F_n(\mathbf{x}) \rightarrow F(\mathbf{x})$  as  $n \rightarrow \infty$  for all points  $\mathbf{x}$  at which  $F(\mathbf{x})$  is continuous. The distribution of  $\mathbf{X}$  is the **limiting distribution** or **asymptotic distribution** of  $\mathbf{X}_n$ .

b)  $\mathbf{X}_n$  **converges in probability** to  $\mathbf{X}$ , written  $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$ , if for every  $\epsilon > 0$ ,  $P(\|\mathbf{X}_n - \mathbf{X}\| > \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$ .

c) Let  $r > 0$  be a real number. Then  $\mathbf{X}_n$  **converges in  $r$ th mean** to  $\mathbf{X}$ , written  $\mathbf{X}_n \xrightarrow{r} \mathbf{X}$ , if  $E(\|\mathbf{X}_n - \mathbf{X}\|^r) \rightarrow 0$  as  $n \rightarrow \infty$ .

d)  $\mathbf{X}_n$  **converges almost everywhere** to  $\mathbf{X}$ , written  $\mathbf{X}_n \xrightarrow{ae} \mathbf{X}$ , if  $P(\lim_{n \rightarrow \infty} \mathbf{X}_n = \mathbf{X}) = 1$ .

Theorems 3.1, 3.2 and 3.3 below are the multivariate extensions of the limit theorems in Section 2.1. When the limiting distribution of  $\mathbf{Z}_n = \sqrt{n}(\mathbf{g}(\mathbf{T}_n) - \mathbf{g}(\boldsymbol{\theta}))$  is multivariate normal  $N_k(\mathbf{0}, \boldsymbol{\Sigma})$ , approximate the joint cdf of  $\mathbf{Z}_n$  with the joint cdf of the  $N_k(\mathbf{0}, \boldsymbol{\Sigma})$  distribution. Thus to find probabilities, manipulate  $\mathbf{Z}_n$  as if  $\mathbf{Z}_n \approx N_k(\mathbf{0}, \boldsymbol{\Sigma})$ . To see that the CLT is a special case of the MCLT below, let  $k = 1$ ,  $E(X) = \mu$  and  $V(X) = \boldsymbol{\Sigma} = \sigma^2$ .

**Theorem 3.1: the Multivariate Central Limit Theorem (MCLT).**

If  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are iid  $k \times 1$  random vectors with  $E(\mathbf{X}) = \boldsymbol{\mu}$  and  $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ , then

$$\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma})$$

where the sample mean

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i.$$

The MCLT is proven after Theorem 3.8. To see that the delta method is a special case of the multivariate delta method, note that if  $T_n$  and parameter  $\theta$  are real valued, then  $\mathbf{D}_{\mathbf{g}}(\theta) = g'(\theta)$ .

**Theorem 3.2: the Multivariate Delta Method.** If

$$\sqrt{n}(\mathbf{T}_n - \boldsymbol{\theta}) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma}),$$

then

$$\sqrt{n}(\mathbf{g}(\mathbf{T}_n) - \mathbf{g}(\boldsymbol{\theta})) \xrightarrow{D} N_d(\mathbf{0}, \mathbf{D}_{\mathbf{g}}(\boldsymbol{\theta}) \boldsymbol{\Sigma} \mathbf{D}_{\mathbf{g}}^T(\boldsymbol{\theta}))$$

if  $\mathbf{D}_{\mathbf{g}}(\boldsymbol{\theta}) \boldsymbol{\Sigma} \mathbf{D}_{\mathbf{g}}^T(\boldsymbol{\theta})$  is nonsingular, where the  $d \times k$  Jacobian matrix of partial derivatives

$$\mathbf{D}_{\mathbf{g}}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} g_1(\boldsymbol{\theta}) & \dots & \frac{\partial}{\partial \theta_k} g_1(\boldsymbol{\theta}) \\ \vdots & & \vdots \\ \frac{\partial}{\partial \theta_1} g_d(\boldsymbol{\theta}) & \dots & \frac{\partial}{\partial \theta_k} g_d(\boldsymbol{\theta}) \end{bmatrix}.$$

Here the mapping  $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^d$  needs to be differentiable in a neighborhood of  $\boldsymbol{\theta} \in \mathbb{R}^k$ .

**Example 3.1.** If  $Y$  has a Weibull distribution,  $Y \sim W(\phi, \lambda)$ , then the pdf of  $Y$  is

$$f(y) = \frac{\phi}{\lambda} y^{\phi-1} e^{-\frac{y^\phi}{\lambda}}$$

where  $\lambda, y$ , and  $\phi$  are all positive. If  $\mu = \lambda^{1/\phi}$  so  $\mu^\phi = \lambda$ , then the Weibull pdf

$$f(y) = \frac{\phi}{\mu} \left(\frac{y}{\mu}\right)^{\phi-1} \exp\left[-\left(\frac{y}{\mu}\right)^\phi\right].$$

Let  $(\hat{\mu}, \hat{\phi})$  be the MLE of  $(\mu, \phi)$ . According to Bain (1978, p. 215),

$$\sqrt{n} \left( \begin{pmatrix} \hat{\mu} \\ \hat{\phi} \end{pmatrix} - \begin{pmatrix} \mu \\ \phi \end{pmatrix} \right) \xrightarrow{D} N \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1.109 \frac{\mu^2}{\phi^2} & 0.257\mu \\ 0.257\mu & 0.608\phi^2 \end{pmatrix} \right)$$

=  $N_2(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\theta}))$  where  $\mathbf{I}(\boldsymbol{\theta})$  is given in Definition 3.2.

