

# MH4501 Multivariate Analysis

- Final Revision -

Naoki Honda

May 2019

## 1 PCA: Principal Component Analysis

### Population PCA (Result 6.1)

Assume the population covariance matrix  $\Sigma$  of random vector  $X = (X_1, X_2, \dots, X_p)^T$  is known, and  $\Sigma$  has eigenvalue-eigenvector pairs  $(\lambda_1, e_1), (\lambda_2, e_2), \dots, (\lambda_p, e_p)$ , subject to  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ . Then the  $k$ th **principal component** is given by

$$Y_k = e_k^T X = e_{k1}X_1 + e_{k2}X_2 + \dots + e_{kp}X_p \quad k = 1, 2, \dots, p$$

In addition,

$$\begin{aligned} \text{Var}(Y_k) &= e_k^T \Sigma e_k = \lambda_k & k = 1, 2, \dots, p \\ \text{Cov}(Y_j, Y_k) &= e_j^T \Sigma e_k = 0 & j \neq k \quad (j, k = 1, 2, \dots, p) \end{aligned}$$

### Proof of Result 6.1

Let  $U = (e_1, e_2, \dots, e_p)$  and  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$ . Thus  $\Sigma = U\Lambda U^T$  is the eigenvalue decomposition of  $\Sigma$ . Denote  $b = U^T a$ . Then

$$a^T \Sigma a = \frac{a^T \Sigma a}{a^T a} = \frac{a^T U \Lambda U^T a}{a^T U U^T a} = \frac{b^T \Lambda b}{b^T b} = \frac{\sum_{j=1}^p \lambda_j b_j^2}{\sum_{j=1}^p b_j^2} \leq \frac{\lambda_1 \sum_{j=1}^p b_j^2}{\sum_{j=1}^p b_j^2} = \lambda_1$$

The maximum is attained when  $a = e_1$ , since it gives  $b = U^T e_1 = (1, 0, 0, \dots, 0)^T$ ,  $b^T b = 1$  and  $b^T \Lambda b = \lambda_1$ .

This step shows that the 1st principal component is given by  $Y_1 = e_1^T X$ .

Secondly, we prove: when the first  $k$  ( $k = 1, 2, \dots, p-1$ ) principal component(s)  $Y_1 = e_1^T X, Y_2 = e_2^T X, \dots, Y_k = e_k^T X$  are decided, to have the  $(k+1)$ th principal component  $Y_{k+1} = a_{k+1}^T X$  orthogonal to every one of  $e_1, e_2, \dots, e_k$ .

For any  $j = 1, 2, \dots, k$ ,

$$\text{Cov}(Y_{k+1}, Y_j) = \text{Cov}(a_{k+1}^T X, e_j^T X) = a_{k+1}^T \Sigma e_j = a_{k+1}^T (\lambda_j e_j) = \lambda_j a_{k+1}^T e_j$$

So to have  $\text{Cov}(Y_{k+1}, Y_j)$  is equivalent to have  $a_{k+1}^T e_j = 0$ .  
This step prepares for the next one.<sup>1</sup>

### Variance (Result 6.2)

$$\sum_{k=1}^p \text{Var}(Y_k) = \sum_{j=1}^p \text{Var}(X_j)$$

Proof:

$$\begin{aligned} \sum_{k=1}^p \text{Var}(Y_k) &= \sum_{k=1}^p \lambda_k = \text{tr}(\Lambda) = \text{tr}(\Lambda U^T U) = \text{tr}(U \Lambda U^T) \\ &= \text{tr}(\Sigma) = \sum_{j=1}^p \sigma_{jj} = \sum_{j=1}^p \text{Var}(X_j) \end{aligned}$$

In practice, some popular quantities of interest are:

% of variance explained by  
the  $k$ th principal component  $Y_k$ :  $\frac{\lambda_k}{\sum_{j=1}^p \lambda_j} \quad k = 1, 2, \dots, p$

Cumulative % of variance explained by  
the first  $k$  principal components:  $\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^p \lambda_j} \quad k = 1, 2, \dots, p$

### Correlation (Result 6.3)

$$\text{Cor}(Y_k, X_j) = \frac{e_{kj} \sqrt{\lambda_k}}{\sqrt{\sigma_{jj}}} \quad k, j = 1, 2, \dots, p$$

Proof:

Denote  $d_j = (0, \dots, 0, 1_{j\text{th}}, 0, \dots, 0)^T$  so that  $X_j = d_j^T X$ . Therefore

$$\begin{aligned} \text{Cov}(Y_k, X_j) &= \text{Cov}(X_j, Y_k) \\ &= \text{Cov}(d_j^T X, e_k^T X) = d_j^T \Sigma e_k = d_j^T (\lambda_k e_k) = \lambda_k (d_j^T e_k) = \lambda_k e_{kj} \end{aligned}$$

