

### L03: Parameters of conditional distributions

#### 1. Parameters of conditional distributions

##### (1) Definitions

Suppose  $\begin{pmatrix} X \\ Y \end{pmatrix}$  has a joint distribution where  $X \in R^p$  and  $Y \in R^q$ .

Then  $E(Y|x) = (E(Y_i|x))_{q \times 1} = (\int \int_{R^q} y_i f_{Y|x}(y) dy_1, \dots, dy_q)_{q \times 1}$

$\text{Cov}(Y|X) = E(YY'|x) - E(Y|x)[E(Y|x)]'$ .

Here  $E(YY'|x) = (E(Y_i Y_j|x))_{q \times q} = (\int \int_{R^q} y_i y_j f_{Y|x}(y) dy_1, \dots, dy_q)_{q \times q}$

**Comments:**  $Y \sim (E(Y), \text{Cov}(Y))$ ,  $Y|X \sim (E(Y|X), \text{Cov}(Y|X))$  where  $E(Y|X)$  and  $\text{Cov}(Y|X)$  are vector-valued and matrix-valued functions of  $X$ , and hence are still random. So one can further consider  $E(Y|X) \sim (E[E(Y|X)], \text{Cov}(E(Y|X)))$

##### (2) Relations

(i)  $E[E(Y|X)] = E(Y)$ .

**Proof.** We show  $E[E(Y_i|X)] = E(Y_i)$ .

$$\begin{aligned} E[E(Y_i|X)] &= \int \int_{R^p} E(Y_i|x) f_X(x_1, \dots, x_p) dx_1, \dots, dx_p \\ &= \int \int_{R^p} [\int \int_{R^q} y_i f_{Y|x}(y_1, \dots, y_q) dy_1, \dots, dy_q] f_X(x_1, \dots, x_p) dx_1, \dots, dx_p \\ &= \int \int_{R^{p+q}} y_i f(x_1, \dots, x_p, y_1, \dots, y_q) dx_1, \dots, dx_p, dy_1, \dots, dy_q \\ &= E(Y_i) \end{aligned}$$

(ii)  $E[\text{Cov}(Y|X)] + \text{Cov}[E(Y|X)] = \text{Cov}(Y)$

$$\begin{aligned} \text{Cov}(Y) &= E(YY') - E(Y)[E(Y)]' \\ E[\text{Cov}(Y|X)] &= E\{E(YY'|X) - E(Y|X)[E(Y|X)]'\} \\ &= E(YY') - E\{E(Y|X)[E(Y|X)]'\} \\ \text{Cov}[E(Y|X)] &= E\{E(Y|X)[E(Y|X)]'\} - E[E(Y|X)]\{E[E(Y|X)]\}' \\ &= E\{E(Y|X)[E(Y|X)]'\} - E(Y)[E(Y)]' \end{aligned}$$

So  $E[\text{Cov}(Y|X)] + \text{Cov}[E(Y|X)] = E(YY') - E(Y)[E(Y)]' = \text{Cov}(Y)$ .

**Ex1:** For  $\begin{pmatrix} X \\ Y \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}\right)$ ,  $Y \sim N(\mu_y, \Sigma_{yy})$

$Y|X \sim N(\mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(X - \mu_x), \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy})$ .

$E(Y|X) = \mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(X - \mu_x) \sim N(\mu_y, \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy})$ .

So  $E[E(Y|X)] = \mu_y = E(Y)$  and

$E[\text{Cov}(Y|X)] + \text{Cov}[E(Y|X)] = (\Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}) + \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy} = \Sigma_{yy} = \text{Cov}(Y)$

#### 2. Independence

##### (1) Independence

$X \in R^p$  and  $Y \in R^q$  are independent  $\stackrel{\text{def}}{\iff} f_X(x) = f_{X|y}(x) \iff f_X(x) = \frac{f(x,y)}{f_Y(y)}$   
 $\iff f(x,y) = f_X(x) f_Y(y) \neq 0$   
 $\iff f_Y(y) = f_{Y|x}(y)$ .

##### (2) Impact of independence on parameters

$X \in R^p$  and  $Y \in R^q$  are independent  $\Rightarrow X$  and  $Y$  are uncorrelated.