Let column vectors  $\boldsymbol{\theta} = (\mu \ \phi)^T$  and  $\boldsymbol{\eta} = (\lambda \ \phi)^T$ . Then

$$\boldsymbol{\eta} = \mathbf{g}(\boldsymbol{\theta}) = \begin{pmatrix} \lambda \\ \phi \end{pmatrix} = \begin{pmatrix} \mu^\phi \\ \phi \end{pmatrix} = \begin{pmatrix} g_1(\boldsymbol{\theta}) \\ g_2(\boldsymbol{\theta}) \end{pmatrix}.$$

So

$$\mathbf{D}_{\mathbf{g}}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} g_1(\boldsymbol{\theta}) & \frac{\partial}{\partial \theta_2} g_1(\boldsymbol{\theta}) \\ \frac{\partial}{\partial \theta_1} g_2(\boldsymbol{\theta}) & \frac{\partial}{\partial \theta_2} g_2(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \mu} \mu^\phi & \frac{\partial}{\partial \phi} \mu^\phi \\ \frac{\partial}{\partial \mu} \phi & \frac{\partial}{\partial \phi} \phi \end{bmatrix} = \begin{bmatrix} \phi \mu^{\phi-1} & \mu^\phi \log(\mu) \\ 0 & 1 \end{bmatrix}.$$

Thus by the multivariate delta method,

$$\sqrt{n} \left( \begin{pmatrix} \hat{\lambda} \\ \hat{\phi} \end{pmatrix} - \begin{pmatrix} \lambda \\ \phi \end{pmatrix} \right) \xrightarrow{D} N_2(\mathbf{0}, \boldsymbol{\Sigma})$$

where (see Definition 3.4 below)

$$\boldsymbol{\Sigma} = \mathbf{I}(\boldsymbol{\eta})^{-1} = [\mathbf{I}(\mathbf{g}(\boldsymbol{\theta}))]^{-1} = \mathbf{D}_{\mathbf{g}}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{D}_{\mathbf{g}}^T(\boldsymbol{\theta}) = \begin{bmatrix} 1.109\lambda^2(1 + 0.4635 \log(\lambda) + 0.5482(\log(\lambda))^2) & 0.257\phi\lambda + 0.608\lambda\phi \log(\lambda) \\ 0.257\phi\lambda + 0.608\lambda\phi \log(\lambda) & 0.608\phi^2 \end{bmatrix}.$$

**Definition 3.2.** Let  $X$  be a random variable with pdf or pmf  $f(x|\boldsymbol{\theta})$ . Then the **information matrix**

$$\mathbf{I}(\boldsymbol{\theta}) = [\mathbf{I}_{i,j}]$$

where

$$\mathbf{I}_{i,j} = E \left[ \frac{\partial}{\partial \theta_i} \log(f(X|\boldsymbol{\theta})) \frac{\partial}{\partial \theta_j} \log(f(X|\boldsymbol{\theta})) \right].$$

**Definition 3.3.** An estimator  $\mathbf{T}_n$  of  $\boldsymbol{\theta}$  is **asymptotically efficient** if

$$\sqrt{n}(\mathbf{T}_n - \boldsymbol{\theta}) \xrightarrow{D} N_k(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\theta})).$$

Following Lehmann (1999, p. 511), if  $\mathbf{T}_n$  is asymptotically efficient and if the estimator  $\mathbf{W}_n$  satisfies

$$\sqrt{n}(\mathbf{W}_n - \boldsymbol{\theta}) \xrightarrow{D} N_k(\mathbf{0}, \mathbf{J}(\boldsymbol{\theta}))$$

where  $\mathbf{J}(\boldsymbol{\theta})$  and  $\mathbf{I}^{-1}(\boldsymbol{\theta})$  are continuous functions of  $\boldsymbol{\theta}$ , then under regularity conditions,  $\mathbf{J}(\boldsymbol{\theta}) - \mathbf{I}^{-1}(\boldsymbol{\theta})$  is a positive semidefinite matrix, and  $\mathbf{T}_n$  is “better” than  $\mathbf{W}_n$ .

**Definition 3.4.** Assume that  $\boldsymbol{\eta} = \mathbf{g}(\boldsymbol{\theta})$ . Then

$$\mathbf{I}(\boldsymbol{\eta}) = \mathbf{I}(\mathbf{g}(\boldsymbol{\theta})) = [\mathbf{D}_{\mathbf{g}}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{D}_{\mathbf{g}}^T(\boldsymbol{\theta})]^{-1}.$$

Notice that this definition agrees with the multivariate delta method if

$$\sqrt{n}(\mathbf{T}_n - \boldsymbol{\theta}) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma})$$

where  $\boldsymbol{\Sigma} = \mathbf{I}^{-1}(\boldsymbol{\theta})$ .