In addition,  $\text{Var}(Y_k) = \lambda_k$  and  $\text{Var}(X_j) = \sigma_{jj}$ . Thus

$$\text{Cor}(Y_k, X_j) = \frac{\text{Cov}(Y_k, X_j)}{\sqrt{\text{Var}(Y_k)} \sqrt{\text{Var}(X_j)}} = \frac{e_{kj} \sqrt{\lambda_k}}{\sqrt{\sigma_{jj}}} \quad k, j = 1, 2, \dots, p$$

<sup>1</sup>The last part of the proof is omitted due to the similarity to the first part. Detail in P.9 - Lec #6

## 1.1 Sample PCA

What we discussed previously still applies in sample PCA. While in addition, since we have a sample  $x_1, x_2, \dots, x_n$ , in PCA, we would as well have a value of every principal component on every observation:  $y_{ik}$  ( $i = 1, 2, \dots, n$  and  $k = 1, 2, \dots, p$ ). These values are called PCA scores.

## 1.2 Sample PCA - Standardized Variables

### Sample PCA on Standardized Variables (Result 6.6)

Assume we have realizations  $x_1, x_2, \dots, x_n$  of random vector  $X = (X_1, X_2, \dots, X_p)^T$ , and the sample correlation matrix  $R$  has eigenvalue-eigenvector pairs  $(\lambda_1^Z, u_1), (\lambda_2^Z, u_2), \dots, (\lambda_p^Z, u_p)$ , subject to  $\lambda_1^Z \geq \lambda_2^Z \geq \dots \geq \lambda_p^Z \geq 0$ . Then the  $k$ th **sample principal component** obtained from standardized variables  $Z = (Z_1, Z_2, \dots, Z_p)^T$  is given by

$$Y_k^Z = u_k^T Z = u_{k1}Z_1 + u_{k2}Z_2 + \dots + u_{kp}Z_p \quad k = 1, 2, \dots, p$$

and the value of  $Y_k^Z$  on the  $i$ th observation is given by

$$y_{ik}^Z = u_{ki}^T z_i \quad i = 1, 2, \dots, n \quad k = 1, 2, \dots, p$$

Again  $\text{Var}(Y_k^Z) = \lambda_k^Z$  ( $k = 1, 2, \dots, p$ ) and  $\text{Cov}(Y_j^Z, Y_k^Z) = 0$  ( $j \neq k$ ). In addition,

$$\sum_{k=1}^p \text{Var}(Y_k^Z) = \sum_{j=1}^p \text{Var}(Z_j) = p$$

and

$$\text{Cor}(Y_k^Z, Z_j) = u_{kj} \sqrt{\lambda_k^Z} \quad k, j = 1, 2, \dots, p$$

## 2 CA: Cluster Analysis

### 2.1 Distance

#### Distance

Minkowski distance:  $d(x, y) = \left[ \sum_{j=1}^p |x_j - y_j|^m \right]^{\frac{1}{m}}$

When  $m = 1$ :  $d(x, y) = \sum_{j=1}^p |x_j - y_j|$  (Manhattan distance)

When  $m = 2$ :  $d(x, y) = \sqrt{\sum_{j=1}^p (x_j - y_j)^2}$  (Euclidean distance)

When  $m \rightarrow \infty$ :  $d(x, y) = \max_{j=1,2,\dots,p} |x_j - y_j|$  (Chebyshev distance)

Mahalanobis distance:  $d(x, y) = \sqrt{(x - y)^T S^{-1} (x - y)}$

(variance of Euclidean distance relative to  $S$ )

Canberra distance:  $d(x, y) = \sum_{j=1}^p \frac{|x_j - y_j|}{(x_j + y_j)}$

(weighted version of Manhattan distance)

Czekanowski distance:  $d(x, y) = 1 - \frac{2 \sum_{j=1}^p \min(x_j, y_j)}{\sum_{j=1}^p (x_j + y_j)}$

### 2.2 K-means Clustering

#### K-means Clustering

Given a collection of observations  $x_1, x_2, \dots, x_n$  and a well defined measure of distance  $d()$ , the K-means clustering works in an iterative fashion as below, with  $K$ , the desired number of clusters, pre-specified:

1. Initialization: assign the  $n$  observations into  $K$  clusters arbitrarily.
2. Find the center  $C_k$  ( $k = 1, 2, \dots, K$ ) of each cluster (based on current assignments), by taking average (mean) of the observations in each cluster respectively.
3. Assign each observation  $x_i$  ( $i = 1, 2, \dots, n$ ) to Cluster  $k(i)$ , by the criterion that

$$k(i) = \arg \min_k d(C_k, x_i)$$

If any reassignment happens, go back to Step 2; Otherwise, the algorithm is completed.

