

## Chapter 3

# Elliptically Contoured Distributions

The multivariate location and dispersion model of Definition 2.1 is in many ways similar to the multiple linear regression model. The data are iid vectors from some distribution such as the multivariate normal (MVN) distribution. The location parameter  $\boldsymbol{\mu}$  of interest may be the mean or the center of symmetry of an elliptically contoured distribution. Hyperellipsoids will be estimated instead of hyperplanes, and Mahalanobis distances will be used instead of absolute residuals to determine if an observation is a potential outlier. Review Section 2.1 for important notation.

Although usually random vectors in this text are denoted by  $\boldsymbol{x}$ ,  $\boldsymbol{y}$  or  $\boldsymbol{z}$ , this chapter will usually use the notation  $\boldsymbol{X} = (X_1, \dots, X_p)^T$  and  $\boldsymbol{Y}$  for the random vectors, and  $\boldsymbol{x} = (x_1, \dots, x_p)^T$  for the observed value of the random vector. This notation will be useful to avoid confusion when studying conditional distributions such as  $\boldsymbol{Y}|\boldsymbol{X} = \boldsymbol{x}$ .

### 3.1 The Multivariate Normal Distribution

**Definition 3.1:** Rao (1965, p. 437). A  $p \times 1$  random vector  $\boldsymbol{X}$  has a  $p$ -dimensional *multivariate normal distribution*  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  iff  $\boldsymbol{t}^T \boldsymbol{X}$  has a univariate normal distribution for any  $p \times 1$  vector  $\boldsymbol{t}$ .

If  $\boldsymbol{\Sigma}$  is positive definite, then  $\boldsymbol{X}$  has a pdf

$$f(\boldsymbol{z}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-(1/2)(\boldsymbol{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{z}-\boldsymbol{\mu})} \quad (3.1)$$

where  $|\boldsymbol{\Sigma}|^{1/2}$  is the square root of the determinant of  $\boldsymbol{\Sigma}$ . Note that if  $p = 1$ , then the quadratic form in the exponent is  $(z - \boldsymbol{\mu})(\sigma^2)^{-1}(z - \boldsymbol{\mu})$  and  $X$  has the univariate  $N(\boldsymbol{\mu}, \sigma^2)$  pdf. If  $\boldsymbol{\Sigma}$  is positive semidefinite but not positive definite, then  $\boldsymbol{x}$  has a degenerate distribution. For example, the univariate  $N(0, 0^2)$  distribution is degenerate (the point mass at 0).

**Definition 3.2.** The *population mean* of a random  $p \times 1$  vector  $\mathbf{X} = (X_1, \dots, X_p)^T$  is

$$E(\mathbf{X}) = (E(X_1), \dots, E(X_p))^T$$

and the  $p \times p$  *population covariance matrix*

$$\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}_{\mathbf{x}} = E(\mathbf{X} - E(\mathbf{X}))(\mathbf{X} - E(\mathbf{X}))^T = ((\sigma_{i,j})).$$

That is, the  $ij$  entry of  $\text{Cov}(\mathbf{X})$  is  $\text{Cov}(X_i, X_j) = \sigma_{i,j}$ .

The covariance matrix is also called the variance–covariance matrix and variance matrix. Sometimes the notation  $\text{Var}(\mathbf{X})$  is used. Note that  $\text{Cov}(\mathbf{X})$  is a symmetric positive semidefinite matrix. If  $\mathbf{X}$  and  $\mathbf{Y}$  are  $p \times 1$  random vectors,  $\mathbf{a}$  a conformable constant vector and  $\mathbf{A}$  and  $\mathbf{B}$  are conformable constant matrices, then

$$E(\mathbf{a} + \mathbf{X}) = \mathbf{a} + E(\mathbf{X}) \quad \text{and} \quad E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}) \quad (3.2)$$

and

$$E(\mathbf{A}\mathbf{X}) = \mathbf{A}E(\mathbf{X}) \quad \text{and} \quad E(\mathbf{A}\mathbf{X}\mathbf{B}) = \mathbf{A}E(\mathbf{X})\mathbf{B}. \quad (3.3)$$

Thus

$$\text{Cov}(\mathbf{a} + \mathbf{A}\mathbf{X}) = \text{Cov}(\mathbf{A}\mathbf{X}) = \mathbf{A}\text{Cov}(\mathbf{X})\mathbf{A}^T. \quad (3.4)$$

Some important properties of MVN distributions are given in the following three propositions. These propositions can be proved using results from Johnson and Wichern (1988, p. 127-132).

**Proposition 3.1.** a) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $E(\mathbf{X}) = \boldsymbol{\mu}$  and

$$\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}.$$

b) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then any linear combination  $\mathbf{t}^T \mathbf{X} = t_1 X_1 + \dots + t_p X_p \sim N_1(\mathbf{t}^T \boldsymbol{\mu}, \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t})$ . Conversely, if  $\mathbf{t}^T \mathbf{X} \sim N_1(\mathbf{t}^T \boldsymbol{\mu}, \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t})$  for every  $p \times 1$  vector  $\mathbf{t}$ , then  $\boldsymbol{x} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

c) **The joint distribution of independent normal random variables is MVN.** If  $X_1, \dots, X_p$  are independent univariate normal  $N(\mu_i, \sigma_i^2)$  random vectors, then  $\mathbf{X} = (X_1, \dots, X_p)^T$  is  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)$  and  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$  (so the off diagonal entries  $\sigma_{i,j} = 0$  while the diagonal entries of  $\boldsymbol{\Sigma}$  are  $\sigma_{i,i} = \sigma_i^2$ ).

d) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and if  $\mathbf{A}$  is a  $q \times p$  matrix, then  $\mathbf{AX} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ . If  $\mathbf{a}$  is a  $p \times 1$  vector of constants, then  $\mathbf{a} + \mathbf{X} \sim N_p(\mathbf{a} + \boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

It will be useful to partition  $\mathbf{X}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\Sigma}$ . Let  $\mathbf{X}_1$  and  $\boldsymbol{\mu}_1$  be  $q \times 1$  vectors, let  $\mathbf{X}_2$  and  $\boldsymbol{\mu}_2$  be  $(p - q) \times 1$  vectors, let  $\boldsymbol{\Sigma}_{11}$  be a  $q \times q$  matrix, let  $\boldsymbol{\Sigma}_{12}$  be a  $q \times (p - q)$  matrix, let  $\boldsymbol{\Sigma}_{21}$  be a  $(p - q) \times q$  matrix, and let  $\boldsymbol{\Sigma}_{22}$  be a  $(p - q) \times (p - q)$  matrix. Then

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \text{and} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

**Proposition 3.2.** a) **All subsets of a MVN are MVN:**  $(X_{k_1}, \dots, X_{k_q})^T \sim N_q(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$  where  $\tilde{\boldsymbol{\mu}}_i = E(X_{k_i})$  and  $\tilde{\boldsymbol{\Sigma}}_{ij} = \text{Cov}(X_{k_i}, X_{k_j})$ . In particular,  $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$  and  $\mathbf{X}_2 \sim N_{p-q}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$ .

b) If  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are independent, then  $\text{Cov}(\mathbf{X}_1, \mathbf{X}_2) = \boldsymbol{\Sigma}_{12} = E[(\mathbf{X}_1 - E(\mathbf{X}_1))(\mathbf{X}_2 - E(\mathbf{X}_2))^T] = \mathbf{0}$ , a  $q \times (p - q)$  matrix of zeroes.

c) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are independent iff  $\boldsymbol{\Sigma}_{12} = \mathbf{0}$ .

d) If  $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$  and  $\mathbf{X}_2 \sim N_{p-q}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$  are independent, then

$$\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} \sim N_p \left( \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \right).$$

**Proposition 3.3.** **The conditional distribution of a MVN is MVN.** If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the conditional distribution of  $\mathbf{X}_1$  given that  $\mathbf{X}_2 = \mathbf{x}_2$  is multivariate normal with mean  $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$  and covariance matrix  $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$ . That is,

$$\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2 \sim N_q(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

**Example 3.1.** Let  $p = 2$  and let  $(Y, X)^T$  have a bivariate normal distribution. That is,

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_Y \\ \mu_X \end{pmatrix}, \begin{pmatrix} \sigma_Y^2 & \text{Cov}(Y, X) \\ \text{Cov}(X, Y) & \sigma_X^2 \end{pmatrix} \right).$$

Also recall that the population correlation between  $X$  and  $Y$  is given by

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{VAR}(X)}\sqrt{\text{VAR}(Y)}} = \frac{\sigma_{X,Y}}{\sigma_X\sigma_Y}$$

if  $\sigma_X > 0$  and  $\sigma_Y > 0$ . Then  $Y|X = x \sim N(E(Y|X = x), \text{VAR}(Y|X = x))$  where the conditional mean

$$E(Y|X = x) = \mu_Y + \text{Cov}(Y, X)\frac{1}{\sigma_X^2}(x - \mu_X) = \mu_Y + \rho(X, Y)\sqrt{\frac{\sigma_Y^2}{\sigma_X^2}}(x - \mu_X)$$

and the conditional variance

$$\begin{aligned} \text{VAR}(Y|X = x) &= \sigma_Y^2 - \text{Cov}(X, Y)\frac{1}{\sigma_X^2}\text{Cov}(X, Y) \\ &= \sigma_Y^2 - \rho(X, Y)\sqrt{\frac{\sigma_Y^2}{\sigma_X^2}}\rho(X, Y)\sqrt{\sigma_X^2}\sqrt{\sigma_Y^2} \\ &= \sigma_Y^2 - \rho^2(X, Y)\sigma_Y^2 = \sigma_Y^2[1 - \rho^2(X, Y)]. \end{aligned}$$

Also  $aX + bY$  is univariate normal with mean  $a\mu_X + b\mu_Y$  and variance

$$a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab \text{Cov}(X, Y).$$

**Remark 3.1.** There are several common misconceptions. First, **it is not true that every linear combination  $t^T \mathbf{X}$  of normal random variables is a normal random variable**, and **it is not true that all uncorrelated normal random variables are independent**. The key condition in Proposition 3.1b and Proposition 3.2c is that the joint distribution of  $\mathbf{X}$  is MVN. It is possible that  $X_1, X_2, \dots, X_p$  each has a marginal distribution that is univariate normal, but the joint distribution of  $\mathbf{X}$  is not MVN. See Seber and Lee (2003, p. 23), Kowalski (1973) and examine the following example from Rohatgi (1976, p. 229). Suppose that the joint pdf of  $X$  and  $Y$  is a mixture of two bivariate normal distributions both with  $EX = EY = 0$  and  $\text{VAR}(X) = \text{VAR}(Y) = 1$ , but  $\text{Cov}(X, Y) = \pm\rho$ . Hence  $f(x, y) =$

$$\begin{aligned} &\frac{1}{2} \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(\frac{-1}{2(1-\rho^2)}(x^2 - 2\rho xy + y^2)\right) + \\ &\frac{1}{2} \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(\frac{-1}{2(1-\rho^2)}(x^2 + 2\rho xy + y^2)\right) \equiv \frac{1}{2}f_1(x, y) + \frac{1}{2}f_2(x, y) \end{aligned}$$

where  $x$  and  $y$  are real and  $0 < \rho < 1$ . Since both marginal distributions of  $f_i(x, y)$  are  $N(0,1)$  for  $i = 1$  and  $2$  by Proposition 3.2 a), the marginal distributions of  $X$  and  $Y$  are  $N(0,1)$ . Since  $\int \int xy f_i(x, y) dx dy = \rho$  for  $i = 1$  and  $-\rho$  for  $i = 2$ ,  $X$  and  $Y$  are uncorrelated, but  $X$  and  $Y$  are not independent since  $f(x, y) \neq f_X(x)f_Y(y)$ .

**Remark 3.2.** In Proposition 3.3, suppose that  $\mathbf{X} = (Y, X_2, \dots, X_p)^T$ . Let  $X_1 = Y$  and  $\mathbf{X}_2 = (X_2, \dots, X_p)^T$ . Then  $E[Y|\mathbf{X}_2] = \beta_1 + \beta_2 X_2 + \dots + \beta_p X_p$  and  $\text{VAR}[Y|\mathbf{X}_2]$  is a constant that does not depend on  $\mathbf{X}_2$ . Hence  $Y = \beta_1 + \beta_2 X_2 + \dots + \beta_p X_p + e$  follows the multiple linear regression model.