Now suppose that  $X_1, \dots, X_n$  are iid random variables from a  $k$ -parameter REF

$$f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp \left[ \sum_{i=1}^k w_i(\boldsymbol{\theta}) t_i(x) \right] \quad (3.1)$$

with natural parameterization

$$f(x|\boldsymbol{\eta}) = h(x)b(\boldsymbol{\eta}) \exp \left[ \sum_{i=1}^k \eta_i t_i(x) \right]. \quad (3.2)$$

Then the complete minimal sufficient statistic is

$$\bar{\mathbf{T}}_n = \frac{1}{n} \left( \sum_{i=1}^n t_1(X_i), \dots, \sum_{i=1}^n t_k(X_i) \right)^T.$$

Let  $\boldsymbol{\mu}_T = (E(t_1(X)), \dots, E(t_k(X)))^T$ . From Theorem 1.31, for  $\boldsymbol{\eta} \in \Omega$ ,

$$E(t_i(X)) = \frac{-\partial}{\partial \eta_i} \log(b(\boldsymbol{\eta})),$$

and

$$\text{Cov}(t_i(X), t_j(X)) \equiv \sigma_{i,j} = \frac{-\partial^2}{\partial \eta_i \partial \eta_j} \log(b(\boldsymbol{\eta})).$$

**Theorem 3.3.** If the random variable  $X$  is a kP-REF with pmf or pdf (3.2), then the information matrix

$$\mathbf{I}(\boldsymbol{\eta}) = [\mathbf{I}_{i,j}]$$

where

$$\mathbf{I}_{i,j} = E \left[ \frac{\partial}{\partial \eta_i} \log(f(X|\boldsymbol{\eta})) \frac{\partial}{\partial \eta_j} \log(f(X|\boldsymbol{\eta})) \right] = -E \left[ \frac{\partial^2}{\partial \eta_i \partial \eta_j} \log(f(X|\boldsymbol{\eta})) \right].$$

Several authors, including Barndorff-Nielsen (1982), have noted that the multivariate CLT can be used to show that  $\sqrt{n}(\bar{\mathbf{T}}_n - \boldsymbol{\mu}_T) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma})$ . The fact that  $\boldsymbol{\Sigma} = \mathbf{I}(\boldsymbol{\eta})$  appears in Lehmann (1983, p. 127).

**Theorem 3.4.** If  $X_1, \dots, X_n$  are iid from a  $k$ -parameter regular exponential family, then

$$\sqrt{n}(\bar{\mathbf{T}}_n - \boldsymbol{\mu}_T) \xrightarrow{D} N_k(\mathbf{0}, \mathbf{I}(\boldsymbol{\eta})).$$

**Proof.** By the multivariate central limit theorem,

$$\sqrt{n}(\bar{\mathbf{T}}_n - \boldsymbol{\mu}_T) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma})$$

where  $\boldsymbol{\Sigma} = [\sigma_{i,j}]$ . Hence the result follows if  $\sigma_{i,j} = \mathbf{I}_{i,j}$ . Since

$$\log(f(x|\boldsymbol{\eta})) = \log(h(x)) + \log(b(\boldsymbol{\eta})) + \sum_{l=1}^k \eta_l t_l(x),$$

$$\frac{\partial}{\partial \eta_i} \log(f(x|\boldsymbol{\eta})) = \frac{\partial}{\partial \eta_i} \log(b(\boldsymbol{\eta})) + t_i(x).$$

Hence

$$-\mathbf{I}_{i,j} = E \left[ \frac{\partial^2}{\partial \eta_i \partial \eta_j} \log(f(X|\boldsymbol{\eta})) \right] = \frac{\partial^2}{\partial \eta_i \partial \eta_j} \log(b(\boldsymbol{\eta})) = -\sigma_{i,j}. \quad \square$$

To obtain standard results, use the multivariate delta method, assume that both  $\boldsymbol{\theta}$  and  $\boldsymbol{\eta}$  are  $k \times 1$  vectors, and assume that  $\boldsymbol{\eta} = \mathbf{g}(\boldsymbol{\theta})$  is a one to one mapping so that the inverse mapping is  $\boldsymbol{\theta} = \mathbf{g}^{-1}(\boldsymbol{\eta})$ . If  $\mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})}$  is nonsingular, then

$$\mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})}^{-1} = \mathbf{D}_{\mathbf{g}^{-1}(\boldsymbol{\eta})} \quad (3.3)$$

(see Searle 1982, p. 339), and

$$\mathbf{I}(\boldsymbol{\eta}) = [\mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})} \mathbf{I}^{-1}(\boldsymbol{\theta}) \mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})}^T]^{-1} = [\mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})}^{-1}]^T \mathbf{I}(\boldsymbol{\theta}) \mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})}^{-1} = \mathbf{D}_{\mathbf{g}^{-1}(\boldsymbol{\eta})}^T \mathbf{I}(\boldsymbol{\theta}) \mathbf{D}_{\mathbf{g}^{-1}(\boldsymbol{\eta})}. \quad (3.4)$$

Compare Lehmann (1999, p. 500) and Lehmann (1983, p. 127).

For example, suppose that  $\boldsymbol{\mu}_T$  and  $\boldsymbol{\eta}$  are  $k \times 1$  vectors, and

$$\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \xrightarrow{D} N_k(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\eta}))$$

where  $\boldsymbol{\mu}_T = \mathbf{g}(\boldsymbol{\eta})$  and  $\boldsymbol{\eta} = \mathbf{g}^{-1}(\boldsymbol{\mu}_T)$ . Also assume that  $\bar{\mathbf{T}}_n = \mathbf{g}(\hat{\boldsymbol{\eta}})$  and  $\hat{\boldsymbol{\eta}} = \mathbf{g}^{-1}(\bar{\mathbf{T}}_n)$ . Then by the multivariate delta method and Theorem 3.4,

$$\sqrt{n}(\bar{\mathbf{T}}_n - \boldsymbol{\mu}_T) = \sqrt{n}(\mathbf{g}(\hat{\boldsymbol{\eta}}) - \mathbf{g}(\boldsymbol{\eta})) \xrightarrow{D} N_k[\mathbf{0}, \mathbf{I}(\boldsymbol{\eta})] = N_k[\mathbf{0}, \mathbf{D}_{\mathbf{g}(\boldsymbol{\eta})} \mathbf{I}^{-1}(\boldsymbol{\eta}) \mathbf{D}_{\mathbf{g}(\boldsymbol{\eta})}^T].$$