## 2.3 Hierarchical Clustering

### Motivation

1. One major drawback of K-means clustering is that the desired number of clusters must be pre-specified, which is not always feasible.
2. Hierarchical clustering instead provides a whole "path", from one extreme situation, namely that each individual observation constitutes a cluster, to the other extreme situation, namely that all observations are in one single cluster. Thus, one is able to review the whole path before specifying  $K$ .
3. There are different algorithms to implement hierarchical clustering, while in this course we only discuss the **agglomerative hierarchical clustering**, which involves merging small clusters which are similar (in some sense) into larger ones.
4. Therefore, besides a measure of distance between two observations, we also need a measure of distance between two clusters, which is conventionally called a **linkage** function.

#### Linkage functions

Single linkage:  $\min_{x \in A, y \in B} d(x, y)$

Complete linkage:  $\max_{x \in A, y \in B} d(x, y)$

Average linkage:  $\frac{1}{|A||B|} \sum_{x \in A} \sum_{y \in B} d(x, y)$

Centroid linkage:  $d(\bar{x}_A, \bar{x}_B)$

where  $\bar{x}_A = \frac{1}{|A|} \sum_{x \in A} x$  and  $\bar{x}_B = \frac{1}{|B|} \sum_{x \in B} x$

Given a collection of observations  $x_1, x_2, \dots, x_n$ , a well defined measure of distance  $d()$ , and a pre-specified linkage function to measure distance between clusters, the agglomerative clustering works in an iterative fashion as below:

#### Hierarchical Clustering (agglomerative)

1. Initialization: every individual observation constitutes a cluster respectively.
2. Calculate the distance (linkage) between every pair of clusters, find the pair that are closest to each other. Merge these two clusters into one cluster.
3. If all observations are in a single cluster, the algorithm is completed. Otherwise, go back to Step 2.

## 3 DA: Discriminant Analysis

### 3.1 Definition Preparation

- **Prior probabilities** measure the chance of an **unspecified** observation ( $x$ ) belonging to the populations:

$$p_k = Pr[x \in \pi_k] \quad \text{where } k = 1, 2, \dots, K$$

Clearly, it requires that  $\sum_{k=1}^K p_k = 1$

- When there is not enough prior knowledge, a naive choice of prior probability is simply  $p_1 = p_2 = \dots = p_K = \frac{1}{K}$
- Alternatively, prior probabilities can be estimated from sample proportions of classes:

$$p_k = \frac{|\{i : y_i = k\}|}{n} \quad k = 1, 2, \dots, K$$

- **Conditional probability**

$$p(t|k) = Pr[x \in R_t | x \in \pi_k] = \int_{R_t} f_k(x) dx \quad t, k = 1, 2, \dots, K$$

where  $f_k(x)$  is the PDF of population  $\pi_k$

- $\sum_{t=1}^K p(t|k) = 1 \quad (k = 1, 2, \dots, K)$
- $Pr[x \in R_t \& x \in \pi_k] = Pr[x \in \pi_k] \times Pr[x \in R_t | x \in \pi_k] = p_k p(t|k)$   
 $(t, k = 1, 2, \dots, K)$
- $\sum_{k=1}^K \sum_{t=1}^K p_k p(t|k) = 1$

- **Misclassification cost**

$c(t|k)$  = loss caused by misclassifying an observation into  $\pi_t$   
when it is actually from  $\pi_k \quad (t, k = 1, 2, \dots, K)$

- In this course, unless otherwise stated,  $c(1|1) = c(2|2) = \dots = c(K|K) = 0$
- In general,  $c(t|k) \neq c(k|t) \quad (t \neq k)$ .<sup>2</sup>

---

<sup>2</sup>For example, misclassifying a person with a certain disease as healthy, v.s. misclassifying a healthy person as diseased

## ECM: Expected Cost of Missclassification

$$ECM = \sum_{k=1}^K \sum_{t=1}^K p_k p(t|k) c(t|k)$$

<sup>a</sup>

<sup>a</sup>Note that ECM is a **function of classification rule**. If  $p_k, f_k(x), p(t|k)$ , and  $c(t|k)$  are all known or can be estimated (from the sample  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ ), the classification rule that minimizes ECM is called the **minimum ECM classification rule**.

### Result 8.1

Explicit solution of the minimum ECM classification rule is given by

$$R_m = \{x : \arg \min_t \sum_{k=1}^K p_k f_k(x) c(t|k) = m\} \quad m = 1, 2, \dots, K$$

That is, to classify  $x$  into  $\pi_m$ , if  $\sum_{k=1}^K p_k f_k(x) c(m|k)$  is the minimum among  $\sum_{k=1}^K p_k f_k(x) c(t|k)$  for all  $t = 1, 2, \dots, K$

## 3.2 $K = 2$ multivariate normal populations

Now suppose  $k = 2$ , and the 2 populations are both multivariate normal:

$$\pi_1 : N_p(\mu_1, \Sigma_1) \quad \& \quad \pi_2 : N_p(\mu_2, \Sigma_2)$$

Recall the ODF of multivariate normal distributions:

$$f_k(x) = \frac{1}{(2\pi)^{p/2} \det(\Sigma_k)^{1/2}} \exp \left[ -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right]$$