## 3.2 Elliptically Contoured Distributions

**Definition 3.3: Johnson (1987, p. 107-108).** A  $p \times 1$  random vector  $\mathbf{X}$  has an *elliptically contoured distribution*, also called an *elliptically symmetric distribution*, if  $\mathbf{X}$  has joint pdf

$$f(\mathbf{z}) = k_p |\Sigma|^{-1/2} g[(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu})], \quad (3.5)$$

and we say  $\mathbf{X}$  has an elliptically contoured  $EC_p(\boldsymbol{\mu}, \Sigma, g)$  distribution.

If  $\mathbf{X}$  has an elliptically contoured (EC) distribution, then the characteristic function of  $\mathbf{X}$  is

$$\phi_{\mathbf{X}}(\mathbf{t}) = \exp(i\mathbf{t}^T \boldsymbol{\mu}) \psi(\mathbf{t}^T \Sigma \mathbf{t}) \quad (3.6)$$

for some function  $\psi$ . If the second moments exist, then

$$E(\mathbf{X}) = \boldsymbol{\mu} \quad (3.7)$$

and

$$\text{Cov}(\mathbf{X}) = c_X \Sigma \quad (3.8)$$

where

$$c_X = -2\psi'(0).$$

**Definition 3.4.** The *population squared Mahalanobis distance*

$$U \equiv D^2 = D^2(\boldsymbol{\mu}, \Sigma) = (\mathbf{X} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{X} - \boldsymbol{\mu}). \quad (3.9)$$

For elliptically contoured distributions,  $U$  has pdf

$$h(u) = \frac{\pi^{p/2}}{\Gamma(p/2)} k_p u^{p/2-1} g(u). \quad (3.10)$$

For  $c > 0$ , an  $EC_p(\boldsymbol{\mu}, c\mathbf{I}, g)$  distribution is *spherical about  $\boldsymbol{\mu}$*  where  $\mathbf{I}$  is the  $p \times p$  identity matrix. The *multivariate normal distribution*  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  has  $k_p = (2\pi)^{-p/2}$ ,  $\psi(u) = g(u) = \exp(-u/2)$  and  $h(u)$  is the  $\chi_p^2$  pdf.

The following lemma is useful for proving properties of EC distributions without using the characteristic function (10.6). See Eaton (1986) and Cook (1998, p. 57, 130).

**Lemma 3.4.** Let  $\mathbf{X}$  be a  $p \times 1$  random vector with 1st moments; ie,  $E(\mathbf{X})$  exists. Let  $\mathbf{B}$  be any constant full rank  $p \times r$  matrix where  $1 \leq r \leq p$ . Then  $\mathbf{X}$  is elliptically contoured iff for all such conforming matrices  $\mathbf{B}$ ,

$$E(\mathbf{X} | \mathbf{B}^T \mathbf{X}) = \boldsymbol{\mu} + \mathbf{M}_B \mathbf{B}^T (\mathbf{X} - \boldsymbol{\mu}) = \mathbf{a}_B + \mathbf{M}_B \mathbf{B}^T \mathbf{X} \quad (3.11)$$

where the  $p \times 1$  constant vector  $\mathbf{a}_B$  and the  $p \times r$  constant matrix  $\mathbf{M}_B$  both depend on  $\mathbf{B}$ .

A useful fact is that  $\mathbf{a}_B$  and  $\mathbf{M}_B$  do not depend on  $g$ :

$$\mathbf{a}_B = \boldsymbol{\mu} - \mathbf{M}_B \mathbf{B}^T \boldsymbol{\mu} = (\mathbf{I}_p - \mathbf{M}_B \mathbf{B}^T) \boldsymbol{\mu},$$

and

$$\mathbf{M}_B = \boldsymbol{\Sigma} \mathbf{B} (\mathbf{B}^T \boldsymbol{\Sigma} \mathbf{B})^{-1}.$$

See Problem 3.11. Notice that in the formula for  $\mathbf{M}_B$ ,  $\boldsymbol{\Sigma}$  can be replaced by  $c\boldsymbol{\Sigma}$  where  $c > 0$  is a constant. In particular, if the EC distribution has 2nd moments,  $\text{Cov}(\mathbf{X})$  can be used instead of  $\boldsymbol{\Sigma}$ .

To use Lemma 3.4 to prove interesting properties, partition  $\mathbf{X}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\Sigma}$ . Let  $\mathbf{X}_1$  and  $\boldsymbol{\mu}_1$  be  $q \times 1$  vectors, let  $\mathbf{X}_2$  and  $\boldsymbol{\mu}_2$  be  $(p-q) \times 1$  vectors. Let  $\boldsymbol{\Sigma}_{11}$  be a  $q \times q$  matrix, let  $\boldsymbol{\Sigma}_{12}$  be a  $q \times (p-q)$  matrix, let  $\boldsymbol{\Sigma}_{21}$  be a  $(p-q) \times q$  matrix, and let  $\boldsymbol{\Sigma}_{22}$  be a  $(p-q) \times (p-q)$  matrix. Then

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \text{and} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

Also assume that the  $(p+1) \times 1$  vector  $(Y, \mathbf{X}^T)^T$  is  $EC_{p+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$  where  $Y$  is a random variable,  $\mathbf{X}$  is a  $p \times 1$  vector, and use

$$\begin{pmatrix} Y \\ \mathbf{X} \end{pmatrix}, \boldsymbol{\mu} = \begin{pmatrix} \mu_Y \\ \boldsymbol{\mu}_X \end{pmatrix}, \text{ and } \boldsymbol{\Sigma} = \begin{pmatrix} \Sigma_{YY} & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_{XX} \end{pmatrix}.$$

**Proposition 3.5.** Let  $\mathbf{X} \sim EC_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$  and assume that  $E(\mathbf{X})$  exists.

- a) Any subset of  $\mathbf{X}$  is EC, in particular  $\mathbf{X}_1$  is EC.
- b) (Cook 1998 p. 131, Kelker 1970). If  $\text{Cov}(\mathbf{X})$  is nonsingular,

$$\text{Cov}(\mathbf{X} | \mathbf{B}^T \mathbf{X}) = d_g(\mathbf{B}^T \mathbf{X}) [\boldsymbol{\Sigma} - \boldsymbol{\Sigma} \mathbf{B} (\mathbf{B}^T \boldsymbol{\Sigma} \mathbf{B})^{-1} \mathbf{B}^T \boldsymbol{\Sigma}]$$

where the real valued function  $d_g(\mathbf{B}^T \mathbf{X})$  is constant iff  $\mathbf{X}$  is MVN.

**Proof** of a). Let  $\mathbf{A}$  be an arbitrary full rank  $q \times r$  matrix where  $1 \leq r \leq q$ . Let

$$\mathbf{B} = \begin{pmatrix} \mathbf{A} \\ \mathbf{0} \end{pmatrix}.$$

Then  $\mathbf{B}^T \mathbf{X} = \mathbf{A}^T \mathbf{X}_1$ , and

$$E[\mathbf{X} | \mathbf{B}^T \mathbf{X}] = E\left[\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} | \mathbf{A}^T \mathbf{X}_1\right] =$$

$$\begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} + \begin{pmatrix} \mathbf{M}_{1B} \\ \mathbf{M}_{2B} \end{pmatrix} \begin{pmatrix} \mathbf{A}^T & \mathbf{0}^T \end{pmatrix} \begin{pmatrix} \mathbf{X}_1 - \boldsymbol{\mu}_1 \\ \mathbf{X}_2 - \boldsymbol{\mu}_2 \end{pmatrix}$$

by Lemma 3.4. Hence  $E[\mathbf{X}_1 | \mathbf{A}^T \mathbf{X}_1] = \boldsymbol{\mu}_1 + \mathbf{M}_{1B} \mathbf{A}^T (\mathbf{X}_1 - \boldsymbol{\mu}_1)$ . Since  $\mathbf{A}$  was arbitrary,  $\mathbf{X}_1$  is EC by Lemma 3.4. Notice that  $\mathbf{M}_B = \boldsymbol{\Sigma} \mathbf{B} (\mathbf{B}^T \boldsymbol{\Sigma} \mathbf{B})^{-1} =$

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \mathbf{A} \\ \mathbf{0} \end{pmatrix} \left[ \begin{pmatrix} \mathbf{A}^T & \mathbf{0}^T \end{pmatrix} \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \mathbf{A} \\ \mathbf{0} \end{pmatrix} \right]^{-1} \\ = \begin{pmatrix} \mathbf{M}_{1B} \\ \mathbf{M}_{2B} \end{pmatrix}.$$

Hence

$$\mathbf{M}_{1B} = \Sigma_{11} \mathbf{A} (\mathbf{A}^T \Sigma_{11} \mathbf{A})^{-1}$$

and  $\mathbf{X}_1$  is EC with location and dispersion parameters  $\boldsymbol{\mu}_1$  and  $\Sigma_{11}$ . QED

**Proposition 3.6.** Let  $(Y, \mathbf{X}^T)^T$  be  $EC_{p+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$  where  $Y$  is a random variable.

a) Assume that  $E[(Y, \mathbf{X}^T)^T]$  exists. Then  $E(Y|\mathbf{X}) = \alpha + \boldsymbol{\beta}^T \mathbf{X}$  where  $\alpha = \mu_Y - \boldsymbol{\beta}^T \boldsymbol{\mu}_X$  and

$$\boldsymbol{\beta} = \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY}.$$

b) Even if the first moment does not exist, the conditional median

$$\text{MED}(Y|\mathbf{X}) = \alpha + \boldsymbol{\beta}^T \mathbf{X}$$

where  $\alpha$  and  $\boldsymbol{\beta}$  are given in a).

**Proof.** a) The trick is to choose  $\mathbf{B}$  so that Lemma 3.4 applies. Let

$$\mathbf{B} = \begin{pmatrix} \mathbf{0}^T \\ \mathbf{I}_p \end{pmatrix}.$$

Then  $\mathbf{B}^T \boldsymbol{\Sigma} \mathbf{B} = \boldsymbol{\Sigma}_{XX}$  and

$$\boldsymbol{\Sigma} \mathbf{B} = \begin{pmatrix} \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XX} \end{pmatrix}.$$

Now

$$\begin{aligned} E\left[\begin{pmatrix} Y \\ \mathbf{X} \end{pmatrix} \mid \mathbf{X}\right] &= E\left[\begin{pmatrix} Y \\ \mathbf{X} \end{pmatrix} \mid \mathbf{B}^T \begin{pmatrix} Y \\ \mathbf{X} \end{pmatrix}\right] \\ &= \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{B} (\mathbf{B}^T \boldsymbol{\Sigma} \mathbf{B})^{-1} \mathbf{B}^T \begin{pmatrix} Y - \mu_Y \\ \mathbf{X} - \boldsymbol{\mu}_X \end{pmatrix} \end{aligned}$$

by Lemma 3.4. The right hand side of the last equation is equal to

$$\boldsymbol{\mu} + \begin{pmatrix} \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XX} \end{pmatrix} \boldsymbol{\Sigma}_{XX}^{-1} (\mathbf{X} - \boldsymbol{\mu}_X) = \begin{pmatrix} \mu_Y - \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\mu}_X + \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \mathbf{X} \\ \mathbf{X} \end{pmatrix}$$

and the result follows since

$$\boldsymbol{\beta}^T = \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1}.$$

b) See Croux, Dehon, Rousseeuw and Van Aelst (2001) for references.



**Example 3.2.** This example illustrates another application of Lemma 3.4. Suppose that  $\mathbf{X}$  comes from a mixture of two multivariate normals with the same mean and proportional covariance matrices. That is, let

$$\mathbf{X} \sim (1 - \gamma)N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + \gamma N_p(\boldsymbol{\mu}, c\boldsymbol{\Sigma})$$

where  $c > 0$  and  $0 < \gamma < 1$ . Since the multivariate normal distribution is elliptically contoured (and see Proposition 1.2c),

$$\begin{aligned} E(\mathbf{X}|\mathbf{B}^T \mathbf{X}) &= (1 - \gamma)[\boldsymbol{\mu} + \mathbf{M}_1 \mathbf{B}^T (\mathbf{X} - \boldsymbol{\mu})] + \gamma[\boldsymbol{\mu} + \mathbf{M}_2 \mathbf{B}^T (\mathbf{X} - \boldsymbol{\mu})] \\ &= \boldsymbol{\mu} + [(1 - \gamma)\mathbf{M}_1 + \gamma\mathbf{M}_2] \mathbf{B}^T (\mathbf{X} - \boldsymbol{\mu}) \equiv \boldsymbol{\mu} + \mathbf{M} \mathbf{B}^T (\mathbf{X} - \boldsymbol{\mu}). \end{aligned}$$

Since  $\mathbf{M}_B$  only depends on  $\mathbf{B}$  and  $\boldsymbol{\Sigma}$ , it follows that  $\mathbf{M}_1 = \mathbf{M}_2 = \mathbf{M} = \mathbf{M}_B$ . Hence  $\mathbf{X}$  has an elliptically contoured distribution by Lemma 3.4.