**Proof.** We show  $\text{cov}(X_i, Y_j) = E(X_i Y_j) - E(X_i)E(Y_j) = 0$ .

$$\begin{aligned} E(X_i Y_j) &= \int \int_{R^{p+q}} x_i y_j f(x_1, \dots, x_p, y_1, \dots, y_q) dx_1, \dots, dx_p, dy_1, \dots, dy_q \\ &= \int \int_{R^{p+q}} x_i y_j f_X(x_1, \dots, x_p) f_Y(y_1, \dots, y_q) dx_1, \dots, dx_p, dy_1, \dots, dy_q \\ &= \int \int_{R^p} x_i f_X(x_1, \dots, x_p) dx_1, \dots, dx_p \int \int_{R^q} y_j f_Y(y_1, \dots, y_q) dy_1, \dots, dy_q \\ &= E(X_i)E(Y_j) \end{aligned}$$

(3) Independence for normal vectors

Suppose  $\begin{pmatrix} X \\ Y \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}\right)$ . Then

$X$  and  $Y$  are independent  $\iff X$  and  $Y$  are uncorrelated

**Proof.** Only give the sketch for  $\Leftarrow$ .

$X$  and  $Y$  are uncorrelated  $\implies \Sigma_{xy} = \text{Cov}(X, Y) = 0 \implies \Sigma_{xy} = 0$  and  $\Sigma_{yx} = \Sigma'_{xy} = 0$ .

With  $\begin{pmatrix} X \\ Y \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & 0 \\ 0 & \Sigma_{yy} \end{pmatrix}\right)$ ,  $X \sim N(\mu_x, \Sigma_{xx})$  and  $Y \sim N(\mu_y, \Sigma_{yy})$  one can check  $f(x_1, \dots, x_p, y_1, \dots, y_q) = f_X(x_1, \dots, x_p)f_Y(y_1, \dots, y_q)$ .

**Ex2:** In 4.3 on page 201  $X = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} \sim N(\mu, \Sigma)$  where  $\Sigma = \begin{pmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix}$ .

(i) Are  $X_1$  and  $X_2$  independent?

$\text{cov}(X_1, X_2) = -2 \neq 0$ . So  $X_1$  and  $X_2$  are not independent.

(ii) Are  $X_1$  and  $X_3$  independent?

$\text{cov}(X_1, X_3) = 0$ . So  $X_1$  and  $X_3$  are independent.

(iii) Are  $X_1$  and  $2X_1 + X_2 + X_3$  independent?

$$\text{Cov}(X_1, 2X_1 + X_2 + X_3) = (1, 0, 0) \begin{pmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} = 0.$$

So  $X_1$  and  $2X_1 + X_2 + X_3$  are independent.

3. Extended definitions for normality and independence

(1) Generalized normal vector

$X \sim N(\mu, \Sigma) \stackrel{\text{def}}{\iff} X = AZ_r + \mu$  where  $Z_r \sim N(0, I_r)$  and  $AA' = \Sigma$

So a vector is normal if it is transformed from a normal vector with pdf by function  $Ax + b$ .

Under this definition there is a convenient transformation for normal vectors,

$$X \sim N(\mu, \Sigma) \iff AX + \beta \sim N(A\mu + \beta, A\Sigma A')$$
 for all  $A$  and  $\beta$

**Caution:** For a  $p$ -dimensional normal vector, its support may not be  $R^p$ .

**Ex3:** For  $Z \sim N(0, 1^2)$ ,  $\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} Z \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}\right)$  has support  $x_1 = x_2$ .

(2) Concept of independence

$g(X)$  and  $h(Y)$  are independent if  $X$  and  $Y$  are independent by pdfs or pmfs.

Under this extended definition,  $X$  and  $Y$  are independent  $\implies X$  and  $Y$  are uncorrelated.

For normal vectors  $(X, Y)'$ ,  $X$  and  $Y$  are independent  $\iff X$  and  $Y$  are uncorrelated.

(3) Independence from normal vectors

Suppose  $X \sim N(\mu, \Sigma)$

(i)  $AX$  and  $BX$ :  $A\Sigma B' = 0 \iff AX$  and  $BX$  are independent.

**Pf:**  $AX$  and  $BX$  are independent  $\implies \text{Cov}(AX, BX) = 0 \iff A\Sigma B' = 0$ .

(ii)  $AX$  and  $X'BX$  where  $B' = B$ :  $A\Sigma B = 0 \implies AX$  and  $X'BX$  are independent.

**Pf:** By the compact form of EVD,  $B = P\Lambda P'$ . So  $A\Sigma B = A\Sigma P\Lambda P' = 0 \implies A\Sigma P = 0$ .

So  $AX$  and  $P'X$  are independent.

Hence  $AX$  and  $X'BX = (P'X)'\Lambda(P'X)$  are independent.

(iii)  $X'AX$  and  $X'BX$  where  $A' = A$  and  $B' = B$

$A\Sigma B = 0 \implies X'AX$  and  $X'BX$  are independent.