The minimum ECM classification rule can be reformulated as:

$$R_1 = \left\{ x : \frac{1}{2} \left( d_2(x) - d_1(x) + \ln \left[ \frac{\det(\Sigma_2)}{\det(\Sigma_1)} \right] \right) \geq \ln \left[ \frac{p_2 c(1|2)}{p_1 c(2|1)} \right] \right\}$$

$$R_2 = \Omega \setminus R_1 \quad \text{where } \Omega \text{ is the sample space}$$

As for quantities used in the rule:

- $\mu_k$  ( $k = 1, 2$ ): if not known, estimated by  $\bar{x}_k$
- $\Sigma_k$  ( $k = 1, 2$ ): if not known, estimated by  $S_k$
- $p_k$  ( $k = 1, 2$ ): if not known, estimated by  $\frac{|\{i: y_i = k\}|}{n}$
- $c(1|2)$  &  $c(2|1)$ : must be pre-specified

If we further have  $\Sigma_1 = \Sigma_2 = \Sigma$ ,

$$\ln \left[ \frac{\det(\Sigma_2)}{\det(\Sigma_1)} \right] = \ln 1 = 0$$

$$d_k(x) = (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) \quad (k = 1, 2)$$

Thus

$$\begin{aligned} d_{12}(x) &= \frac{1}{2} [d_2(x) - d_1(x)] \\ &= \frac{1}{2} [(x - \mu_2)^T \Sigma^{-1} (x - \mu_2) - (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)] \\ &= (\mu_1 - \mu_2)^T \Sigma^{-1} x - \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 + \mu_2) \end{aligned}$$

The minimum EMC classification rule can be reformulated as

$$R_1 = \left\{ x : d_{12}(x) \geq \ln \left[ \frac{p_2 c(1|2)}{p_1 c(2|1)} \right] \right\}$$

$$R_2 = \Omega \setminus R_1$$

In the rule,  $\Sigma$ , if unknown, can be estimated by

$$S_{pool} = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2}$$

In addition, if  $p_1 = p_2$  and  $c(1|2) = c(2|1)$ ,

$$R_1 = \{x : d_{12}(x) \geq 0\}$$

$$R_2 = \{x : d_{12}(x) < 0\}$$

### 3.3 $K \geq 3$ multivariate normal populations

- If  $p_1 = p_2 = \dots = p_K$  and  $c(t|k)$  is constant for all pairs  $t \neq k$ , to find  $m$  that maximizes  $f_m(x)$ , is equivalent with to find  $m$  that minimizes

$$-2\ln[f_m(x)] - p\ln[2\pi] = d_m(x) + \ln[\det(\Sigma_m)]$$

where  $d_m(x) = (x - \mu_m)^T \Sigma_m^{-1} (x - \mu_m)$

- If  $c(t|k)$  is constant for all pairs  $t \neq k$  while  $p_1, p_2, \dots, p_K$  are general, to find  $m$  that maximizes  $p_m f_m(x)$ , is equivalent with to find  $m$  that minimizes

$$-2\ln[p_m f_m(x)] - p\ln[2\pi] = d_m(x) + \ln[\det(\Sigma_m)] - 2\ln p_m$$



## 4 FA: Factor Analysis

The essential purpose of FA is to describe, if possible, the covariabce relationships among many variables in terms of a fer underlying, but **unobservable**, random quantities (latent variables) called factors.

### 4.1 The Orthogonal Factor Model

#### Notations

$$\begin{aligned} X_1 - \mu_1 &= l_{11}F_1 + l_{12}F_2 + \cdots + l_{1m}F_m + \epsilon_1 \\ X_2 - \mu_2 &= l_{21}F_1 + l_{22}F_2 + \cdots + l_{2m}F_m + \epsilon_2 \\ &\vdots \\ X_p - \mu_p &= l_{p1}F_1 + l_{p2}F_2 + \cdots + l_{pm}F_m + \epsilon_p \end{aligned}$$

$X_j$ :  $j$ th observable random variable ( $j = 1, 2, \dots, p$ )

$\mu_j$ : mean of  $X_j$  ( $j = 1, 2, \dots, p$ )

$F_k$ :  $k$ th **common factor** ( $k = 1, 2, \dots, m$ )

$l_{jk}$ : **loading** of the  $j$ th variable on the  $k$ th (common) factor ( $j = 1, 2, \dots, p$  and  $k = 1, 2, \dots, m$ )

$\epsilon_j$ :  $j$ th **specific factor** or **error** ( $j = 1, 2, \dots, p$ )

#### Assumptions

- $\mathbb{E}[F] = \mathbf{0}_{m \times 1}$  and  $\text{Cov}(F) = l_m$

- $\mathbb{E}[\epsilon] = \mathbf{0}_{p \times 1}$  and

$$\text{Cov}(\epsilon) = \Psi_{p \times p} = \begin{pmatrix} \psi_1 & 0 & \cdots & 0 \\ 0 & \psi_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \psi_p \end{pmatrix}$$