Hence

$$\mathbf{I}(\boldsymbol{\eta}) = \mathbf{D}_{\mathbf{g}(\boldsymbol{\eta})} \mathbf{I}^{-1}(\boldsymbol{\eta}) \mathbf{D}_{\mathbf{g}(\boldsymbol{\eta})}^T.$$

Similarly,

$$\sqrt{n}(\mathbf{g}^{-1}(\bar{\mathbf{T}}_n) - \mathbf{g}^{-1}(\boldsymbol{\mu}_T)) = \sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \xrightarrow{D} N_k[\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\eta})] =$$

$$N_k[\mathbf{0}, D\mathbf{g}^{-1}(\boldsymbol{\mu}_T)\mathbf{I}(\boldsymbol{\eta})D\mathbf{g}^{-1}(\boldsymbol{\mu}_T)^T].$$

Thus

$$\mathbf{I}^{-1}(\boldsymbol{\eta}) = D\mathbf{g}^{-1}(\boldsymbol{\mu}_T)\mathbf{I}(\boldsymbol{\eta})D\mathbf{g}^{-1}(\boldsymbol{\mu}_T)^T = D\mathbf{g}^{-1}(\boldsymbol{\mu}_T)D\mathbf{g}(\boldsymbol{\eta})\mathbf{I}^{-1}(\boldsymbol{\eta})D\mathbf{g}(\boldsymbol{\eta})D\mathbf{g}^{-1}(\boldsymbol{\mu}_T)^T$$

as expected by Equation (3.4). Typically  $\hat{\boldsymbol{\theta}}$  is a function of the sufficient statistic  $\mathbf{T}_n$  and is the unique MLE of  $\boldsymbol{\theta}$ . Replacing  $\boldsymbol{\eta}$  by  $\boldsymbol{\theta}$  in the above discussion shows that  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{D} N_k(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\theta}))$  is equivalent to  $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}_T) \xrightarrow{D} N_k(\mathbf{0}, \mathbf{I}(\boldsymbol{\theta}))$  provided that  $D\mathbf{g}(\boldsymbol{\theta})$  is nonsingular.

## 3.2 More Multivariate Results

**Definition 3.5.** If the estimator  $\mathbf{g}(\mathbf{T}_n) \xrightarrow{P} \mathbf{g}(\boldsymbol{\theta})$  for all  $\boldsymbol{\theta} \in \Theta$ , then  $\mathbf{g}(\mathbf{T}_n)$  is a **consistent estimator** of  $\mathbf{g}(\boldsymbol{\theta})$ .

**Theorem 3.5.** If  $0 < \delta \leq 1$ ,  $\mathbf{X}$  is a random vector, and

$$n^\delta(\mathbf{g}(\mathbf{T}_n) - \mathbf{g}(\boldsymbol{\theta})) \xrightarrow{D} \mathbf{X},$$

then  $\mathbf{g}(\mathbf{T}_n) \xrightarrow{P} \mathbf{g}(\boldsymbol{\theta})$ .

**Theorem 3.6.** If  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are iid,  $E(\|\mathbf{X}\|) < \infty$  and  $E(\mathbf{X}) = \boldsymbol{\mu}$ , then

- a) WLLN:  $\bar{\mathbf{X}}_n \xrightarrow{P} \boldsymbol{\mu}$  and
- b) SLLN:  $\bar{\mathbf{X}}_n \xrightarrow{ae} \boldsymbol{\mu}$ .

**Theorem 3.7: Continuity Theorem.** Let  $\mathbf{X}_n$  be a sequence of  $k \times 1$  random vectors with characteristic function  $c_n(\mathbf{t})$  and let  $\mathbf{X}$  be a  $k \times 1$  random vector with cf  $c(\mathbf{t})$ . Then

$$\mathbf{X}_n \xrightarrow{D} \mathbf{X} \text{ iff } c_n(\mathbf{t}) \rightarrow c(\mathbf{t})$$

for all  $\mathbf{t} \in \mathbb{R}^k$ .

**Theorem 3.8: Cramér Wold Device.** Let  $\mathbf{X}_n$  be a sequence of  $k \times 1$  random vectors and let  $\mathbf{X}$  be a  $k \times 1$  random vector. Then

$$\mathbf{X}_n \xrightarrow{D} \mathbf{X} \text{ iff } \mathbf{t}^T \mathbf{X}_n \xrightarrow{D} \mathbf{t}^T \mathbf{X}$$

for all  $\mathbf{t} \in \mathbb{R}^k$ .