**Pf:** Skipped.

## L04: Principal components

### 1. Principal components for $X$

(1) Components of  $X \sim (\mu, \Sigma)$

The importance of the components of  $X \in R^p$  is usually measured by their variances. The larger the variance, the more information provided by the component.

The information provided by correlated components overlapped.

So an ideal vector would have  $\Sigma = V_x$  with non-increasing diagonal elements.

(2) Principal components of  $X$

Among all linear combinations of the components of  $X$  with unit vector coefficients, select one with the largest variance and call it the first principal component of  $X$ .

Once the first  $k$  principal components of  $X$  are selected, among all linear combinations of the components of  $X$  with unit vector coefficients and uncorrelated to the first  $k$  principal components already selected, select one with the largest variance and call it the  $k + 1$  the principal component of  $X$ .

### 2. Eigenvalue decomposition and principal components

(1) Eigenvalue decomposition of  $\Sigma$

If  $\Sigma v = \lambda v$  where  $0 \neq v \in R^p$ , then  $\lambda$  is an eigenvalue of  $\Sigma$ ,  $v$  is an eigenvector of  $\Sigma$  belonging to the eigenvalue  $\lambda$ .

$\Sigma > 0$  has  $p$  eigenvalues, all are positive numbers. Those eigenvalues are the roots of the characteristic polynomial  $|\Sigma - \lambda I|$ , i.e., the solutions to the characteristic equation  $|\Sigma - \lambda I| = 0$ .

All vectors except 0 in the eigenspace  $\mathcal{N}(\Sigma - \lambda I)$  are eigenvectors belonging to  $\lambda$ .

For real symmetric  $\Sigma$ , the eigenvectors belonging to different eigenvalues are orthogonal.

So one can find eigenvalue matrix  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$  with  $\lambda_1 \geq \dots \geq \lambda_p$  such that all diagonal elements of  $\Lambda$  are eigenvalues of  $\Sigma$ ; and eigenvector matrix  $P \in R^{p \times p}$  such that the  $P_i$ ,  $i$ th column of  $P$ , is the eigenvector for  $\lambda_i$ ,  $\|P_i\| = 1$ , and all columns of  $P$  are orthogonal. Thus  $P$  is an orthogonal matrix, i.e.,  $P' = P^{-1}$ .

Clearly  $\Sigma P = P \Lambda \iff \Sigma = P \Lambda P'$ , called the EVD for  $\Sigma$ .

$|\Sigma - \lambda I|$  is a polynomial of  $\lambda$  called characteristic polynomial of  $\Sigma$ .

Equation  $|\Sigma - \lambda I| = 0$  is called the characteristic equation for  $\Sigma$ .

The characteristic equation has  $p$  solutions  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  called the first, second, ..., the  $p$ th largest eigenvalue of  $\Sigma$ .

Solution  $x \neq 0$  to  $(\Sigma - \lambda_i I)x = 0$  is an eigenvector of  $\Sigma$  belonging to the eigenvalue  $\lambda_i$ .

It can be shown that  $\Sigma = P \Lambda P'$  where  $P = (P_1, \dots, P_p) \in R^{p \times p}$  is an orthogonal matrix such that  $P' P = I$ , and  $P_i$  is an eigenvector of  $\Sigma$  belonging to  $\lambda_i$ .

$\Sigma = P \Lambda P'$  is called an eigenvalue decomposition of  $\Sigma$ .

**Ex1:** Find eigenvalue decomposition for  $\Sigma = \begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix}$ .

$$|\Sigma - \lambda I| = \begin{vmatrix} 3 - \lambda & -1 \\ -1 & 3 - \lambda \end{vmatrix} = (\lambda - 3)^2 - 1 = (\lambda - 2)(\lambda - 4).$$

$$|\Sigma - \lambda I| = 0, \lambda_1 = 4 \text{ and } \lambda_2 = 2. \text{ So } \Lambda = \begin{pmatrix} 4 & 0 \\ 0 & 2 \end{pmatrix}.$$

$$(\Sigma - \lambda_1)x = 0, \begin{pmatrix} -1 & -1 \\ -1 & -1 \end{pmatrix} x = 0, x_1 = -x_2, x = c \begin{pmatrix} 1 \\ -1 \end{pmatrix}, P_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}.$$

$$(\Sigma - \lambda_2)x = 0, \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} x = 0, x_1 = x_2, x = c \begin{pmatrix} 1 \\ 1 \end{pmatrix}, P_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

$$P = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}. \Sigma = P \Lambda P' = \left[ \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \right] \begin{pmatrix} 4 & 0 \\ 0 & 2 \end{pmatrix} \left[ \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \right]'$$

(2) EVD and principal components

For  $X \sim (\mu, \Sigma)$  with EVD  $\Sigma = P\Lambda P'$  with non-decreasing eigenvalues, the components of  $Y = P'X$  are the principal components of  $X$ .