Let  $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $y \sim \chi_d^2$  be independent. Let  $w_i = x_i/(y/d)^{1/2}$  for  $i = 1, \dots, p$ . Then  $\mathbf{w}$  has a multivariate t-distribution with parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  and degrees of freedom  $d$ , an important elliptically contoured distribution. Cornish (1954) shows that the covariance matrix of  $\mathbf{w}$  is  $\text{Cov}(\mathbf{w}) = \frac{d}{d-2} \boldsymbol{\Sigma}$  for  $d > 2$ . The case  $d = 1$  is known as a multivariate Cauchy distribution. The joint pdf of  $\mathbf{w}$  is

$$f(\mathbf{z}) = \frac{\Gamma((d+p)/2) |\boldsymbol{\Sigma}|^{-1/2}}{(\pi d)^{p/2} \Gamma(d/2)} [1 + d^{-1}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})]^{-(d+p)/2}.$$

See Mardia, Kent and Bibby (1979, p. 43, 57). See Johnson and Kotz (1972, p. 134) for the special case where the  $x_i \sim N(0, 1)$ .

If  $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $u_i = \exp(x_i)$  for  $i = 1, \dots, p$ , then  $\mathbf{u}$  has a multivariate lognormal distribution with parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . This distribution is not an elliptically contoured distribution.

### 3.3 Sample Mahalanobis Distances

In the multivariate location and dispersion model, sample Mahalanobis distances play a role similar to that of residuals in multiple linear regression. The observed data  $\mathbf{X}_i = \mathbf{x}_i$  for  $i = 1, \dots, n$  is collected in an  $n \times p$  matrix  $\mathbf{W}$  with  $n$  rows  $\mathbf{x}_1^T, \dots, \mathbf{x}_n^T$ . Let the  $p \times 1$  column vector  $T(\mathbf{W})$  be a multivariate location estimator, and let the  $p \times p$  symmetric positive definite matrix  $\mathbf{C}(\mathbf{W})$  be a dispersion estimator.

**Definition 3.5.** The *ith squared Mahalanobis distance* is

$$D_i^2 = D_i^2(T(\mathbf{W}), \mathbf{C}(\mathbf{W})) = (\mathbf{X}_i - T(\mathbf{W}))^T \mathbf{C}^{-1}(\mathbf{W})(\mathbf{X}_i - T(\mathbf{W})) \quad (3.12)$$

for each point  $\mathbf{X}_i$ . Notice that  $D_i^2$  is a random variable (scalar valued).

Notice that the population squared Mahalanobis distance is

$$D_{\mathbf{X}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) \quad (3.13)$$

and that the term  $\boldsymbol{\Sigma}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$  is the  $p$ -dimensional analog to the  $z$ -score used to transform a univariate  $N(\mu, \sigma^2)$  random variable into a  $N(0, 1)$  random variable. Hence the sample Mahalanobis distance  $D_i = \sqrt{D_i^2}$  is an analog of the absolute value  $|Z_i|$  of the sample  $Z$ -score  $Z_i = (X_i - \bar{X})/\hat{\sigma}$ . Also notice that the Euclidean distance of  $\mathbf{x}_i$  from the estimate of center  $T(\mathbf{W})$  is  $D_i(T(\mathbf{W}), \mathbf{I}_p)$  where  $\mathbf{I}_p$  is the  $p \times p$  identity matrix.

**Example 3.3.** The contours of constant density for the  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution are ellipsoids defined by  $\mathbf{x}$  such that  $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = a^2$ . An  $\alpha$ -density region  $R_\alpha$  is a set such that  $P(\mathbf{X} \in R_\alpha) = \alpha$ , and for the  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution, the regions of highest density are sets of the form

$$\{\mathbf{x} : (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \leq \chi_p^2(\alpha)\} = \{\mathbf{x} : D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \leq \chi_p^2(\alpha)\}$$

where  $P(W \leq \chi_p^2(\alpha)) = \alpha$  if  $W \sim \chi_p^2$ . If the  $\mathbf{X}_i$  are  $n$  iid random vectors each with a  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  pdf, then a scatterplot of  $X_{i,k}$  versus  $X_{i,j}$  should be ellipsoidal for  $k \neq j$ . Similar statements hold if  $\mathbf{X}$  is  $EC_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$ , but the  $\alpha$ -density region will use a constant  $U_\alpha$  obtained from Equation (3.10).

The classical Mahalanobis distance corresponds to the sample mean and sample covariance matrix

$$T(\mathbf{W}) = \bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i,$$

and

$$\mathbf{C}(\mathbf{W}) = \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^T$$

and will be denoted by  $MD_i$ . When  $T(\mathbf{W})$  and  $\mathbf{C}(\mathbf{W})$  are estimators other than the sample mean and covariance,  $D_i = \sqrt{D_i^2}$  will sometimes be denoted by  $RD_i$ .

## 3.4 Large Sample Theory

The first three subsections will review large sample theory for the univariate case, then multivariate theory will be given.

### 3.4.1 The CLT and the Delta Method

Large sample theory, also called asymptotic theory, is used to approximate the distribution of an estimator when the sample size  $n$  is large. This theory is extremely useful if the exact sampling distribution of the estimator is complicated or unknown. To use this theory, one must determine what the estimator is estimating, the rate of convergence, the asymptotic distribution, and how large  $n$  must be for the approximation to be useful. Moreover, the (asymptotic) standard error (SE), an estimator of the asymptotic standard deviation, must be computable if the estimator is to be useful for inference.

**Theorem 3.7: the Central Limit Theorem (CLT).** Let  $Y_1, \dots, Y_n$  be iid with  $E(Y) = \mu$  and  $\text{VAR}(Y) = \sigma^2$ . Let the sample mean  $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$ . Then

$$\sqrt{n}(\bar{Y}_n - \mu) \xrightarrow{D} N(0, \sigma^2).$$

Hence

$$\sqrt{n} \left( \frac{\bar{Y}_n - \mu}{\sigma} \right) = \sqrt{n} \left( \frac{\sum_{i=1}^n Y_i - n\mu}{n\sigma} \right) \xrightarrow{D} N(0, 1).$$

Note that the sample mean is estimating the *population mean*  $\mu$  with a  $\sqrt{n}$  convergence rate, the asymptotic distribution is normal, and the SE =  $S/\sqrt{n}$  where  $S$  is the *sample standard deviation*. For many distributions the central limit theorem provides a good approximation if the sample size  $n > 30$ . A special case of the CLT is proven after Theorem 3.20.

**Notation.** The notation  $X \sim Y$  and  $X \stackrel{D}{=} Y$  both mean that the random variables  $X$  and  $Y$  have the same distribution. Hence  $F_X(x) = F_Y(y)$  for all real  $y$ . The notation  $Y_n \xrightarrow{D} X$  means that for large  $n$  we can approximate the cdf of  $Y_n$  by the cdf of  $X$ . The distribution of  $X$  is the limiting distribution or asymptotic distribution of  $Y_n$ . For the CLT, notice that

$$Z_n = \sqrt{n} \left( \frac{\bar{Y}_n - \mu}{\sigma} \right) = \left( \frac{\bar{Y}_n - \mu}{\sigma/\sqrt{n}} \right)$$

is the z-score of  $\bar{Y}$ . If  $Z_n \xrightarrow{D} N(0, 1)$ , then the notation  $Z_n \approx N(0, 1)$ , also written as  $Z_n \sim AN(0, 1)$ , means approximate the cdf of  $Z_n$  by the standard normal cdf. Similarly, the notation

$$\bar{Y}_n \approx N(\mu, \sigma^2/n),$$

also written as  $\bar{Y}_n \sim AN(\mu, \sigma^2/n)$ , means approximate the cdf of  $\bar{Y}_n$  as if  $\bar{Y}_n \sim N(\mu, \sigma^2/n)$ .

The two main applications of the CLT are to give the limiting distribution of  $\sqrt{n}(\bar{Y}_n - \mu)$  and the limiting distribution of  $\sqrt{n}(Y_n/n - \mu_X)$  for a random variable  $Y_n$  such that  $Y_n = \sum_{i=1}^n X_i$  where the  $X_i$  are iid with  $E(X) = \mu_X$  and  $\text{VAR}(X) = \sigma_X^2$ .

**Example 3.4.** a) Let  $Y_1, \dots, Y_n$  be iid  $\text{Ber}(\rho)$ . Then  $E(Y) = \rho$  and  $\text{VAR}(Y) = \rho(1 - \rho)$ . Hence

$$\sqrt{n}(\bar{Y}_n - \rho) \xrightarrow{D} N(0, \rho(1 - \rho))$$

by the CLT.

b) Now suppose that  $Y_n \sim \text{BIN}(n, \rho)$ . Then  $Y_n \stackrel{D}{=} \sum_{i=1}^n X_i$  where  $X_1, \dots, X_n$  are iid  $\text{Ber}(\rho)$ . Hence

$$\sqrt{n}\left(\frac{Y_n}{n} - \rho\right) \xrightarrow{D} N(0, \rho(1 - \rho))$$

since

$$\sqrt{n}\left(\frac{Y_n}{n} - \rho\right) \stackrel{D}{=} \sqrt{n}(\bar{X}_n - \rho) \xrightarrow{D} N(0, \rho(1 - \rho))$$

by a).

c) Now suppose that  $Y_n \sim \text{BIN}(k_n, \rho)$  where  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Then

$$\sqrt{k_n}\left(\frac{Y_n}{k_n} - \rho\right) \approx N(0, \rho(1 - \rho))$$

or

$$\frac{Y_n}{k_n} \approx N\left(\rho, \frac{\rho(1 - \rho)}{k_n}\right) \quad \text{or} \quad Y_n \approx N(k_n\rho, k_n\rho(1 - \rho)).$$

**Theorem 3.8: the Delta Method.** If  $g'(\theta) \neq 0$  and

$$\sqrt{n}(T_n - \theta) \xrightarrow{D} N(0, \sigma^2),$$

then

$$\sqrt{n}(g(T_n) - g(\theta)) \xrightarrow{D} N(0, \sigma^2[g'(\theta)]^2).$$

**Example 3.5.** Let  $Y_1, \dots, Y_n$  be iid with  $E(Y) = \mu$  and  $\text{VAR}(Y) = \sigma^2$ . Then by the CLT,

$$\sqrt{n}(\bar{Y}_n - \mu) \xrightarrow{D} N(0, \sigma^2).$$

Let  $g(\mu) = \mu^2$ . Then  $g'(\mu) = 2\mu \neq 0$  for  $\mu \neq 0$ . Hence

$$\sqrt{n}((\bar{Y}_n)^2 - \mu^2) \xrightarrow{D} N(0, 4\sigma^2\mu^2)$$

for  $\mu \neq 0$  by the delta method.

**Example 3.6.** Let  $X \sim \text{Binomial}(n, p)$  where the positive integer  $n$  is large and  $0 < p < 1$ . Find the limiting distribution of  $\sqrt{n} \left[ \left( \frac{X}{n} \right)^2 - p^2 \right]$ .

Solution. Example 3.4b gives the limiting distribution of  $\sqrt{n}(\frac{X}{n} - p)$ . Let  $g(p) = p^2$ . Then  $g'(p) = 2p$  and by the delta method,

$$\sqrt{n} \left[ \left( \frac{X}{n} \right)^2 - p^2 \right] = \sqrt{n} \left( g\left(\frac{X}{n}\right) - g(p) \right) \xrightarrow{D}$$

$$N(0, p(1-p)(g'(p))^2) = N(0, p(1-p)4p^2) = N(0, 4p^3(1-p)).$$

**Example 3.7.** Let  $X_n \sim \text{Poisson}(n\lambda)$  where the positive integer  $n$  is large and  $0 < \lambda$ .

a) Find the limiting distribution of  $\sqrt{n} \left( \frac{X_n}{n} - \lambda \right)$ .

b) Find the limiting distribution of  $\sqrt{n} \left[ \sqrt{\frac{X_n}{n}} - \sqrt{\lambda} \right]$ .