**Application: Proof of the MCLT Theorem 3.1.** Note that for fixed  $\mathbf{t}$ , the  $\mathbf{t}^T \mathbf{X}_i$  are iid random variables with mean  $\mathbf{t}^T \boldsymbol{\mu}$  and variance  $\mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}$ . Hence by the CLT,  $\mathbf{t}^T \sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} N(0, \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t})$ . The right hand side has

distribution  $\mathbf{t}^T \mathbf{X}$  where  $\mathbf{X} \sim N_k(\mathbf{0}, \boldsymbol{\Sigma})$ . Hence by the Cramér Wold Device,  $\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma})$ .  $\square$

**Theorem 3.9.** a) If  $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$ , then  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ .

b)

$$\mathbf{X}_n \xrightarrow{P} \mathbf{g}(\boldsymbol{\theta}) \text{ iff } \mathbf{X}_n \xrightarrow{D} \mathbf{g}(\boldsymbol{\theta}).$$

Let  $g(n) \geq 1$  be an increasing function of the sample size  $n$ :  $g(n) \uparrow \infty$ , e.g.  $g(n) = \sqrt{n}$ . See White (1984, p. 15). If a  $k \times 1$  random vector  $\mathbf{T}_n - \boldsymbol{\mu}$  converges to a nondegenerate multivariate normal distribution with convergence rate  $\sqrt{n}$ , then  $\mathbf{T}_n$  has (tightness) rate  $\sqrt{n}$ .

**Definition 3.6.** Let  $\mathbf{A}_n = [a_{i,j}(n)]$  be an  $r \times c$  random matrix.

a)  $\mathbf{A}_n = O_P(\mathbf{X}_n)$  if  $a_{i,j}(n) = O_P(\mathbf{X}_n)$  for  $1 \leq i \leq r$  and  $1 \leq j \leq c$ .

b)  $\mathbf{A}_n = o_P(\mathbf{X}_n)$  if  $a_{i,j}(n) = o_P(\mathbf{X}_n)$  for  $1 \leq i \leq r$  and  $1 \leq j \leq c$ .

c)  $\mathbf{A}_n \asymp_P (1/g(n))$  if  $a_{i,j}(n) \asymp_P (1/g(n))$  for  $1 \leq i \leq r$  and  $1 \leq j \leq c$ .

d) Let  $\mathbf{A}_{1,n} = \mathbf{T}_n - \boldsymbol{\mu}$  and  $\mathbf{A}_{2,n} = \mathbf{C}_n - c\boldsymbol{\Sigma}$  for some constant  $c > 0$ . If  $\mathbf{A}_{1,n} \asymp_P (1/g(n))$  and  $\mathbf{A}_{2,n} \asymp_P (1/g(n))$ , then  $(\mathbf{T}_n, \mathbf{C}_n)$  has (tightness) rate  $g(n)$ .

**Theorem 3.10.** Let  $W_n, X_n, Y_n$  and  $Z_n$  be sequences of random variables such that  $Y_n > 0$  and  $Z_n > 0$ . (Often  $Y_n$  and  $Z_n$  are deterministic, e.g.  $Y_n = n^{-1/2}$ .)

a) If  $W_n = O_P(1)$  and  $X_n = O_P(1)$ , then  $W_n + X_n = O_P(1)$  and  $W_n X_n = O_P(1)$ , thus  $O_P(1) + O_P(1) = O_P(1)$  and  $O_P(1)O_P(1) = O_P(1)$ .

b) If  $W_n = O_P(1)$  and  $X_n = o_P(1)$ , then  $W_n + X_n = O_P(1)$  and  $W_n X_n = o_P(1)$ , thus  $O_P(1) + o_P(1) = O_P(1)$  and  $O_P(1)o_P(1) = o_P(1)$ .

c) If  $W_n = O_P(Y_n)$  and  $X_n = O_P(Z_n)$ , then  $W_n + X_n = O_P(\max(Y_n, Z_n))$  and  $W_n X_n = O_P(Y_n Z_n)$ , thus  $O_P(Y_n) + O_P(Z_n) = O_P(\max(Y_n, Z_n))$  and  $O_P(Y_n)O_P(Z_n) = O_P(Y_n Z_n)$ .

Recall that the smallest integer function  $[x]$  rounds up, e.g.  $[7.7] = 8$ .

**Definition 3.7.** The *sample  $\rho$  quantile*  $\hat{y}_{n,\rho} = \hat{\xi}_{n,\rho} = Y_{([\!n\rho])}$ . The *population quantile*  $y_\rho = \xi_\rho = Q(\rho) = \inf\{y : F(y) \geq \rho\}$ .

There are many other ways to define sample quantiles, and the different estimators tend to be asymptotically equivalent. If the inverse  $F^{-1}$  of the cdf exists, then  $Q(u) = F^{-1}(u)$ .  $Q(u) \leq x$  iff  $u \leq F(x)$ .  $F(y_\rho) = P(Y \leq y_\rho) \geq \rho$  and  $P(Y \geq y_\rho) \geq 1 - \rho$ . Let the observed data be  $Y_1, \dots, Y_n$ , and let  $\hat{F}(y) = \frac{1}{n} \sum_{i=1}^n I(Y_i \leq y)$ . Then  $\hat{Q}(\rho) = \inf\{y : \hat{F}(y) \geq \rho\} = \hat{y}_{n,\rho} = Y_{([\!n\rho])}$ . (An alternative definition of the population quantile that is often used is that  $y_\rho$  is any real number satisfying  $P(Y \leq y_\rho) \geq \rho$  and  $P(Y \geq y_\rho) \geq 1 - \rho$ . Then  $y_\rho$  is not necessarily unique. Definition 3.7 makes the population quantile

unique. The regularity conditions in Theorem 3.11 make  $y_\rho$  unique if the alternative definition is used.)