- $F$  and  $\epsilon$  are independent, thus  $\text{Cov}(F, \epsilon) = \mathbf{0}_{m \times p}$

### 4.2 Covariance Structure; OFM

$$X = \mu + LF + \epsilon$$

implies

$$\begin{aligned} (X - \mu)(X - \mu)^T &= (LF + \epsilon)(LF + \epsilon)^T \\ &= (LF)(LF)^T + \epsilon(LF)^T + (LF)\epsilon^T + \epsilon\epsilon^T \\ &= LFF^TL^T + \epsilon F^T L^T + LF\epsilon^T + \epsilon\epsilon^T \end{aligned}$$

Thus

$$\begin{aligned}
\Sigma &= E[(X - \mu)(X - \mu)^T] \\
&= E[LF F^T L^T + \epsilon F^T L^T + LF \epsilon^T + \epsilon \epsilon^T] \\
&= L E[FF^T] L^T + E[\epsilon] E[F^T] L^T + L E[F] E[\epsilon^T] + E[\epsilon \epsilon^T] \\
&= L \text{Cov}(F) L^T + \mathbf{0}_{p \times p} + \mathbf{0}_{p \times p} + \text{Cov}(\epsilon) \\
&= LL^T + \Psi
\end{aligned}$$

which implies

$$\sigma_{jj} = l_{j1}^2 + l_{j2}^2 + \cdots + l_{jm}^2 + \psi_j \quad j = 1, 2, \dots, p$$

$$\begin{aligned}
h_j^2 &= l_{j1}^2 + l_{j2}^2 + \cdots + l_{jm}^2 & \sigma_{jj}: & \text{variance of } X_j \\
& & j\text{th } \mathbf{communality} & \\
& & \text{portion of } \sigma_{jj} \text{ explained by the } m \text{ common factors} & \Sigma = LL^T + \\
\psi_j & & j\text{th } \mathbf{specific variance or uniqueness} & \\
& & \text{portion of } \sigma_{jj} \text{ explained by } \epsilon_j &
\end{aligned}$$

$\Psi$  also implies

$$\sigma_{ij} = l_{i1}l_{j1} + l_{i2}l_{j2} + \cdots + l_{im}l_{jm}$$

where  $i, j = 1, 2, \dots, p$  and  $i \neq j$

Back to

$$X = \mu + LF + \epsilon$$

it also implies

$$(X - \mu)F^T = (LF + \epsilon)F^T = LFF^T + \epsilon F^T$$

Thus

$$\begin{aligned}
\text{Cov}(X, F) &= E[(X - \mu)F^T] = E[LFF^T + \epsilon F^T] \\
&= L E[FF^T] + E[\epsilon] E[F^T] = L \text{Cov}(F) + \mathbf{0}_{p \times m} \\
&= L
\end{aligned}$$

i.e.

$$\text{Cov}(X_j, F_k) = l_{jk} \quad \text{where} \quad j = 1, 2, \dots, p \quad \text{and} \quad k = 1, 2, \dots, m$$

### 4.3 Principal Component Method

#### Principal Component Method

In PCA, we had

$$\begin{aligned}\Sigma &= U\Lambda U^T = U\Lambda^{1/2}\Lambda^{1/2}U^T \\ &= (U\Lambda^{1/2})(U\Lambda^{1/2})^T\end{aligned}$$

In consequence, the solution to FA is given by

$$\begin{aligned}L &= U\Lambda^{1/2} \\ \Psi &= \mathbf{0}_{p \times p}\end{aligned}$$

This solution is feasible in theory but definitely suboptimal, since  $m = p$ . A strategy to address this is to use only the largest  $m$  ( $m < p$  or even  $m \ll p$ ) eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_m$ , and their corresponding eigenvectors  $e_1, e_2, \dots, e_m$ .

1. Let  $U_m = (e_1, e_2, \dots, e_m)_{p \times m}$  and  $\Lambda_m = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)_{m \times m}$ , then

$$\tilde{L}_{p \times m} = U_m \Lambda_m^{1/2} = (\sqrt{\lambda_1}e_1, \sqrt{\lambda_2}e_2, \dots, \sqrt{\lambda_m}e_m)$$

2. After  $\tilde{L}$  is determined, let

$$\tilde{\Psi} = \text{diag}(\Sigma - \tilde{L}\tilde{L}^T)$$

which is a diagonal matrix with the same diagonal elements as  $\Sigma - \tilde{L}\tilde{L}^T$

3. At this stage, we think, more or less,

$$\Sigma \approx \tilde{L}\tilde{L}^T + \tilde{\Psi}$$

#### Principal Component Method; Residual Matrix

$$\Sigma - (\tilde{L}\tilde{L}^T + \tilde{\Psi})$$

is thus called the **residual matrix**, which has the following properties:

1. Symmetric matrix.
2. All diagonal elements are zero.

## 4.4 Covariance Structure; PCM

As for the (estimated) factor loading matrix

$$\tilde{L} = \begin{pmatrix} \tilde{l}_{11} & \tilde{l}_{12} & \cdots & \tilde{l}_{1m} \\ \tilde{l}_{21} & \tilde{l}_{22} & \cdots & \tilde{l}_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \tilde{l}_{p1} & \tilde{l}_{p2} & \cdots & \tilde{l}_{pm} \end{pmatrix}$$

- By row, we have interpreted it as

1. Portion of  $\sigma_{jj}$  explained by the  $m$  common factors:  $\tilde{h}_j^2 = \tilde{l}_{j1}^2 + \tilde{l}_{j2}^2 + \cdots + \tilde{l}_{jm}^2$
2. Portion of  $\sigma_{jj}$  explained by  $\epsilon_j$ :  $\tilde{\psi}_j = \sigma_{jj} - \tilde{h}_j^2$

- In addition, by column, we interpreted it as

$$\text{Portion of total variance explained by } F_k = \sum_{j=1}^p \tilde{l}_{jk}^2$$

1. In consequence,

$$\text{Portion of total variance explained by } F_k = \frac{\sum_{j=1}^p \tilde{l}_{jk}^2}{\sum_{j=1}^p \sigma_{jj}}$$

Cumulative portion of total variance

$$\text{explained by the first } k \text{ common factors} = \frac{\sum_{i=1}^k \sum_{j=1}^p \tilde{l}_{ji}^2}{\sum_{j=1}^p \sigma_{jj}}$$

2. In principle component method of estimation for FA,

$$\sum_{j=1}^p \tilde{l}_{jk}^2 = (\sqrt{\lambda_k} e_k)^T (\sqrt{\lambda_k} e_k) = \lambda_k e_k^T e_k = \lambda_k$$

## 4.5 The Maximum Likelihood Method

If  $F$  and  $\epsilon$  both can be assumed to be normally distributed, so is  $X = \mu + LF + \epsilon$ .

If we have a collection of observations  $x_1, x_2, \dots, x_n$ , the likelihood function would be

$$L(\mu, \Sigma | \{x_i\}_{i=1}^n) = \frac{1}{(2\pi)^{np/2} \det(\Sigma)^{n/2}} e^{-\sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)/2}$$

Substitute  $\Sigma = LL^T + \Psi$  in we have

$$L(\mu, L, \Psi | \{x_i\}_{i=1}^n) = \frac{1}{(2\pi)^{np/2} \det(LL^T + \Psi)^{n/2}} e^{-\sum_{i=1}^n (x_i - \mu)^T (LL^T + \Psi)^{-1} (x_i - \mu)/2}$$

Maximum likelihood solution,  $\hat{L}$  and  $\hat{\Psi}$ , can be obtained by maximizing the likelihood function, with assistance of computers.

## 4.6 Test of Sufficiency; MLM

### Test of Sufficiency

Whenever a FA model is build using sample covariance matrix  $S$ , a common question of interest is whether the model constructed fits the observed data well enough. That is, we want to test the (null) hypothesis

$$H_0 : \Sigma = \hat{L}\hat{L}^T + \hat{\Psi}$$

If the maximum likelihood method of estimation is employed, a testing statistic is proposed:

$$TS = \left( n - 1 - \frac{2p + 4m + 5}{6} \right) \ln \frac{\det(\hat{L}\hat{L}^T + \hat{\Psi})}{\det(\hat{\Sigma})}$$

Remarks:

1.  $\hat{L}$  and  $\hat{\Psi}$  are the maximum likelihood estimate of  $L$  and  $\Psi$  respectively
2.  $\hat{\Sigma}$  is the maximum likelihood estimate of  $\Sigma$ , that is,

$$\hat{\Sigma} = \frac{1}{n}SSCP = \frac{n-1}{n}S$$

3. If  $n$  is large, under  $H_0$ ,  $TS$  follows a  $\chi^2$  distribution with degree of freedom

$$d.f. = \frac{(p-m)^2}{2} - \frac{p+m}{2}$$

To have a positive  $d.f.$ , it must be the case that

$$m < p - \frac{\sqrt{8p+1} - 1}{2}$$

## 4.7 Factor Rotation

In a FA model

$$\Sigma = LL^T + \Psi$$

it is noticed that, for any  $m \times m$  orthogonal matrix  $G$ ,

$$(LG)(LG)^T = LGG^TL^T = LI_mL^T = LL^T$$

That is,

$$L^* = LG$$

$$\Psi^* = \Psi$$

is also a feasible solution  $\Sigma = L^*L^{*T} + \Psi^*$ .

In linear algebra, to right-multiply  $G$  to  $L$  can be interpreted as to **rotate**  $L$ .