Solution. a)  $X_n \stackrel{D}{=} \sum_{i=1}^n Y_i$  where the  $Y_i$  are iid  $\text{Poisson}(\lambda)$ . Hence  $E(Y) = \lambda = \text{Var}(Y)$ . Thus by the CLT,

$$\sqrt{n} \left( \frac{X_n}{n} - \lambda \right) \stackrel{D}{=} \sqrt{n} \left( \frac{\sum_{i=1}^n Y_i}{n} - \lambda \right) \xrightarrow{D} N(0, \lambda).$$

b) Let  $g(\lambda) = \sqrt{\lambda}$ . Then  $g'(\lambda) = \frac{1}{2\sqrt{\lambda}}$  and by the delta method,

$$\sqrt{n} \left[ \sqrt{\frac{X_n}{n}} - \sqrt{\lambda} \right] = \sqrt{n} \left( g\left(\frac{X_n}{n}\right) - g(\lambda) \right) \xrightarrow{D} N\left(0, \lambda (g'(\lambda))^2\right) = N\left(0, \lambda \frac{1}{4\lambda}\right) = N\left(0, \frac{1}{4}\right).$$

**Example 3.8.** Let  $Y_1, \dots, Y_n$  be independent and identically distributed (iid) from a Gamma( $\alpha, \beta$ ) distribution.

a) Find the limiting distribution of  $\sqrt{n} (\bar{Y} - \alpha\beta)$ .

b) Find the limiting distribution of  $\sqrt{n} ((\bar{Y})^2 - c)$  for appropriate constant  $c$ .

Solution: a) Since  $E(Y) = \alpha\beta$  and  $V(Y) = \alpha\beta^2$ , by the CLT

$$\sqrt{n} (\bar{Y} - \alpha\beta) \xrightarrow{D} N(0, \alpha\beta^2).$$

b) Let  $\mu = \alpha\beta$  and  $\sigma^2 = \alpha\beta^2$ . Let  $g(\mu) = \mu^2$  so  $g'(\mu) = 2\mu$  and  $[g'(\mu)]^2 = 4\mu^2 = 4\alpha^2\beta^2$ . Then by the delta method,  $\sqrt{n} ((\bar{Y})^2 - c) \xrightarrow{D} N(0, \sigma^2[g'(\mu)]^2) = N(0, 4\alpha^3\beta^4)$  where  $c = \mu^2 = \alpha^2\beta^2$ .

### 3.4.2 Modes of Convergence and Consistency

**Definition 3.6.** Let  $\{Z_n, n = 1, 2, \dots\}$  be a sequence of random variables with cdfs  $F_n$ , and let  $X$  be a random variable with cdf  $F$ . Then  $Z_n$  **converges in distribution to  $X$** , written

$$Z_n \xrightarrow{D} X,$$

or  $Z_n$  *converges in law to  $X$* , written  $Z_n \xrightarrow{L} X$ , if

$$\lim_{n \rightarrow \infty} F_n(t) = F(t)$$

at each continuity point  $t$  of  $F$ . The distribution of  $X$  is called the **limiting distribution** or the **asymptotic distribution** of  $Z_n$ .

An important fact is that **the limiting distribution does not depend on the sample size  $n$** . Notice that the CLT and delta method give the limiting distributions of  $Z_n = \sqrt{n}(\bar{Y}_n - \mu)$  and  $Z_n = \sqrt{n}(g(T_n) - g(\theta))$ , respectively.

Convergence in distribution is useful because if the distribution of  $X_n$  is unknown or complicated and the distribution of  $X$  is easy to use, then for large  $n$  we can approximate the probability that  $X_n$  is in an interval by the probability that  $X$  is in the interval. To see this, notice that if  $X_n \xrightarrow{D} X$ , then  $P(a < X_n \leq b) = F_n(b) - F_n(a) \rightarrow F(b) - F(a) = P(a < X \leq b)$  if  $F$  is continuous at  $a$  and  $b$ .

Warning: convergence in distribution says that the cdf  $F_n(t)$  of  $X_n$  gets close to the cdf of  $F(t)$  of  $X$  as  $n \rightarrow \infty$  provided that  $t$  is a continuity point of  $F$ . Hence for any  $\epsilon > 0$  there exists  $N_t$  such that if  $n > N_t$ , then  $|F_n(t) - F(t)| < \epsilon$ . Notice that  $N_t$  depends on the value of  $t$ . Convergence in distribution does not imply that the random variables  $X_n \equiv X_n(\omega)$  converge to the random variable  $X \equiv X(\omega)$  for all  $\omega$ .

**Example 3.8.** Suppose that  $X_n \sim U(-1/n, 1/n)$ . Then the cdf  $F_n(x)$  of  $X_n$  is

$$F_n(x) = \begin{cases} 0, & x \leq -\frac{1}{n} \\ \frac{nx}{2} + \frac{1}{2}, & -\frac{1}{n} \leq x \leq \frac{1}{n} \\ 1, & x \geq \frac{1}{n}. \end{cases}$$

Sketching  $F_n(x)$  shows that it has a line segment rising from 0 at  $x = -1/n$  to 1 at  $x = 1/n$  and that  $F_n(0) = 0.5$  for all  $n \geq 1$ . Examining the cases  $x < 0$ ,  $x = 0$  and  $x > 0$  shows that as  $n \rightarrow \infty$ ,

$$F_n(x) \rightarrow \begin{cases} 0, & x < 0 \\ \frac{1}{2}, & x = 0 \\ 1, & x > 0. \end{cases}$$

Notice that if  $X$  is a random variable such that  $P(X = 0) = 1$ , then  $X$  has cdf

$$F_X(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0. \end{cases}$$

Since  $x = 0$  is the only discontinuity point of  $F_X(x)$  and since  $F_n(x) \rightarrow F_X(x)$  for all continuity points of  $F_X(x)$  (ie for  $x \neq 0$ ),

$$X_n \xrightarrow{D} X.$$

**Example 3.9.** Suppose  $Y_n \sim U(0, n)$ . Then  $F_n(t) = t/n$  for  $0 < t \leq n$  and  $F_n(t) = 0$  for  $t \leq 0$ . Hence  $\lim_{n \rightarrow \infty} F_n(t) = 0$  for  $t \leq 0$ . If  $t > 0$  and

$n > t$ , then  $F_n(t) = t/n \rightarrow 0$  as  $n \rightarrow \infty$ . Thus  $\lim_{n \rightarrow \infty} F_n(t) = 0$  for all  $t$  and  $Y_n$  does not converge in distribution to any random variable  $Y$  since  $H(t) \equiv 0$  is not a cdf.

**Definition 3.7.** A sequence of random variables  $X_n$  converges in distribution to a constant  $\tau(\theta)$ , written

$$X_n \xrightarrow{D} \tau(\theta), \text{ if } X_n \xrightarrow{D} X$$

where  $P(X = \tau(\theta)) = 1$ . The distribution of the random variable  $X$  is said to be degenerate at  $\tau(\theta)$ .

**Definition 3.8.** A sequence of random variables  $X_n$  converges in probability to a constant  $\tau(\theta)$ , written

$$X_n \xrightarrow{P} \tau(\theta),$$

if for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|X_n - \tau(\theta)| < \epsilon) = 1 \text{ or, equivalently, } \lim_{n \rightarrow \infty} P(|X_n - \tau(\theta)| \geq \epsilon) = 0.$$

The sequence  $X_n$  **converges in probability to  $X$** , written

$$X_n \xrightarrow{P} X,$$

if  $X_n - X \xrightarrow{P} 0$ .

Notice that  $X_n \xrightarrow{P} X$  if for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1, \text{ or, equivalently, } \lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0.$$

**Definition 3.9.** A sequence of estimators  $T_n$  of  $\tau(\theta)$  is **consistent** for  $\tau(\theta)$  if

$$T_n \xrightarrow{P} \tau(\theta)$$

for every  $\theta \in \Theta$ . If  $T_n$  is consistent for  $\tau(\theta)$ , then  $T_n$  is a **consistent estimator** of  $\tau(\theta)$ .

Consistency is a weak property that is usually satisfied by good estimators.  $T_n$  is a consistent estimator for  $\tau(\theta)$  if the probability that  $T_n$  falls in any neighborhood of  $\tau(\theta)$  goes to one, regardless of the value of  $\theta \in \Theta$ .



**Definition 3.10.** For a real number  $r > 0$ ,  $Y_n$  converges in  $r$ th mean to a random variable  $Y$ , written

$$Y_n \xrightarrow{r} Y,$$

if

$$E(|Y_n - Y|^r) \rightarrow 0$$

as  $n \rightarrow \infty$ . In particular, if  $r = 2$ ,  $Y_n$  **converges in quadratic mean** to  $Y$ , written

$$Y_n \xrightarrow{2} Y \quad \text{or} \quad Y_n \xrightarrow{\text{qm}} Y,$$

if

$$E[(Y_n - Y)^2] \rightarrow 0$$

as  $n \rightarrow \infty$ .

**Lemma 3.9: Generalized Chebyshev's Inequality.** Let  $u : \mathfrak{R} \rightarrow [0, \infty)$  be a nonnegative function. If  $E[u(Y)]$  exists then for any  $c > 0$ ,

$$P[u(Y) \geq c] \leq \frac{E[u(Y)]}{c}.$$

If  $\mu = E(Y)$  exists, then taking  $u(y) = |y - \mu|^r$  and  $\tilde{c} = c^r$  gives **Markov's Inequality:** for  $r > 0$  and any  $c > 0$ ,

$$P(|Y - \mu| \geq c) = P(|Y - \mu|^r \geq c^r) \leq \frac{E[|Y - \mu|^r]}{c^r}.$$

If  $r = 2$  and  $\sigma^2 = \text{VAR}(Y)$  exists, then we obtain

**Chebyshev's Inequality:**

$$P(|Y - \mu| \geq c) \leq \frac{\text{VAR}(Y)}{c^2}.$$

**Proof.** The proof is given for pdfs. For pmfs, replace the integrals by sums. Now

$$\begin{aligned} E[u(Y)] &= \int_{\mathfrak{R}} u(y)f(y)dy = \int_{\{y:u(y) \geq c\}} u(y)f(y)dy + \int_{\{y:u(y) < c\}} u(y)f(y)dy \\ &\geq \int_{\{y:u(y) \geq c\}} u(y)f(y)dy \end{aligned}$$

since the integrand  $u(y)f(y) \geq 0$ . Hence

$$E[u(Y)] \geq c \int_{\{y:u(y) \geq c\}} f(y)dy = cP[u(Y) \geq c]. \quad QED$$

The following proposition gives sufficient conditions for  $T_n$  to be a consistent estimator of  $\tau(\theta)$ . Notice that  $MSE_{\tau(\theta)}(T_n) \rightarrow 0$  for all  $\theta \in \Theta$  is equivalent to  $T_n \xrightarrow{qm} \tau(\theta)$  for all  $\theta \in \Theta$ .

**Proposition 3.10.** a) If

$$\lim_{n \rightarrow \infty} MSE_{\tau(\theta)}(T_n) = 0$$

for all  $\theta \in \Theta$ , then  $T_n$  is a consistent estimator of  $\tau(\theta)$ .

b) If

$$\lim_{n \rightarrow \infty} \text{VAR}_{\theta}(T_n) = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} E_{\theta}(T_n) = \tau(\theta)$$

for all  $\theta \in \Theta$ , then  $T_n$  is a consistent estimator of  $\tau(\theta)$ .

**Proof.** a) Using Lemma 3.9 with  $Y = T_n$ ,  $u(T_n) = (T_n - \tau(\theta))^2$  and  $c = \epsilon^2$  shows that for any  $\epsilon > 0$ ,

$$P_{\theta}(|T_n - \tau(\theta)| \geq \epsilon) = P_{\theta}[(T_n - \tau(\theta))^2 \geq \epsilon^2] \leq \frac{E_{\theta}[(T_n - \tau(\theta))^2]}{\epsilon^2}.$$

Hence

$$\lim_{n \rightarrow \infty} E_{\theta}[(T_n - \tau(\theta))^2] = \lim_{n \rightarrow \infty} MSE_{\tau(\theta)}(T_n) \rightarrow 0$$

is a sufficient condition for  $T_n$  to be a consistent estimator of  $\tau(\theta)$ .

b) Recall that

$$MSE_{\tau(\theta)}(T_n) = \text{VAR}_{\theta}(T_n) + [\text{Bias}_{\tau(\theta)}(T_n)]^2$$

where  $\text{Bias}_{\tau(\theta)}(T_n) = E_{\theta}(T_n) - \tau(\theta)$ . Since  $MSE_{\tau(\theta)}(T_n) \rightarrow 0$  if both  $\text{VAR}_{\theta}(T_n) \rightarrow 0$  and  $\text{Bias}_{\tau(\theta)}(T_n) = E_{\theta}(T_n) - \tau(\theta) \rightarrow 0$ , the result follows from a).  $QED$

The following result shows estimators that converge at a  $\sqrt{n}$  rate are consistent. Use this result and the delta method to show that  $g(T_n)$  is a consistent estimator of  $g(\theta)$ . Note that b) follows from a) with  $X_{\theta} \sim N(0, v(\theta))$ .