**Theorem 3.11: Serfling (1980, p. 80).** Let  $0 < \rho_1 < \rho_2 < \cdots < \rho_k < 1$ . Suppose that  $F$  has a density  $f$  that is positive and continuous in neighborhoods of  $\xi_{\rho_1}, \dots, \xi_{\rho_k}$ . Then

$$\sqrt{n}[(\hat{\xi}_{n,\rho_1}, \dots, \hat{\xi}_{n,\rho_k})^T - (\xi_{\rho_1}, \dots, \xi_{\rho_k})^T] \xrightarrow{D} N_k(\mathbf{0}, \Sigma)$$

where  $\Sigma = (\sigma_{ij})$  and

$$\sigma_{ij} = \frac{\rho_i(1 - \rho_j)}{f(\xi_{\rho_i})f(\xi_{\rho_j})}$$

for  $i \leq j$  and  $\sigma_{ij} = \sigma_{ji}$  for  $i > j$ .

**Theorem 3.12: Continuous Mapping Theorem.** Let  $\mathbf{X}_n \in \mathbb{R}^k$ . If  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$  and if the function  $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^j$  is continuous and does not depend on  $n$ , then  $\mathbf{g}(\mathbf{X}_n) \xrightarrow{D} \mathbf{g}(\mathbf{X})$ .

The following two theorems are taken from Severini (2005, pp. 345-349, 354).

**Theorem 3.13.** Let  $\mathbf{X}_n = (X_{1n}, \dots, X_{kn})^T$  be a sequence of  $k \times 1$  random vectors, let  $\mathbf{Y}_n$  be a sequence of  $k \times 1$  random vectors, and let  $\mathbf{X} = (X_1, \dots, X_k)^T$  be a  $k \times 1$  random vector. Let  $\mathbf{W}_n$  be a sequence of  $k \times k$  nonsingular random matrices, and let  $\mathbf{C}$  be a  $k \times k$  constant nonsingular matrix.

- a)  $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$  iff  $X_{in} \xrightarrow{P} X_i$  for  $i = 1, \dots, k$ .
- b) **Slutsky's Theorem:** If  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$  and  $\mathbf{Y}_n \xrightarrow{P} \mathbf{c}$  for some constant  $k \times 1$  vector  $\mathbf{c}$ , then i)  $\mathbf{X}_n + \mathbf{Y}_n \xrightarrow{D} \mathbf{X} + \mathbf{c}$  and ii)  $\mathbf{Y}_n^T \mathbf{X}_n \xrightarrow{D} \mathbf{c}^T \mathbf{X}$ .
- c) If  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$  and  $\mathbf{W}_n \xrightarrow{P} \mathbf{C}$ , then  $\mathbf{W}_n \mathbf{X}_n \xrightarrow{D} \mathbf{C} \mathbf{X}$ ,  $\mathbf{X}_n^T \mathbf{W}_n \xrightarrow{D} \mathbf{X}^T \mathbf{C}$ ,  $\mathbf{W}_n^{-1} \mathbf{X}_n \xrightarrow{D} \mathbf{C}^{-1} \mathbf{X}$ , and  $\mathbf{X}_n^T \mathbf{W}_n^{-1} \xrightarrow{D} \mathbf{X}^T \mathbf{C}^{-1}$ .

**Theorem 3.14.** Let  $W_n, X_n, Y_n$ , and  $Z_n$  be sequences of random variables such that  $Y_n > 0$  and  $Z_n > 0$ . (Often  $Y_n$  and  $Z_n$  are deterministic, e.g.  $Y_n = n^{-1/2}$ .)

- a) If  $W_n = O_P(1)$  and  $X_n = O_P(1)$ , then  $W_n + X_n = O_P(1)$  and  $W_n X_n = O_P(1)$ , thus  $O_P(1) + O_P(1) = O_P(1)$  and  $O_P(1)O_P(1) = O_P(1)$ .
- b) If  $W_n = O_P(1)$  and  $X_n = o_P(1)$ , then  $W_n + X_n = O_P(1)$  and  $W_n X_n = o_P(1)$ , thus  $O_P(1) + o_P(1) = O_P(1)$  and  $O_P(1)o_P(1) = o_P(1)$ .
- c) If  $W_n = O_P(Y_n)$  and  $X_n = O_P(Z_n)$ , then  $W_n + X_n = O_P(\max(Y_n, Z_n))$  and  $W_n X_n = O_P(Y_n Z_n)$ , thus  $O_P(Y_n) + O_P(Z_n) = O_P(\max(Y_n, Z_n))$  and  $O_P(Y_n)O_P(Z_n) = O_P(Y_n Z_n)$ .