Since any orthogonal matrix will keep  $L^*$  a feasible solution, an idea is proposed to keep rotating  $L$  until a simple structure is achieved. Here by "simple", it means every factor loading  $l_{jk}$  is either large or close to 0. The so-called **varimax criterion** is thus proposed:

### Varimax Criterion

$$V = \sum_{k=1}^m \left[ \sum_{j=1}^p l_{jk}^4 - \frac{1}{p} \left( \sum_{j=1}^p l_{jk}^2 \right)^2 \right]$$

Although it looks a little bit forbidding, the criterion in fact has a quite direct interpretation, since

$$\sum_{j=1}^p l_{jk}^4 - \frac{1}{p} \left( \sum_{j=1}^p l_{jk}^2 \right)^2 = \text{variance of } \{l_{1k}^2, l_{2k}^2, \dots, l_{pk}^2\}$$

Therefore, to maximize the varimax criterion, is effectively "spreading out"  $l_{jk}^2$  as much as possible, so that some of them are large, while the others are close to 0.

A weighted version of varimax criterion is defined to be

$$V_W = \sum_{k=1}^m \left[ \sum_{j=1}^p \left( \frac{l_{jk}^2}{h_j^2} \right) - \frac{1}{p} \left( \sum_{j=1}^p \frac{l_{jk}^2}{h_j^2} \right)^2 \right]$$

What are different before and after factor rotation:

1. The factor loadings:  $l_{jk}$  ( $j = 1, 2, \dots, p$  and  $k = 1, 2, \dots, m$ )
2. The variance explained by each factor:  $\sum_{j=1}^p l_{jk}^2$  ( $k = 1, 2, \dots, m$ )

While what are the same:

1. The communalities:  $h_j^{*2} = \sum_{k=1}^m l_{jk}^{*2} = (L^* L^{*T})_{jj} = (L L^T)_{jj} = h_j^2$  ( $j = 1, 2, \dots, p$ )
2. The specific variances:  $\psi_j^* = 1 - h_j^{*2} = 1 - h_j^2 = \psi_j$  ( $j = 1, 2, \dots, p$ )
3. The cumulative variance explained by **all** the  $m$  common factors:

$$\sum_{k=1}^m \sum_{j=1}^p l_{jk}^{*2} = \sum_{j=1}^p \sum_{k=1}^m l_{jk}^{*2} = \sum_{j=1}^p h_j^{*2} = \sum_{j=1}^p h_j^2 = \sum_{k=1}^m \sum_{j=1}^p l_{jk}^2$$

## 4.8 Factor Scores

When sample FA is performed, similar to the component scores in sample PCA, every factor is supposed to have a value on every observation.

However, due to different approaches of PCA and FA, the factor scores are not able to be directly calculated. Here we mention two methods of estimation:

1. The weighted least squares method (the Bartlett method):

$$\hat{f}_i = (\hat{L}^T \hat{\Psi}^{-1} \hat{L})^{-1} \hat{L}^T \hat{\Psi}^{-1} (x_i - \bar{x})$$

where  $\hat{f}_i = (\hat{f}_{i1}, \hat{f}_{i2}, \dots, \hat{f}_{im})^T$  are estimated values of  $F_1, F_2, \dots, F_m$  on  $x_i$ .

2. The regression method:

$$\hat{f}_i = \hat{L}^T S^{-1} (x_i - \bar{x})$$

again,  $\hat{f}_i = (\hat{f}_{i1}, \hat{f}_{i2}, \dots, \hat{f}_{im})^T$  are estimated values of  $F_1, F_2, \dots, F_m$  on  $x_i$ .

## 5 CCA: Canonical Correlation Analysis

CCA seeks to identify and quantify the associations between two sets of variables.

### 5.1 Population CCA

With the partitioned r.v., let

$$\begin{aligned} U &= a^T X^{(1)} \\ V &= b^T X^{(2)} \end{aligned}$$

where  $a_{p \times 1}$  and  $b_{q \times 1}$  are a pair of coefficient vectors. In consequence,

$$\begin{aligned} \text{Var}(U) &= a^T \Sigma_{11} a \\ \text{Var}(V) &= b^T \Sigma_{22} b \\ \text{Cov}(U, V) &= a^T \Sigma_{12} b = b^T \Sigma_{21} a \end{aligned}$$

and

$$\text{Cor}(U, V) = \frac{a^T \Sigma_{12} b}{\sqrt{a^T \Sigma_{11} a} \sqrt{b^T \Sigma_{22} b}}$$

#### Result 10.1

- $e_1, e_2, \dots, e_p$  (each with unit length) to be the eigenvectors of  $p \times p$  matrix  $\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}$ , with corresponding eigenvalues  $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$
- $f_1, f_2, \dots, f_p$  (each with unit length) to be the **first p** ( $p \leq q$ ) eigenvectors of  $q \times q$  matrix  $\Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}$ , with corresponding eigenvalues  $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$

Then the  $k$ th pair of canonical variables ( $k = 1, 2, \dots, p$ ) is explicitly given by

$$\begin{aligned} U_k &= e_k^T \Sigma_{11}^{-1/2} X^{(1)}, & \text{that is} & & a_k &= \Sigma_{11}^{-1/2} e_k \\ V_k &= f_k^T \Sigma_{22}^{-1/2} X^{(2)}, & \text{that is} & & b_k &= \Sigma_{22}^{-1/2} f_k \end{aligned}$$