The WLLN shows that  $\bar{Y}$  is a consistent estimator of  $E(Y) = \mu$  if  $E(Y)$  exists.

**Proposition 3.11.** a) Let  $X$  be a random variable and  $0 < \delta \leq 1$ . If

$$n^\delta(T_n - \tau(\theta)) \xrightarrow{D} X$$

then  $T_n \xrightarrow{P} \tau(\theta)$ .

b) If

$$\sqrt{n}(T_n - \tau(\theta)) \xrightarrow{D} N(0, v(\theta))$$

for all  $\theta \in \Theta$ , then  $T_n$  is a consistent estimator of  $\tau(\theta)$ .

**Definition 3.11.** A sequence of random variables  $X_n$  *converges almost everywhere* (or *almost surely*, or *with probability 1*) to  $X$  if

$$P(\lim_{n \rightarrow \infty} X_n = X) = 1.$$

This type of convergence will be denoted by

$$X_n \xrightarrow{ae} X.$$

Notation such as “ $X_n$  converges to  $X$  ae” will also be used. Sometimes “ae” will be replaced with “as” or “wp1.” We say that  $X_n$  *converges almost everywhere* to  $\tau(\theta)$ , written

$$X_n \xrightarrow{ae} \tau(\theta),$$

if  $P(\lim_{n \rightarrow \infty} X_n = \tau(\theta)) = 1$ .

**Theorem 3.12.** Let  $Y_n$  be a sequence of iid random variables with  $E(Y_i) = \mu$ . Then

a) **Strong Law of Large Numbers (SLLN):**  $\bar{Y}_n \xrightarrow{ae} \mu$ , and

b) **Weak Law of Large Numbers (WLLN):**  $\bar{Y}_n \xrightarrow{P} \mu$ .

**Proof of WLLN when  $V(Y_i) = \sigma^2$ :** By Chebyshev’s inequality, for every  $\epsilon > 0$ ,

$$P(|\bar{Y}_n - \mu| \geq \epsilon) \leq \frac{V(\bar{Y}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \rightarrow 0$$

as  $n \rightarrow \infty$ . QED

In proving consistency results, there is an infinite sequence of estimators that depend on the sample size  $n$ . Hence the subscript  $n$  will be added to the estimators.

**Definition 3.12.** Lehmann (1999, p. 53-54): a) A sequence of random variables  $W_n$  is *tight* or *bounded in probability*, written  $W_n = O_P(1)$ , if for every  $\epsilon > 0$  there exist positive constants  $D_\epsilon$  and  $N_\epsilon$  such that

$$P(|W_n| \leq D_\epsilon) \geq 1 - \epsilon$$

for all  $n \geq N_\epsilon$ . Also  $W_n = O_P(X_n)$  if  $|W_n/X_n| = O_P(1)$ .

b) The sequence  $W_n = o_P(n^{-\delta})$  if  $n^\delta W_n = o_P(1)$  which means that

$$n^\delta W_n \xrightarrow{P} 0.$$

c)  $W_n$  has the *same order as  $X_n$  in probability*, written  $W_n \asymp_P X_n$ , if for every  $\epsilon > 0$  there exist positive constants  $N_\epsilon$  and  $0 < d_\epsilon < D_\epsilon$  such that

$$P(d_\epsilon \leq \left| \frac{W_n}{X_n} \right| \leq D_\epsilon) = P\left(\frac{1}{D_\epsilon} \leq \left| \frac{X_n}{W_n} \right| \leq \frac{1}{d_\epsilon}\right) \geq 1 - \epsilon$$

for all  $n \geq N_\epsilon$ .

d) Similar notation is used for a  $k \times r$  matrix  $\mathbf{A} = [a_{i,j}]$  if each element  $a_{i,j}$  has the desired property. For example,  $\mathbf{A} = O_P(n^{-1/2})$  if each  $a_{i,j} = O_P(n^{-1/2})$ .

**Definition 3.13.** Let  $W_n = \|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|$ .

a) If  $W_n \asymp_P n^{-\delta}$  for some  $\delta > 0$ , then both  $W_n$  and  $\hat{\boldsymbol{\mu}}_n$  have (tightness) **rate**  $n^\delta$ .

b) If there exists a constant  $\kappa$  such that

$$n^\delta(W_n - \kappa) \xrightarrow{D} X$$

for some nondegenerate random variable  $X$ , then both  $W_n$  and  $\hat{\boldsymbol{\mu}}_n$  have *convergence rate*  $n^\delta$ .

If  $W_n$  has convergence rate  $n^\delta$ , then  $W_n$  has tightness rate  $n^\delta$ , and the term “tightness” will often be omitted. Notice that if  $W_n \asymp_P X_n$ , then  $X_n \asymp_P W_n$ ,  $W_n = O_P(X_n)$  and  $X_n = O_P(W_n)$ . Notice that if  $W_n = O_P(n^{-\delta})$ , then  $n^\delta$  is a lower bound on the rate of  $W_n$ .

**Proposition 3.13.** Suppose there exists a constant  $\kappa$  such that

$$n^\delta(W_n - \kappa) \xrightarrow{D} X.$$

- a) Then  $W_n = O_P(n^{-\delta})$ .  
b) If  $X$  is not degenerate, then  $W_n \asymp_P n^{-\delta}$ .

The above result implies that if  $W_n$  has convergence rate  $n^\delta$ , then  $W_n$  has tightness rate  $n^\delta$ , and the term “tightness” will often be omitted. Part a) is proved, for example, in Lehmann (1999, p. 67).

The following result shows that if  $W_n \asymp_P X_n$ , then  $X_n \asymp_P W_n$ ,  $W_n = O_P(X_n)$  and  $X_n = O_P(W_n)$ . Notice that if  $W_n = O_P(n^{-\delta})$ , then  $n^\delta$  is a lower bound on the rate of  $W_n$ . As an example, if the CLT holds then  $\bar{Y}_n = O_P(n^{-1/3})$ , but  $\bar{Y}_n \asymp_P n^{-1/2}$ .

- Proposition 3.14.** a) If  $W_n \asymp_P X_n$  then  $X_n \asymp_P W_n$ .  
b) If  $W_n \asymp_P X_n$  then  $W_n = O_P(X_n)$ .  
c) If  $W_n \asymp_P X_n$  then  $X_n = O_P(W_n)$ .  
d)  $W_n \asymp_P X_n$  iff  $W_n = O_P(X_n)$  and  $X_n = O_P(W_n)$ .

**Proof.** a) Since  $W_n \asymp_P X_n$ ,

$$P(d_\epsilon \leq \left| \frac{W_n}{X_n} \right| \leq D_\epsilon) = P\left(\frac{1}{D_\epsilon} \leq \left| \frac{X_n}{W_n} \right| \leq \frac{1}{d_\epsilon}\right) \geq 1 - \epsilon$$

for all  $n \geq N_\epsilon$ . Hence  $X_n \asymp_P W_n$ .

b) Since  $W_n \asymp_P X_n$ ,

$$P(|W_n| \leq |X_n D_\epsilon|) \geq P(d_\epsilon \leq \left| \frac{W_n}{X_n} \right| \leq D_\epsilon) \geq 1 - \epsilon$$

for all  $n \geq N_\epsilon$ . Hence  $W_n = O_P(X_n)$ .

c) Follows by a) and b).

d) If  $W_n \asymp_P X_n$ , then  $W_n = O_P(X_n)$  and  $X_n = O_P(W_n)$  by b) and c). Now suppose  $W_n = O_P(X_n)$  and  $X_n = O_P(W_n)$ . Then

$$P(|W_n| \leq |X_n| D_{\epsilon/2}) \geq 1 - \epsilon/2$$

for all  $n \geq N_1$ , and

$$P(|X_n| \leq |W_n| 1/d_{\epsilon/2}) \geq 1 - \epsilon/2$$

for all  $n \geq N_2$ . Hence

$$P(A) \equiv P\left(\left| \frac{W_n}{X_n} \right| \leq D_{\epsilon/2}\right) \geq 1 - \epsilon/2$$

and

$$P(B) \equiv P(d_{\epsilon/2} \leq \left| \frac{W_n}{X_n} \right|) \geq 1 - \epsilon/2$$

for all  $n \geq N = \max(N_1, N_2)$ . Since  $P(A \cap B) = P(A) + P(B) - P(A \cup B) \geq P(A) + P(B) - 1$ ,

$$P(A \cap B) = P(d_{\epsilon/2} \leq \left| \frac{W_n}{X_n} \right| \leq D_{\epsilon/2}) \geq 1 - \epsilon/2 + 1 - \epsilon/2 - 1 = 1 - \epsilon$$

for all  $n \geq N$ . Hence  $W_n \asymp_P X_n$ . QED

The following result is used to prove the following Theorem 3.16 which says that if there are  $K$  estimators  $T_{j,n}$  of a parameter  $\beta$ , such that  $\|T_{j,n} - \beta\| = O_P(n^{-\delta})$  where  $0 < \delta \leq 1$ , and if  $T_n^*$  picks one of these estimators, then  $\|T_n^* - \beta\| = O_P(n^{-\delta})$ .

**Proposition 3.15: Pratt (1959).** Let  $X_{1,n}, \dots, X_{K,n}$  each be  $O_P(1)$  where  $K$  is fixed. Suppose  $W_n = X_{i_n,n}$  for some  $i_n \in \{1, \dots, K\}$ . Then

$$W_n = O_P(1). \quad (3.14)$$

**Proof.**

$$P(\max\{X_{1,n}, \dots, X_{K,n}\} \leq x) = P(X_{1,n} \leq x, \dots, X_{K,n} \leq x) \leq$$

$$F_{W_n}(x) \leq P(\min\{X_{1,n}, \dots, X_{K,n}\} \leq x) = 1 - P(X_{1,n} > x, \dots, X_{K,n} > x).$$

Since  $K$  is finite, there exists  $B > 0$  and  $N$  such that  $P(X_{i,n} \leq B) > 1 - \epsilon/2K$  and  $P(X_{i,n} > -B) > 1 - \epsilon/2K$  for all  $n > N$  and  $i = 1, \dots, K$ . Bonferroni's inequality states that  $P(\cap_{i=1}^K A_i) \geq \sum_{i=1}^K P(A_i) - (K - 1)$ . Thus

$$F_{W_n}(B) \geq P(X_{1,n} \leq B, \dots, X_{K,n} \leq B) \geq$$

$$K(1 - \epsilon/2K) - (K - 1) = K - \epsilon/2 - K + 1 = 1 - \epsilon/2$$

and

$$\begin{aligned} -F_{W_n}(-B) &\geq -1 + P(X_{1,n} > -B, \dots, X_{K,n} > -B) \geq \\ -1 + K(1 - \epsilon/2K) - (K - 1) &= -1 + K - \epsilon/2 - K + 1 = -\epsilon/2. \end{aligned}$$

Hence

$$F_{W_n}(B) - F_{W_n}(-B) \geq 1 - \epsilon \text{ for } n > N. \text{ QED}$$

**Theorem 3.16.** Suppose  $\|T_{j,n} - \beta\| = O_P(n^{-\delta})$  for  $j = 1, \dots, K$  where  $0 < \delta \leq 1$ . Let  $T_n^* = T_{i_n,n}$  for some  $i_n \in \{1, \dots, K\}$  where, for example,  $T_{i_n,n}$  is the  $T_{j,n}$  that minimized some criterion function. Then

$$\|T_n^* - \beta\| = O_P(n^{-\delta}). \quad (3.15)$$

**Proof.** Let  $X_{j,n} = n^\delta \|T_{j,n} - \beta\|$ . Then  $X_{j,n} = O_P(1)$  so by Proposition 3.15,  $n^\delta \|T_n^* - \beta\| = O_P(1)$ . Hence  $\|T_n^* - \beta\| = O_P(n^{-\delta})$ . QED

### 3.4.3 Slutsky's Theorem and Related Results

**Theorem 3.17: Slutsky's Theorem.** Suppose  $Y_n \xrightarrow{D} Y$  and  $W_n \xrightarrow{P} w$  for some constant  $w$ . Then

- a)  $Y_n + W_n \xrightarrow{D} Y + w$ ,
- b)  $Y_n W_n \xrightarrow{D} wY$ , and
- c)  $Y_n/W_n \xrightarrow{D} Y/w$  if  $w \neq 0$ .