**Theorem 3.15.** i) Suppose  $\sqrt{n}(T_n - \boldsymbol{\mu}) \xrightarrow{D} N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$ . Let  $\mathbf{A}$  be a  $q \times p$  constant matrix. Then  $\mathbf{A}\sqrt{n}(T_n - \boldsymbol{\mu}) = \sqrt{n}(\mathbf{A}T_n - \mathbf{A}\boldsymbol{\mu}) \xrightarrow{D} N_q(\mathbf{A}\boldsymbol{\theta}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .

ii) Let  $\boldsymbol{\Sigma} > 0$ . If  $(T, \mathbf{C})$  is a consistent estimator of  $(\boldsymbol{\mu}, s \boldsymbol{\Sigma})$  where  $s > 0$  is some constant, then  $D_{\mathbf{x}}^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1}(\mathbf{x} - T) = s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + o_P(1)$ , so  $D_{\mathbf{x}}^2(T, \mathbf{C})$  is a consistent estimator of  $s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

iii) Let  $\boldsymbol{\Sigma} > 0$ . If  $\sqrt{n}(T - \boldsymbol{\mu}) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma})$  and if  $\mathbf{C}$  is a consistent estimator of  $\boldsymbol{\Sigma}$ , then  $n(T - \boldsymbol{\mu})^T \mathbf{C}^{-1}(T - \boldsymbol{\mu}) \xrightarrow{D} \chi_p^2$ . In particular,

$$n(\bar{\mathbf{x}} - \boldsymbol{\mu})^T \mathbf{S}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}) \xrightarrow{D} \chi_p^2.$$

**Proof:** ii)  $D_{\mathbf{x}}^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1}(\mathbf{x} - T) =$   
 $(\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T)^T [\mathbf{C}^{-1} - s^{-1} \boldsymbol{\Sigma}^{-1} + s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T)$   
 $= (\mathbf{x} - \boldsymbol{\mu})^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu}) + (\mathbf{x} - T)^T [\mathbf{C}^{-1} - s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - T)$   
 $+ (\mathbf{x} - \boldsymbol{\mu})^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\boldsymbol{\mu} - T) + (\boldsymbol{\mu} - T)^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu})$   
 $+ (\boldsymbol{\mu} - T)^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\boldsymbol{\mu} - T) = s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + O_P(1).$

(Note that  $D_{\mathbf{x}}^2(T, \mathbf{C}) = s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + O_P(n^{-\delta})$  if  $(T, \mathbf{C})$  is a consistent estimator of  $(\boldsymbol{\mu}, s \boldsymbol{\Sigma})$  with rate  $n^\delta$  where  $0 < \delta \leq 0.5$  if  $[\mathbf{C}^{-1} - s^{-1} \boldsymbol{\Sigma}^{-1}] = O_P(n^{-\delta})$ .)

Alternatively,  $D_{\mathbf{x}}^2(T, \mathbf{C})$  is a continuous function of  $(T, \mathbf{C})$  if  $\mathbf{C} > 0$  for  $n > 10p$ . Hence  $D_{\mathbf{x}}^2(T, \mathbf{C}) \xrightarrow{P} D_{\mathbf{x}}^2(\boldsymbol{\mu}, s \boldsymbol{\Sigma})$ .

iii) Note that  $\mathbf{Z}_n = \sqrt{n} \boldsymbol{\Sigma}^{-1/2}(T - \boldsymbol{\mu}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{I}_p)$ . Thus  $\mathbf{Z}_n^T \mathbf{Z}_n = n(T - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(T - \boldsymbol{\mu}) \xrightarrow{D} \chi_p^2$ . Now  $n(T - \boldsymbol{\mu})^T \mathbf{C}^{-1}(T - \boldsymbol{\mu}) = n(T - \boldsymbol{\mu})^T [\mathbf{C}^{-1} - \boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-1}] (T - \boldsymbol{\mu}) = n(T - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(T - \boldsymbol{\mu}) + n(T - \boldsymbol{\mu})^T [\mathbf{C}^{-1} - \boldsymbol{\Sigma}^{-1}] (T - \boldsymbol{\mu}) = n(T - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(T - \boldsymbol{\mu}) + o_P(1) \xrightarrow{D} \chi_p^2$  since  $\sqrt{n}(T - \boldsymbol{\mu})^T [\mathbf{C}^{-1} - \boldsymbol{\Sigma}^{-1}] \sqrt{n}(T - \boldsymbol{\mu}) = O_P(1) o_P(1) O_P(1) = o_P(1)$ .  $\square$

**Example 3.2.** Suppose that  $\mathbf{x}_n \perp \mathbf{y}_n$  for  $n = 1, 2, \dots$ . Suppose  $\mathbf{x}_n \xrightarrow{D} \mathbf{x}$ , and  $\mathbf{y}_n \xrightarrow{D} \mathbf{y}$  where  $\mathbf{x} \perp \mathbf{y}$ . Then

$$\begin{bmatrix} \mathbf{x}_n \\ \mathbf{y}_n \end{bmatrix} \xrightarrow{D} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$$

by Theorem 3.7. To see this, let  $\mathbf{t} = (t_1^T, t_2^T)^T$ ,  $\mathbf{z}_n = (\mathbf{x}_n^T, \mathbf{y}_n^T)^T$ , and  $\mathbf{z} = (\mathbf{x}^T, \mathbf{y}^T)^T$ . Since  $\mathbf{x}_n \perp \mathbf{y}_n$  and  $\mathbf{x} \perp \mathbf{y}$ , the characteristic function

$$\phi_{\mathbf{z}_n}(\mathbf{t}) = \phi_{\mathbf{x}_n}(t_1) \phi_{\mathbf{y}_n}(t_2) \rightarrow \phi_{\mathbf{x}}(t_1) \phi_{\mathbf{y}}(t_2) = \phi_{\mathbf{z}}(\mathbf{t}).$$

Hence  $\mathbf{g}(\mathbf{z}_n) \xrightarrow{D} \mathbf{g}(\mathbf{z})$  by Theorem 3.12.

**Remark 3.1.** In the above example, we can show  $\mathbf{x} \perp \mathbf{y}$  instead of assuming  $\mathbf{x} \perp \mathbf{y}$ . See Ferguson (1996, p. 42).