For  $k = 1, 2, \dots, p$

$$\begin{aligned} f_k &= \frac{1}{\rho_k} \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1/2} e_k \\ e_k &= \frac{1}{\rho_k} \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2} f_k \end{aligned}$$

For  $k, l = 1, 2, \dots, p$  and  $k \neq l$

$$\begin{aligned} \text{Var}(U_k) &= 1, & \text{Cor}(U_k, U_l) &= 0 \\ \text{Var}(V_k) &= 1, & \text{Cor}(V_k, V_l) &= 0 \\ \text{Cor}(U_k, V_k) &= \rho_k, & \text{Cor}(U_k, V_l) &= 0 \end{aligned}$$



## 5.2 Variation explained by canonical variables

Now denote

$$A = \begin{pmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_p^T \end{pmatrix} \quad \text{and} \quad U = \begin{pmatrix} U_1 \\ U_2 \\ \vdots \\ U_p \end{pmatrix}$$

Then  $U = AX^{(1)}$ . Therefore,

$$\text{Cov}(U, X^{(1)}) = \text{Cov}(AX^{(1)}, X^{(1)}) = A \cdot \text{Cov}(X^{(1)}) = A\Sigma_{11}$$

At the meantime,  $X^{(1)} = A^{-1}U$ , thus

$$\text{Cov}(U, X^{(1)}) = \text{Cov}(U, A^{-1}U) = \text{Cov}(U)(A^{-1})^T = (A^{-1})^T$$

That is,

$$A\Sigma_{11} = (A^{-1})^T, \quad \text{or equivalently,} \quad A\Sigma_{11}A^T = I_p$$

### Variation explained by canonical variables

We have <sup>a</sup>

$$\text{Total variation in } X^{(1)} = \sum_{j=1}^p \sigma_{jj}$$

$$\text{Portion of variation in } X^{(1)} \text{ explained by } U_k = \sum_{j=1}^p \text{Cov}(U_k, X_j^{(1)})$$

where  $k = 1, 2, \dots, p$ .

$\text{Cov}(U_k, X_j^{(1)})$  is given by the  $kj$ -th element of matrix  $A\Sigma_{11} = (A^{-1})^T$ .

---

<sup>a</sup>The similar also applies to  $X^{(2)}, V, B$  and  $\Sigma_{22}$ .  
But notice that  $V_{p+1}, V_{p+2}, \dots, V_q$  are not canonical variables.

## 5.3 Sample CCA

Can be obtain by replacing  $\Sigma \rightarrow S$ .

In addition, every sample canonical variable has an observed value on each of the observations:

$$U_{ik} = a_k^T x_i^{(1)} = e_k^T S_{11}^{-1/2} x_i^{(1)}$$

$$V_{ik} = b_k^T x_i^{(2)} = f_k^T S_{22}^{-1/2} x_i^{(2)}$$

where  $i = 1, 2, \dots, n$  and  $k = 1, 2, \dots, p$ .

#### Test for $\Sigma_{12}$

$$H_0 : \Sigma_{12} = \mathbf{0}_{p \times q}$$

If we can assume that  $X$  follows the multivariate normal distribution  $N_{p+q}(\mu, \Sigma)$ , and the sample size  $n$  is large, a testing statistic is proposed

$$TS = -\left(n - 1 - \frac{1}{2}(p + q + 1)\right) \ln \prod_{k=1}^p (1 - \rho_k^2) \\ \sim \chi^2(pq)$$

where  $\rho_k$  ( $k = 1, 2, \dots, p$ ) are the canonical correlations obtained from sample CCA.

## 5.4 Sample CCA - Standardized variables

Sample CCA performed on original variables and sample CCA performed on standardized variables are related in a simple manner:

#### Result 10.3

- The canonical correlations are the same up to a change of sign:

$$\rho_k^2 = (\rho_k^Z)^2 \quad (k = 1, 2, \dots, p)$$

- If  $\bar{x}_j^{(1)} = 0$  (for all  $j = 1, 2, \dots, p$ ) and  $\bar{x}_j^{(2)} = 0$  (for all  $j = 1, 2, \dots, q$ )<sup>a</sup>, then the canonical variables are the same:

$$U_k = U_k^Z \quad \text{and} \quad V_k = V_k^Z \quad (k = 1, 2, \dots, p)$$

- Denote

$$\text{diag}(S) = D = \begin{pmatrix} D_1 & \mathbf{0} \\ \mathbf{0} & D_2 \end{pmatrix}$$

Then

$$a_k^Z = D_1^{1/2} a_k \quad \text{and} \quad b_k^Z = D_2^{1/2} b_k \quad (k = 1, 2, \dots, p)$$

---

<sup>a</sup>Or equivalently "both set of variables have mean 0".