**Theorem 3.18.** a) If  $X_n \xrightarrow{P} X$  then  $X_n \xrightarrow{D} X$ .

b) If  $X_n \xrightarrow{ae} X$  then  $X_n \xrightarrow{P} X$  and  $X_n \xrightarrow{D} X$ .

c) If  $X_n \xrightarrow{r} X$  then  $X_n \xrightarrow{P} X$  and  $X_n \xrightarrow{D} X$ .

d)  $X_n \xrightarrow{P} \tau(\theta)$  iff  $X_n \xrightarrow{D} \tau(\theta)$ .

e) If  $X_n \xrightarrow{P} \theta$  and  $\tau$  is continuous at  $\theta$ , then  $\tau(X_n) \xrightarrow{P} \tau(\theta)$ .

f) If  $X_n \xrightarrow{D} \theta$  and  $\tau$  is continuous at  $\theta$ , then  $\tau(X_n) \xrightarrow{D} \tau(\theta)$ .

Suppose that for all  $\theta \in \Theta$ ,  $T_n \xrightarrow{D} \tau(\theta)$ ,  $T_n \xrightarrow{r} \tau(\theta)$  or  $T_n \xrightarrow{ae} \tau(\theta)$ . Then  $T_n$  is a consistent estimator of  $\tau(\theta)$  by Theorem 3.18.

**Example 3.10.** Let  $Y_1, \dots, Y_n$  be iid with mean  $E(Y_i) = \mu$  and variance  $V(Y_i) = \sigma^2$ . Then the sample mean  $\bar{Y}_n$  is a consistent estimator of  $\mu$  since i) the SLLN holds (use Theorem 3.12 and 3.18), ii) the WLLN holds and iii) the CLT holds (use Proposition 3.11). Since

$$\lim_{n \rightarrow \infty} \text{VAR}_\mu(\bar{Y}_n) = \lim_{n \rightarrow \infty} \sigma^2/n = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} E_\mu(\bar{Y}_n) = \mu,$$

$\bar{Y}_n$  is also a consistent estimator of  $\mu$  by Proposition 3.10b. By the delta method and Proposition 3.11b,  $T_n = g(\bar{Y}_n)$  is a consistent estimator of  $g(\mu)$  if  $g'(\mu) \neq 0$  for all  $\mu \in \Theta$ . By Theorem 3.18e,  $g(\bar{Y}_n)$  is a consistent estimator of  $g(\mu)$  if  $g$  is continuous at  $\mu$  for all  $\mu \in \Theta$ .

**Theorem 3.19.** a) **Generalized Continuous Mapping Theorem:** If  $X_n \xrightarrow{D} X$  and the function  $g$  is such that  $P[X \in C(g)] = 1$  where  $C(g)$  is the set of points where  $g$  is continuous, then  $g(X_n) \xrightarrow{D} g(X)$ .

b) **Continuous Mapping Theorem:** If  $X_n \xrightarrow{D} X$  and the function  $g$  is continuous, then  $g(X_n) \xrightarrow{D} g(X)$ .

**Remark 3.3.** For Theorem 3.18, a) follows from Slutsky's Theorem by taking  $Y_n \equiv X = Y$  and  $W_n = X_n - X$ . Then  $Y_n \xrightarrow{D} Y = X$  and  $W_n \xrightarrow{P} 0$ . Hence  $X_n = Y_n + W_n \xrightarrow{D} Y + 0 = X$ . The convergence in distribution parts of b) and c) follow from a). Part f) follows from d) and e). Part e) implies that if  $T_n$  is a consistent estimator of  $\theta$  and  $\tau$  is a continuous function, then  $\tau(T_n)$  is a consistent estimator of  $\tau(\theta)$ . Theorem 3.19 says that convergence in distribution is preserved by continuous functions, and even some discontinuities are allowed as long as the set of continuity points is assigned probability 1 by the asymptotic distribution. Equivalently, the set of discontinuity points is assigned probability 0.

**Example 3.11.** (Ferguson 1996, p. 40): If  $X_n \xrightarrow{D} X$  then  $1/X_n \xrightarrow{D} 1/X$  if  $X$  is a continuous random variable since  $P(X = 0) = 0$  and  $x = 0$  is the only discontinuity point of  $g(x) = 1/x$ .

**Example 3.12.** Show that if  $Y_n \sim t_n$ , a  $t$  distribution with  $n$  degrees of freedom, then  $Y_n \xrightarrow{D} Z$  where  $Z \sim N(0, 1)$ .

Solution:  $Y_n \stackrel{D}{=} Z/\sqrt{V_n/n}$  where  $Z \perp V_n \sim \chi_n^2$ . If  $W_n = \sqrt{V_n/n} \xrightarrow{P} 1$ , then the result follows by Slutsky's Theorem. But  $V_n \stackrel{D}{=} \sum_{i=1}^n X_i^2$  where the iid  $X_i \sim \chi_1^2$ . Hence  $V_n/n \xrightarrow{P} 1$  by the WLLN and  $\sqrt{V_n/n} \xrightarrow{P} 1$  by Theorem 3.14e.

**Theorem 3.20: Continuity Theorem.** Let  $Y_n$  be sequence of random variables with characteristic functions  $\phi_n(t)$ . Let  $Y$  be a random variable with cf  $\phi(t)$ .

a)

$$Y_n \xrightarrow{D} Y \text{ iff } \phi_n(t) \rightarrow \phi(t) \forall t \in \mathfrak{R}.$$

b) Also assume that  $Y_n$  has mgf  $m_n$  and  $Y$  has mgf  $m$ . Assume that all of the mgfs  $m_n$  and  $m$  are defined on  $|t| \leq d$  for some  $d > 0$ . Then if  $m_n(t) \rightarrow m(t)$  as  $n \rightarrow \infty$  for all  $|t| < c$  where  $0 < c < d$ , then  $Y_n \xrightarrow{D} Y$ .



**Application: Proof of a Special Case of the CLT.** Following Rohatgi (1984, p. 569-9), let  $Y_1, \dots, Y_n$  be iid with mean  $\mu$ , variance  $\sigma^2$  and mgf  $m_Y(t)$  for  $|t| < t_o$ . Then

$$Z_i = \frac{Y_i - \mu}{\sigma}$$

has mean 0, variance 1 and mgf  $m_Z(t) = \exp(-t\mu/\sigma)m_Y(t/\sigma)$  for  $|t| < \sigma t_o$ . Want to show that

$$W_n = \sqrt{n} \left( \frac{\bar{Y}_n - \mu}{\sigma} \right) \xrightarrow{D} N(0, 1).$$

Notice that  $W_n =$

$$n^{-1/2} \sum_{i=1}^n Z_i = n^{-1/2} \sum_{i=1}^n \left( \frac{Y_i - \mu}{\sigma} \right) = n^{-1/2} \frac{\sum_{i=1}^n Y_i - n\mu}{\sigma} = \frac{n^{-1/2}}{\frac{1}{n}} \frac{\bar{Y}_n - \mu}{\sigma}.$$

Thus

$$\begin{aligned} m_{W_n}(t) &= E(e^{tW_n}) = E[\exp(tn^{-1/2} \sum_{i=1}^n Z_i)] = E[\exp(\sum_{i=1}^n tZ_i/\sqrt{n})] \\ &= \prod_{i=1}^n E[e^{tZ_i/\sqrt{n}}] = \prod_{i=1}^n m_Z(t/\sqrt{n}) = [m_Z(t/\sqrt{n})]^n. \end{aligned}$$

Set  $\psi(x) = \log(m_Z(x))$ . Then

$$\log[m_{W_n}(t)] = n \log[m_Z(t/\sqrt{n})] = n\psi(t/\sqrt{n}) = \frac{\psi(t/\sqrt{n})}{\frac{1}{n}}.$$

Now  $\psi(0) = \log[m_Z(0)] = \log(1) = 0$ . Thus by L'Hôpital's rule (where the derivative is with respect to  $n$ ),  $\lim_{n \rightarrow \infty} \log[m_{W_n}(t)] =$

$$\lim_{n \rightarrow \infty} \frac{\psi(t/\sqrt{n})}{\frac{1}{n}} = \lim_{n \rightarrow \infty} \frac{\psi'(t/\sqrt{n}) \left[ \frac{-t/2}{n^{3/2}} \right]}{\left( \frac{-1}{n^2} \right)} = \frac{t}{2} \lim_{n \rightarrow \infty} \frac{\psi'(t/\sqrt{n})}{\frac{1}{\sqrt{n}}}.$$

Now

$$\psi'(0) = \frac{m'_Z(0)}{m_Z(0)} = E(Z_i)/1 = 0,$$

so L'Hôpital's rule can be applied again, giving  $\lim_{n \rightarrow \infty} \log[m_{W_n}(t)] =$

$$\frac{t}{2} \lim_{n \rightarrow \infty} \frac{\psi''(t/\sqrt{n}) \left[ \frac{-t}{2n^{3/2}} \right]}{\left( \frac{-1}{2n^{3/2}} \right)} = \frac{t^2}{2} \lim_{n \rightarrow \infty} \psi''(t/\sqrt{n}) = \frac{t^2}{2} \psi''(0).$$

Now

$$\psi''(t) = \frac{d}{dt} \frac{m'_Z(t)}{m_Z(t)} = \frac{m''_Z(t)m_Z(t) - (m'_Z(t))^2}{[m_Z(t)]^2}.$$

So

$$\psi''(0) = m''_Z(0) - [m'_Z(0)]^2 = E(Z_i^2) - [E(Z_i)]^2 = 1.$$

Hence  $\lim_{n \rightarrow \infty} \log[m_{W_n}(t)] = t^2/2$  and

$$\lim_{n \rightarrow \infty} m_{W_n}(t) = \exp(t^2/2)$$

which is the  $N(0,1)$  mgf. Thus by the continuity theorem,

$$W_n = \sqrt{n} \left( \frac{\bar{Y}_n - \mu}{\sigma} \right) \xrightarrow{D} N(0, 1).$$

### 3.4.4 Multivariate Limit Theorems

Many of the univariate results of the previous 3 subsections 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 3.14.** 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 3.21 and 3.22 below are the multivariate extensions of the limit theorems in subsection 3.4.1. 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  $V(X) = \boldsymbol{\Sigma}x = \sigma^2$ .

**Theorem 3.21: 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  $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}x$ , then

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

where the sample mean

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

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  $\mathbf{D}_{\mathbf{g}(\boldsymbol{\theta})} = g'(\theta)$ .

**Theorem 3.22: the Multivariate Delta Method.** If

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

then

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

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

**Definition 3.15.** 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  $\mathbf{g}(\boldsymbol{\theta})$ .

**Proposition 3.23.** 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 3.24.** 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 3.25: Continuity Theorem.** Let  $\mathbf{X}_n$  be a sequence of  $k \times 1$  random vectors with characteristic function  $\phi_n(\mathbf{t})$  and let  $\mathbf{X}$  be a  $k \times 1$  random vector with cf  $\phi(\mathbf{t})$ . Then

$$\mathbf{X}_n \xrightarrow{D} \mathbf{X} \text{ iff } \phi_n(\mathbf{t}) \rightarrow \phi(\mathbf{t})$$

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

**Theorem 3.26: 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

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

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

- Theorem 3.27:** 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  $n$ :  $g(n) \uparrow \infty$ , eg  $g(n) = \sqrt{n}$ . 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 3.16.** Let  $\mathbf{A}_n = [a_{i,j}(n)]$  be an  $r \times c$  random matrix.

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

Recall that the smallest integer function  $\lceil x \rceil$  rounds up, eg  $\lceil 7.7 \rceil = 8$ .

**Theorem 3.28: Continuous Mapping Theorem.** Let  $\mathbf{X}_n \in \mathfrak{R}^k$ . If  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$  and if the function  $\mathbf{g} : \mathfrak{R}^k \rightarrow \mathfrak{R}^j$  is continuous, then  $\mathbf{g}(\mathbf{X}_n) \xrightarrow{D} \mathbf{g}(\mathbf{X})$ .