### 3.3 Summary

### 3.4 Complements

Suppose  $\boldsymbol{\theta} = \mathbf{g}^{-1}(\boldsymbol{\eta})$ . In analysis, the fact that

$$\mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})}^{-1} = \mathbf{D}\mathbf{g}^{-1}(\boldsymbol{\eta})$$

is a corollary of the inverse mapping theorem (or of the inverse function theorem). See Apostol (1957, p. 146), Bickel and Doksum (2007, p. 517), Marsden and Hoffman (1993, p. 393 ) and Wade (2000, p. 353).

### 3.5 Problems

**3.1.** Many multiple linear regression estimators  $\hat{\boldsymbol{\beta}}$  satisfy

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(0, V(\hat{\boldsymbol{\beta}}, F) \mathbf{W}) \quad (3.5)$$

when

$$\frac{\mathbf{X}^T \mathbf{X}}{n} \xrightarrow{P} \mathbf{W}^{-1}, \quad (3.6)$$

and when the errors  $e_i$  are iid with a cdf  $F$  and a unimodal pdf  $f$  that is symmetric with a unique maximum at 0. When the variance  $V(e_i)$  exists,

$$V(OLS, F) = V(e_i) = \sigma^2 \quad \text{while} \quad V(L_1, F) = \frac{1}{4[f(0)]^2}.$$

In the multiple linear regression model,

$$Y_i = x_{i,1}\beta_1 + x_{i,2}\beta_2 + \cdots + x_{i,p}\beta_p + e_i = \mathbf{x}_i^T \boldsymbol{\beta} + e_i \quad (3.7)$$

for  $i = 1, \dots, n$ . In matrix notation, these  $n$  equations become

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad (3.8)$$

where  $\mathbf{Y}$  is an  $n \times 1$  vector of dependent variables,  $\mathbf{X}$  is an  $n \times p$  matrix of predictors,  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of unknown coefficients, and  $\mathbf{e}$  is an  $n \times 1$  vector of unknown errors.

a) What is the  $ij$ th element of the matrix

$$\frac{\mathbf{X}^T \mathbf{X}}{n}?$$

b) Suppose  $x_{k,1} = 1$  and that  $x_{k,j} \sim X_j$  are iid with  $E(X_j) = 0$  and  $V(X_j) = 1$  for  $k = 1, \dots, n$  and  $j = 2, \dots, p$ . Assume that  $X_i$  and  $X_j$  are independent for  $i \neq j$ ,  $i > 1$  and  $j > 1$ . (Often  $x_{k,j} \sim N(0, 1)$  in simulations.) Then what is  $\mathbf{W}^{-1}$  for model (3.7)?

c) Suppose  $p = 2$  and  $Y_i = \alpha + \beta X_i + e_i$ . Show

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \frac{\sum X_i^2}{n \sum (X_i - \bar{X})^2} & \frac{-\sum X_i}{n \sum (X_i - \bar{X})^2} \\ \frac{-\sum X_i}{n \sum (X_i - \bar{X})^2} & \frac{n}{n \sum (X_i - \bar{X})^2} \end{bmatrix}.$$

d) Under the conditions of c), let  $S_x^2 = \sum (X_i - \bar{X})^2 / n$ . Show that

$$n(\mathbf{X}^T \mathbf{X})^{-1} = \left( \frac{\mathbf{X}^T \mathbf{X}}{n} \right)^{-1} = \begin{bmatrix} \frac{\frac{1}{n} \sum X_i^2}{S_x^2} & \frac{-\bar{X}}{S_x^2} \\ \frac{-\bar{X}}{S_x^2} & \frac{1}{S_x^2} \end{bmatrix}.$$

e) If the  $X_i$  are iid with variance  $V(X)$  then  $n(\mathbf{X}^T \mathbf{X})^{-1} \xrightarrow{P} \mathbf{W}$ . What is  $\mathbf{W}$ ?

f) Now suppose that  $n$  is divisible by 5 and the  $n/5$  of  $X_i$  are at 0.1,  $n/5$  at 0.3,  $n/5$  at 0.5,  $n/5$  at 0.7 and  $n/5$  at 0.9. (Hence if  $n = 100$ , 20 of the  $X_i$  are at 0.1, 0.3, 0.5, 0.7 and 0.9.)

Find  $\sum X_i^2 / n$ ,  $\bar{X}$  and  $S_x^2$ . (Your answers should not depend on  $n$ .)

g) Under the conditions of f), estimate  $V(\hat{\alpha})$  and  $V(\hat{\beta})$  if  $L_1$  is used and if the  $e_i$  are iid  $N(0, 0.01)$ .

Hint: Estimate  $\mathbf{W}$  with  $n(\mathbf{X}^T \mathbf{X})^{-1}$  and  $V(\hat{\beta}, F) = V(L_1, F) = \frac{1}{4[f'(0)]^2}$ . Hence

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} \approx N_2 \left[ \begin{pmatrix} \alpha \\ \beta \end{pmatrix}, \frac{1}{n} \frac{1}{4[f'(0)]^2} \begin{pmatrix} \frac{\frac{1}{n} \sum X_i^2}{S_x^2} & \frac{-\bar{X}}{S_x^2} \\ \frac{-\bar{X}}{S_x^2} & \frac{1}{S_x^2} \end{pmatrix} \right].$$

You should get an answer like  $0.0648/n$ .