**Proof.** First, note that  $\text{Cov}(Y) = P'\Sigma P = P'P\Lambda P'P = \Lambda$ . So  $Y_i = P'_i X$  with  $\|P_i\| = 1$ ;  $\text{Cov}(Y_i, Y_j) = 0$  for  $i \neq j$ , and  $\text{var}(Y_i) = \lambda_i$  with  $\lambda_1 \geq \dots \geq \lambda_p$ .

Secondly we show that it is the best choice. When creating  $Y_1$ , in  $\{\alpha'X : \|\alpha\| = 1\}$ ,  $\text{var}(\alpha'X) = \alpha'P\Lambda P'\alpha = \beta'\Lambda\beta = \beta_1^2\lambda_1 + \dots + \beta_p^2\lambda_p$ . Here  $\beta_1^2 + \dots + \beta_p^2 = \alpha'P P'\alpha = 1$ .

When creating  $Y_2$ , in  $\{\alpha'X : \|\alpha\| = 1 \text{ and } \text{cov}(\alpha'X, Y_1) = 0\}$ ,  $\text{var}(\alpha'X) = \beta_1^2\lambda_1 + \dots + \beta_p^2\lambda_p$  where  $\beta_1^2 + \dots + \beta_p^2 = 1$ . But

$$0 = \text{cov}(\alpha'X, Y_1) = \text{cov}(\alpha'X, P'_1 X) = \alpha'P\Sigma P'_1 = \beta'\Lambda e_i = \beta_1\lambda_1.$$

So  $\beta_1 = 0$ . Hence  $\max[\text{var}(\alpha'X) : \|\alpha\| = 1, \text{cov}(\alpha'X, Y_1) = 0] = \lambda_2$ .

Similarly one can show  $\max[\text{var}(\alpha'X) : \|\alpha\| = 1 \text{ and } \text{cov}(\alpha'X, Y_i) = 0 \text{ for } i = 1, \dots, k]$  is  $\lambda_{k+1}$ .

**Ex2:** For  $X$  in Ex1,  $\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = P'X = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} X_1 - X_2 \\ X_1 + X_2 \end{pmatrix}$ . Here  $Y_1$  and  $Y_2$  are the first and the second principal components of  $X$ , and  $\text{Cov}(Y) = \begin{pmatrix} 4 & 0 \\ 0 & 2 \end{pmatrix}$ .

3. Properties

(1) Total variance

Total variance in  $X$  and total variance in  $Y$ , the principal component vector of  $X$ , are equal.

**Proof.**  $\text{Cov}(X) = \Sigma = P\Lambda P'$ . So  $\text{tr}(\Sigma) = \text{tr}(P\Lambda P') = \text{tr}(P'P\Lambda) = \text{tr}(\Lambda) = \lambda_1 + \dots + \lambda_p$ .

Here  $\text{tr}(\Sigma) = \text{var}(X_1) + \dots + \text{var}(X_p)$ , the total variance in  $X$ .

With principal component vector  $Y = P'X$ , the total variance in  $Y$  is

$$\text{var}(Y_1) + \dots + \text{var}(Y_p) = \lambda_1 + \dots + \lambda_p.$$

**Ex3:** In Ex1,  $\text{var}(X_1) + \text{var}(X_2) = 3 + 3 = 6$ . In Ex2,  $\text{var}(Y_1) + \text{var}(Y_2) = 4 + 2 = 6$ .

(2) Variations explained by principal components

The proportion of total variation in the original  $X$  explained by the  $i$ th principal component is

$$\frac{\text{var}(Y_i)}{\text{Total variance in } X} = \frac{\lambda_i}{\lambda_1 + \dots + \lambda_p}.$$

The proportion of total variance explained by the first  $k$  principal components is

$$\frac{\text{var}(Y_1) + \dots + \text{var}(Y_k)}{\text{Total variance in } X} = \frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p}.$$

Using principal components to achieve certain proportion of total variation explained, one can reduce the number of components in the vector.

**Ex4:** In Ex1, if we want to explain 60% of total variations in  $X$ , we have to use both  $X_1$  and  $X_2$ . But if using principal components, from Ex2, we see we only need the first principal component since  $\frac{4}{6} = 66.7\% > 60\%$ .

**Comment:** The concept of principal component is a good mathematical work. But the new component might be  $\frac{1}{\sqrt{3}}(\text{Age}) + \frac{1}{\sqrt{3}}(\text{Height}) + \frac{1}{\sqrt{3}}(\text{Number of hours in study})$  that does not have clear meaning.