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

**Theorem 3.29:** 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  $\mathbf{W}_n \xrightarrow{D} \mathbf{C}$ , 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 3.30:** 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, eg  $Y_n = n^{-1/2}$ .)

- 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 3.31.** 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) If  $(T, \mathbf{C})$  is a consistent estimator of  $(\boldsymbol{\mu}, s \boldsymbol{\Sigma})$  with rate  $n^\delta$  where  $s > 0$

is some constant and  $0 < \delta \leq 0.5$ , then  $D_{\mathbf{x}}^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1} (\mathbf{x} - T) =$

$$s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + O_P(n^{-\delta}).$$

iii) 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  $s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

iv) 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)  $D_{\mathbf{x}}^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1} (\mathbf{x} - T) = (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T)^T [\mathbf{C}^{-1} - s^{-1} \boldsymbol{\Sigma}^{-1} + s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T) = (\mathbf{x} - \boldsymbol{\mu})^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu}) + (\mathbf{x} - T)^T [\mathbf{C}^{-1} - s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - T) + (\mathbf{x} - \boldsymbol{\mu})^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\boldsymbol{\mu} - T) + (\boldsymbol{\mu} - T)^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu}) + (\boldsymbol{\mu} - T)^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\boldsymbol{\mu} - T) = s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + O_P(n^{-\delta})$ .

iii) Following the proof for ii),  $D_{\mathbf{x}}^2(T, \mathbf{C}) = s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + o_P(1)$ . Alternatively,  $D_{\mathbf{x}}^2(T, \mathbf{C})$  is a continuous function of  $(T, \mathbf{C})$  if  $\mathbf{C} > 0$  for  $n > 10p$ . Hence  $D_{\mathbf{x}}^2(T, \mathbf{C}) \xrightarrow{P} D_{\mathbf{x}}^2(\boldsymbol{\mu}, s \boldsymbol{\Sigma})$ .

iv) 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  $\sqrt{n}(T - \boldsymbol{\mu})^T [\mathbf{C}^{-1} - \boldsymbol{\Sigma}^{-1}] \sqrt{n}(T - \boldsymbol{\mu}) = O_P(1) o_P(1) O_P(1) = o_P(1)$ .

### 3.5 Summary

1) If  $\mathbf{X}$  and  $\mathbf{Y}$  are  $p \times 1$  random vectors,  $\mathbf{a}$  a conformable constant vector, and  $\mathbf{A}$  and  $\mathbf{B}$  are conformable constant matrices, then

$$E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}), \quad E(\mathbf{a} + \mathbf{Y}) = \mathbf{a} + E(\mathbf{Y}), \quad \& \quad E(\mathbf{A} \mathbf{X} \mathbf{B}) = \mathbf{A} E(\mathbf{X}) \mathbf{B}.$$

Also

$$\text{Cov}(\mathbf{a} + \mathbf{A} \mathbf{X}) = \text{Cov}(\mathbf{A} \mathbf{X}) = \mathbf{A} \text{Cov}(\mathbf{X}) \mathbf{A}^T.$$

Note that  $E(\mathbf{A} \mathbf{Y}) = \mathbf{A} E(\mathbf{Y})$  and  $\text{Cov}(\mathbf{A} \mathbf{Y}) = \mathbf{A} \text{Cov}(\mathbf{Y}) \mathbf{A}^T$ .

2) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $E(\mathbf{X}) = \boldsymbol{\mu}$  and  $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ .

3) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and if  $\mathbf{A}$  is a  $q \times p$  matrix, then  $\mathbf{AX} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ . If  $\mathbf{a}$  is a  $p \times 1$  vector of constants, then  $\mathbf{X} + \mathbf{a} \sim N_p(\boldsymbol{\mu} + \mathbf{a}, \boldsymbol{\Sigma})$ . See Q2, HW2 E.

$$\text{Let } \mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \text{ and } \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

4) **All subsets of a MVN are MVN:**  $(X_{k_1}, \dots, X_{k_q})^T \sim N_q(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$  where  $\tilde{\boldsymbol{\mu}}_i = E(X_{k_i})$  and  $\tilde{\boldsymbol{\Sigma}}_{ij} = \text{Cov}(X_{k_i}, X_{k_j})$ . In particular,  $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$  and  $\mathbf{X}_2 \sim N_{p-q}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$ . If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are independent iff  $\boldsymbol{\Sigma}_{12} = \mathbf{0}$ .

5)

$$\text{Let } \begin{pmatrix} Y \\ X \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_Y \\ \mu_X \end{pmatrix}, \begin{pmatrix} \sigma_Y^2 & \text{Cov}(Y, X) \\ \text{Cov}(X, Y) & \sigma_X^2 \end{pmatrix} \right).$$

Also recall that the *population correlation* between  $X$  and  $Y$  is given by

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{VAR}(X)}\sqrt{\text{VAR}(Y)}} = \frac{\sigma_{X,Y}}{\sigma_X\sigma_Y}$$

if  $\sigma_X > 0$  and  $\sigma_Y > 0$ .

6) The conditional distribution of a MVN is MVN. If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the conditional distribution of  $\mathbf{X}_1$  given that  $\mathbf{X}_2 = \mathbf{x}_2$  is multivariate normal with mean  $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$  and covariance matrix  $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$ . That is,

$$\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2 \sim N_q(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

7) Notation:

$$\mathbf{X}_1 | \mathbf{X}_2 \sim N_q(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{X}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

8) Be able to compute the above quantities if  $X_1$  and  $X_2$  are scalars.

9) A  $p \times 1$  random vector  $\mathbf{X}$  has an *elliptically contoured distribution*, if  $\mathbf{X}$  has density

$$f(\mathbf{z}) = k_p |\boldsymbol{\Sigma}|^{-1/2} g[(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})], \quad (3.16)$$

and we say  $\mathbf{X}$  has an elliptically contoured  $EC_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$  distribution. If the second moments exist, then

$$E(\mathbf{X}) = \boldsymbol{\mu} \quad (3.17)$$

and

$$\text{Cov}(\mathbf{X}) = c_X \boldsymbol{\Sigma} \quad (3.18)$$

for some constant  $c_X > 0$ .

10) The *population squared Mahalanobis distance*

$$U \equiv D^2 = D^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}). \quad (3.19)$$

For elliptically contoured distributions,  $U$  has pdf

$$h(u) = \frac{\pi^{p/2}}{\Gamma(p/2)} k_p u^{p/2-1} g(u). \quad (3.20)$$

$U \sim \chi_p^2$  if  $\mathbf{x}$  has a multivariate normal  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution.

11) The classical estimator  $(\bar{\mathbf{x}}, \mathbf{S})$  of multivariate location and dispersion is the sample mean and sample covariance matrix where

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \text{and} \quad \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T.$$

12) Let the  $p \times 1$  column vector  $T(\mathbf{W})$  be a multivariate location estimator, and let the  $p \times p$  symmetric positive definite matrix  $\mathbf{C}(\mathbf{W})$  be a dispersion estimator. Then the  $i$ th *squared sample Mahalanobis distance* is the scalar

$$D_i^2 = D_i^2(T(\mathbf{W}), \mathbf{C}(\mathbf{W})) = (\mathbf{x}_i - T(\mathbf{W}))^T \mathbf{C}^{-1}(\mathbf{W})(\mathbf{x}_i - T(\mathbf{W})) \quad (3.21)$$

for each observation  $\mathbf{x}_i$ . Notice that the Euclidean distance of  $\mathbf{x}_i$  from the estimate of center  $T(\mathbf{W})$  is  $D_i(T(\mathbf{W}), \mathbf{I}_p)$ . The classical Mahalanobis distance uses  $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ .

13) If  $p$  random variables come from an elliptically contoured distribution, then the subplots in the scatterplot matrix should be linear.

14) 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  $n \rightarrow \infty$ .

15) 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  $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}_{\mathbf{x}}$ , then

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

where the sample mean

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

16) 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  $\mathbf{A}\sqrt{n}(T_n - \boldsymbol{\mu}) = \sqrt{n}(\mathbf{A}T_n - \mathbf{A}\boldsymbol{\mu}) \xrightarrow{D} N_q(\mathbf{A}\boldsymbol{\theta}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .

17) Suppose  $\mathbf{A}$  is a conformable constant matrix and  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ . Then  $\mathbf{A}\mathbf{X}_n \xrightarrow{D} \mathbf{A}\mathbf{X}$ .

### 3.6 Complements

Johnson and Wichern (1988) and Mardia, Kent and Bibby (1979) are good references for multivariate statistical analysis based on the multivariate normal distribution. The elliptically contoured distributions generalize the multivariate normal distribution and are discussed (in increasing order of difficulty) in Johnson (1987), Fang, Kotz and Ng (1990), Fang and Anderson (1990), and Gupta and Varga (1993). Fang, Kotz and Ng (1990) sketch the history of elliptically contoured distributions while Gupta and Varga (1993) discuss matrix valued elliptically contoured distributions. Cambanis, Huang and Simons (1981), Chmielewski (1981) and Eaton (1986) are also important references. Also see Muirhead (1982, p. 30–42).

There are several PhD level texts on large sample theory including, in roughly increasing order of difficulty, Lehmann (1999), Ferguson (1996), Sen and Singer (1993), and Serfling (1980). Cramér (1946) is also an important reference, and White (1984) considers asymptotic theory for econometric applications. Also see DasGupta (2008), Davidson (1994), Jiang (2010), Polansky (2011), Sen, Singer and Pedrosa De Lima (2010) and van der Vaart (1998). Section 3.4 followed Olive (2012b, ch. 8) closely.

In analysis, convergence in probability is a special case of convergence in measure and convergence in distribution is a special case of weak convergence.

See Ash (1972, p. 322) and Sen and Singer (1993, p. 39). Almost sure convergence is also known as strong convergence. See Sen and Singer (1993, p. 34). Since  $\bar{Y} \xrightarrow{P} \mu$  iff  $\bar{Y} \xrightarrow{D} \mu$ , the WLLN refers to weak convergence. Technically the  $X_n$  and  $X$  need to share a common probability space for convergence in probability and almost sure convergence.

## 3.7 Problems

**PROBLEMS WITH AN ASTERISK \* ARE ESPECIALLY USEFUL.**

**3.1\***. Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} \sim N_4 \left( \begin{pmatrix} 49 \\ 100 \\ 17 \\ 7 \end{pmatrix}, \begin{pmatrix} 3 & 1 & -1 & 0 \\ 1 & 6 & 1 & -1 \\ -1 & 1 & 4 & 0 \\ 0 & -1 & 0 & 2 \end{pmatrix} \right).$$

- Find the distribution of  $X_2$ .
- Find the distribution of  $(X_1, X_3)^T$ .
- Which pairs of random variables  $X_i$  and  $X_j$  are independent?
- Find the correlation  $\rho(X_1, X_3)$ .

**3.2\***. Recall that if  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the conditional distribution of  $\mathbf{X}_1$  given that  $\mathbf{X}_2 = \mathbf{x}_2$  is multivariate normal with mean  $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$  and covariance matrix  $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$ .

Let  $\sigma_{12} = \text{Cov}(Y, X)$  and suppose  $Y$  and  $X$  follow a bivariate normal distribution

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim N_2 \left( \begin{pmatrix} 49 \\ 100 \end{pmatrix}, \begin{pmatrix} 16 & \sigma_{12} \\ \sigma_{12} & 25 \end{pmatrix} \right).$$

- If  $\sigma_{12} = 0$ , find  $Y|X$ . Explain your reasoning.
- If  $\sigma_{12} = 10$  find  $E(Y|X)$ .
- If  $\sigma_{12} = 10$ , find  $\text{Var}(Y|X)$ .

**3.3.** Let  $\sigma_{12} = \text{Cov}(Y, X)$  and suppose  $Y$  and  $X$  follow a bivariate normal distribution

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim N_2 \left( \begin{pmatrix} 15 \\ 20 \end{pmatrix}, \begin{pmatrix} 64 & \sigma_{12} \\ \sigma_{12} & 81 \end{pmatrix} \right).$$

- a) If  $\sigma_{12} = 10$  find  $E(Y|X)$ .
- b) If  $\sigma_{12} = 10$ , find  $\text{Var}(Y|X)$ .
- c) If  $\sigma_{12} = 10$ , find  $\rho(Y, X)$ , the correlation between  $Y$  and  $X$ .

**3.4.** Suppose that

$$\mathbf{X} \sim (1 - \gamma)EC_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g_1) + \gamma EC_p(\boldsymbol{\mu}, c\boldsymbol{\Sigma}, g_2)$$

where  $c > 0$  and  $0 < \gamma < 1$ . Following Example 3.2, show that  $\mathbf{X}$  has an elliptically contoured distribution assuming that all relevant expectations exist.

**3.5.** In Proposition 3.5b, show that if the second moments exist, then  $\boldsymbol{\Sigma}$  can be replaced by  $\text{Cov}(\mathbf{X})$ .

crancap	hdlen	hdht	Data for 3.6
1485	175	132	
1450	191	117	
1460	186	122	
1425	191	125	
1430	178	120	
1290	180	117	
90	75	51	

**3.6\*.** The table ( $\mathbf{W}$ ) above represents 3 head measurements on 6 people and one ape. Let  $X_1 = \text{cranial capacity}$ ,  $X_2 = \text{head length}$  and  $X_3 = \text{head height}$ . Let  $\mathbf{x} = (X_1, X_2, X_3)^T$ . Several multivariate location estimators, including the coordinatewise median and sample mean, are found by applying a univariate location estimator to each random variable and then collecting the results into a vector. a) Find the coordinatewise median  $\text{MED}(\mathbf{W})$ .

- b) Find the sample mean  $\bar{\mathbf{x}}$ .

**3.7.** Using the notation in Proposition 3.6, show that if the second moments exist, then

$$\Sigma_{XX}^{-1} \Sigma_{XY} = [\text{Cov}(\mathbf{X})]^{-1} \text{Cov}(\mathbf{X}, Y).$$

**3.8.** Using the notation under Lemma 3.4, show that if  $\mathbf{X}$  is elliptically contoured, then the conditional distribution of  $\mathbf{X}_1$  given that  $\mathbf{X}_2 = \mathbf{x}_2$  is also elliptically contoured.

**3.9\*.** Suppose  $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ . Find the distribution of  $(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$  if  $\mathbf{X}$  is an  $n \times p$  full rank constant matrix and  $\boldsymbol{\beta}$  is a  $p \times 1$  constant vector.

**3.10.** Recall that  $\text{Cov}(\mathbf{X}, \mathbf{Y}) = E[(\mathbf{X} - E(\mathbf{X}))(\mathbf{Y} - E(\mathbf{Y}))^T]$ . Using the notation of Proposition 3.6, let  $(Y, \mathbf{X}^T)^T$  be  $EC_{p+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$  where  $Y$  is a random variable. Let the covariance matrix of  $(Y, \mathbf{X}^T)$  be

$$\text{Cov}((Y, \mathbf{X}^T)^T) = c \begin{pmatrix} \Sigma_{YY} & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_{XX} \end{pmatrix} = \begin{pmatrix} \text{VAR}(Y) & \text{Cov}(Y, \mathbf{X}) \\ \text{Cov}(\mathbf{X}, Y) & \text{Cov}(\mathbf{X}) \end{pmatrix}$$

where  $c$  is some positive constant. Show that  $E(Y|\mathbf{X}) = \alpha + \boldsymbol{\beta}^T \mathbf{X}$  where

$$\alpha = \mu_Y - \boldsymbol{\beta}^T \boldsymbol{\mu}_X \quad \text{and}$$

$$\boldsymbol{\beta} = [\text{Cov}(\mathbf{X})]^{-1} \text{Cov}(\mathbf{X}, Y).$$

**3.11.** (Due to R.D. Cook.) Let  $\mathbf{X}$  be a  $p \times 1$  random vector with  $E(\mathbf{X}) = \mathbf{0}$  and  $\text{Cov}(\mathbf{X}) = \Sigma$ . Let  $\mathbf{B}$  be any constant full rank  $p \times r$  matrix where  $1 \leq r \leq p$ . Suppose that for all such conforming matrices  $\mathbf{B}$ ,

$$E(\mathbf{X}|\mathbf{B}^T \mathbf{X}) = \mathbf{M}_B \mathbf{B}^T \mathbf{X}$$

where  $\mathbf{M}_B$  a  $p \times r$  constant matrix that depend on  $\mathbf{B}$ .

Using the fact that  $\Sigma \mathbf{B} = \text{Cov}(\mathbf{X}, \mathbf{B}^T \mathbf{X}) = E(\mathbf{X} \mathbf{X}^T \mathbf{B}) = E[E(\mathbf{X} \mathbf{X}^T \mathbf{B} | \mathbf{B}^T \mathbf{X})]$ , compute  $\Sigma \mathbf{B}$  and show that  $\mathbf{M}_B = \Sigma \mathbf{B} (\mathbf{B}^T \Sigma \mathbf{B})^{-1}$ . Hint: what acts as a constant in the inner expectation?

**3.12.** Let  $\mathbf{x}$  be a  $p \times 1$  random vector with covariance matrix  $\text{Cov}(\mathbf{x})$ . Let  $\mathbf{A}$  be an  $r \times p$  constant matrix and let  $\mathbf{B}$  be a  $q \times p$  constant matrix. Find  $\text{Cov}(\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{x})$  in terms of  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\text{Cov}(\mathbf{x})$ .

**3.13.** The table  $\mathbf{W}$  shown below represents 4 measurements on 5 people.

age	breadth	cephalic	size
39.00	149.5	81.9	3738
35.00	152.5	75.9	4261
35.00	145.5	75.4	3777
19.00	146.0	78.1	3904
0.06	88.5	77.6	933

- Find the sample mean  $\bar{\mathbf{x}}$ .
- Find the coordinatewise median  $\text{MED}(\mathbf{W})$ .

**3.14.** Suppose  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are iid  $p \times 1$  random vectors from a multivariate t-distribution with parameters  $\boldsymbol{\mu}$  and  $\Sigma$  with  $d$  degrees of freedom. Then  $E(\mathbf{x}_i) = \boldsymbol{\mu}$  and  $\text{Cov}(\mathbf{x}) = \frac{d}{d-2} \Sigma$  for  $d > 2$ . Assuming  $d > 2$ , find the limiting distribution of  $\sqrt{n}(\bar{\mathbf{x}} - \mathbf{c})$  for appropriate vector  $\mathbf{c}$ .

**3.15.** Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} \sim N_4 \left( \begin{pmatrix} 9 \\ 16 \\ 4 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.8 & -0.4 & 0 \\ 0.8 & 1 & -0.56 & 0 \\ -0.4 & -0.56 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \right)$$

- Find the distribution of  $X_3$ .

- b) Find the distribution of  $(X_2, X_4)^T$ .  
 c) Which pairs of random variables  $X_i$  and  $X_j$  are independent?  
 d) Find the correlation  $\rho(X_1, X_3)$ .

**3.16.** Suppose  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are iid  $p \times 1$  random vectors where

$$\mathbf{x}_i \sim (1 - \gamma)N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + \gamma N_p(\boldsymbol{\mu}, c\boldsymbol{\Sigma})$$

with  $0 < \gamma < 1$  and  $c > 0$ . Then  $E(\mathbf{x}_i) = \boldsymbol{\mu}$  and  $\text{Cov}(\mathbf{x}_i) = [1 + \gamma(c - 1)]\boldsymbol{\Sigma}$ . Find the limiting distribution of  $\sqrt{n}(\bar{\mathbf{x}} - \mathbf{c})$  for appropriate vector  $\mathbf{c}$ .

Let  $\mathbf{X}$  be an  $n \times p$  constant matrix and let  $\boldsymbol{\beta}$  be a  $p \times 1$  constant vector. Suppose  $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ . Find the distribution of  $\mathbf{H}\mathbf{Y}$  if  $\mathbf{H}^T = \mathbf{H} = \mathbf{H}^2$  is an  $n \times n$  matrix and if  $\mathbf{H}\mathbf{X} = \mathbf{X}$ . Simplify.

**3.17.** Recall that if  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the conditional distribution of  $\mathbf{X}_1$  given that  $\mathbf{X}_2 = \mathbf{x}_2$  is multivariate normal with mean  $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$  and covariance matrix  $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$ . Let  $Y$  and  $X$  follow a bivariate normal distribution

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim N_2 \left( \begin{pmatrix} 134 \\ 96 \end{pmatrix}, \begin{pmatrix} 24.5 & 1.1 \\ 1.1 & 23.0 \end{pmatrix} \right).$$

- a) Find  $E(Y|X)$ .  
 b) Find  $\text{Var}(Y|X)$ .

**3.18.** Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} \sim N_4 \left( \begin{pmatrix} 1 \\ 7 \\ 3 \\ 0 \end{pmatrix}, \begin{pmatrix} 4 & 0 & 2 & 1 \\ 0 & 1 & 0 & 0 \\ 2 & 0 & 3 & 1 \\ 1 & 0 & 1 & 5 \end{pmatrix} \right).$$

- a) Find the distribution of  $(X_1, X_4)^T$ .  
 b) Which pairs of random variables  $X_i$  and  $X_j$  are independent?  
 c) Find the correlation  $\rho(X_1, X_4)$ .

**3.19.** Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} \sim N_4 \left( \begin{pmatrix} 3 \\ 4 \\ 2 \\ 3 \end{pmatrix}, \begin{pmatrix} 3 & 2 & 1 & 1 \\ 2 & 4 & 1 & 0 \\ 1 & 1 & 2 & 0 \\ 1 & 0 & 0 & 3 \end{pmatrix} \right).$$

- a) Find the distribution of  $(X_1, X_3)^T$ .
- b) Which pairs of random variables  $X_i$  and  $X_j$  are independent?
- c) Find the correlation  $\rho(X_1, X_3)$ .

**3.20.** Suppose  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are iid  $p \times 1$  random vectors where  $E(\mathbf{x}_i) = e^{0.5}\mathbf{1}$  and  $\text{Cov}(\mathbf{x}_i) = (e^2 - e)\mathbf{I}_p$ . Find the limiting distribution of  $\sqrt{n}(\bar{\mathbf{x}} - \mathbf{c})$  for appropriate vector  $\mathbf{c}$ .

**3.21.** Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} \sim N_4 \left( \begin{pmatrix} 49 \\ 25 \\ 9 \\ 4 \end{pmatrix}, \begin{pmatrix} 2 & -1 & 3 & 0 \\ -1 & 5 & -3 & 0 \\ 3 & -3 & 5 & 0 \\ 0 & 0 & 0 & 4 \end{pmatrix} \right).$$

- a) Find the distribution of  $(X_1, X_3)^T$ .
- b) Which pairs of random variables  $X_i$  and  $X_j$  are independent?
- c) Find the correlation  $\rho(X_1, X_3)$ .

**3.22.** Recall that if  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the conditional distribution of  $\mathbf{X}_1$  given that  $\mathbf{X}_2 = \mathbf{x}_2$  is multivariate normal with mean  $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$  and covariance matrix  $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$ . Let  $Y$  and  $X$  follow a bivariate normal distribution

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim N_2 \left( \begin{pmatrix} 49 \\ 17 \end{pmatrix}, \begin{pmatrix} 3 & -1 \\ -1 & 4 \end{pmatrix} \right).$$

- a) Find  $E(Y|X)$ .
- b) Find  $\text{Var}(Y|X)$ .

**3.23.** Suppose  $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ . Find the distribution of  $(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$  if  $\mathbf{X}$  is an  $n \times p$  full rank constant matrix and  $\boldsymbol{\beta}$  is a  $p \times 1$  constant vector. Simplify.

**3.24.** Suppose  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are iid  $2 \times 1$  random vectors from a multivariate lognormal  $\text{LN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution. Let  $\mathbf{x}_i = (X_{i1}, X_{i2})^T$ . Following Press (2005, p. 149-150),

$E(X_{ij}) = \exp(\mu_j + \sigma_j^2/2)$ ,  $V(X_{ij}) = \exp(\sigma_j^2)[\exp(\sigma_j^2) - 1] \exp(2\mu_j)$  for  $j = 1, 2$ ,  
and  
 $\text{Cov}(X_{i1}, X_{i2}) = \exp[\mu_1 + \mu_2 + 0.5(\sigma_1^2 + \sigma_2^2) + \sigma_{12}][\exp(\sigma_{12}) - 1]$ . Find the  
limiting distribution of  $\sqrt{n}(\bar{\mathbf{x}} - \mathbf{c})$  for appropriate vector  $\mathbf{c}$ .