

Lecture 2

Estimated covariance matrix

$$S = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^T$$

$\downarrow p \text{ vectors}$

If $p \geq n$ then S is singular (S^{-1} does not exist)

Why? Need to show S is not positive definite

$$\underline{\alpha}^T S \underline{\alpha} = 0 \text{ for some } \underline{\alpha} \neq \underline{0}$$

$$\begin{matrix} & p \text{ cols} \\ \begin{matrix} n \text{ rows} \\ \left(\begin{array}{c} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_n^T \end{array} \right) \end{matrix} & \left(\begin{array}{c} \underline{\alpha}_1 \\ \vdots \\ \underline{\alpha}_p \end{array} \right) = \left(\begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right) \end{matrix} \quad p \geq n \Rightarrow \text{there exists } \underline{\alpha} = \left(\begin{array}{c} \alpha_1 \\ \vdots \\ \alpha_p \end{array} \right) \text{ satisfying the equation}$$

$$\underline{x}_i^T \underline{\alpha} = 1 \text{ for all } i=1, \dots, n \Rightarrow \underline{\alpha}^T S \underline{\alpha} = 0$$

Simple sol'n (how to fix)

Add a positive definite matrix to S .

$$\begin{aligned} \tilde{S} &= S + \varepsilon I \\ \underline{\alpha}^T \tilde{S} \underline{\alpha} &= \underbrace{\underline{\alpha}^T S \underline{\alpha}}_{\geq 0} + \underbrace{\varepsilon \underline{\alpha}^T \underline{\alpha}}_{\text{small}} > 0 \text{ for } \underline{\alpha} \neq \underline{0} \end{aligned}$$

Correlation Matrix

$$R = \left(\begin{array}{cccc} 1 & \frac{\sigma_{12}}{\sigma_1 \sigma_2} & \cdots & \frac{\sigma_{1p}}{\sigma_1 \sigma_p} \\ \frac{\sigma_{21}}{\sigma_2 \sigma_1} & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\sigma_{p1}}{\sigma_p \sigma_1} & \cdots & \cdots & 1 \end{array} \right) \quad \begin{aligned} \sigma_{ij} &= \text{cov}(X_i, X_j) \\ \frac{\sigma_{ij}}{\sigma_i \sigma_j} &= \text{corr}(X_i, X_j) \end{aligned}$$

$$= D^{-1/2} C D^{-1/2} \quad \text{diagonal matrix with diagonal elements of } C.$$

- contains same information as C but R is dimensionless (no units)
 - independent of the scaling of X_1, \dots, X_p

- R is also the covariance of standardized variables

$$x'_j = \frac{x_j - E(x_j)}{\sqrt{\text{Var}(x_j)}}$$

- very important in multivariate analysis

- guarantees similar range for all (transformed) variables.

- Taking logs of positive variables often achieves the same goal.

- Estimate of R :

$$\hat{R} = \hat{D}^{-1/2} \hat{S} \hat{D}^{-1/2} \quad \text{diagonal element of } \hat{S}$$

Distance between points (observations)

Idea: Define clusters of observations

- If 2 multivariate observations are "close" then they should belong to the same cluster

How to define distance?

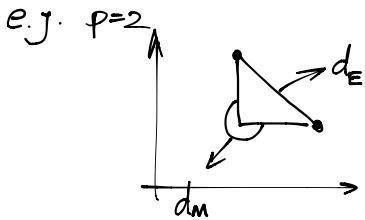
$$\underline{x_i} = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} \quad \underline{x_j} = \begin{pmatrix} x_{j1} \\ \vdots \\ x_{jp} \end{pmatrix}$$

Naire: ① Euclidean distance

$$d_E(\underline{x_i}, \underline{x_j}) = \left\{ \sum_{k=1}^p (x_{ik} - x_{jk})^2 \right\}^{\frac{1}{2}}$$

② Manhattan (L_1) distance:

$$d_M(\underline{x_i}, \underline{x_j}) = \sum_{k=1}^p |x_{ik} - x_{jk}|$$



What's wrong here?

- adding apples to oranges! (with diff units!)
- to fix this, scale all the variables to be dimensionless.

For example, define $s_k^2 = \overbrace{\frac{1}{n-1} \sum_{i=1}^n (x_{ik} - \bar{x}_k)^2}^{k\text{-th diagonal of } S}$

Take $d_E(\underline{x_i}, \underline{x_j}) = \left\{ \sum_{k=1}^p \left(\frac{x_{ik} - x_{jk}}{s_k} \right)^2 \right\}^{\frac{1}{2}}$

$$d_M(\underline{x_i}, \underline{x_j}) = \sum_{k=1}^p \left| \frac{x_{ik} - x_{jk}}{s_k} \right|$$

Distance matrix: Symmetric $n \times n$ matrix

- diagonal elements are 0
 - (i, j) element $d_{ij}(\underline{x_i}, \underline{x_j})$
- $\downarrow M \text{ or } E$

Applications (Later)

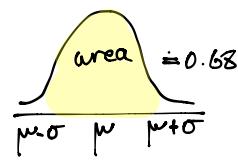
① Cluster analysis (find groups of similar observations)

② Multidimensional scaling: obtain a 2 dimensional representation of data with similar distance structure.

Multivariate Normal Distn

Review: (Univariate) Normal Distribution

$$X \sim N(\mu, \sigma^2), f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{1}{2\sigma^2} (x - \mu)^2 \right]$$



Given data x_1, \dots, x_n , a simple graphical check for normality is a normal probability (normal quantile-quantile plot)

Ordered data: $x_{(1)} \leq x_{(2)} \leq x_{(3)} \leq \dots \leq x_{(n)}$

Normal scores: $E[\bar{Z}_{(i)}] = \Phi^{-1}\left(\frac{i-\frac{1}{2}}{n}\right)$

$\uparrow \quad \downarrow$
 $N(0,1)$ order statistics quantile of $N(0,1)$

Plot $x_{(i)}$ vs $\Phi^{-1}\left(\frac{i-\frac{1}{2}}{n}\right)$ for $i=1, \dots, n$

- If data are normal, points should lie close to a straight line
- R function: `qqnorm`
- Shapiro-Wilk test (R function: `shapiro.test`)

How to extend to multivariate case?

- random vector $\underline{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_p \end{pmatrix}$

- want all one-dimensional projections
 $\underline{Q}^T \underline{X}$ to be univariate normal.

Multivariate normal

$\underline{X} \sim N_p(\underline{\mu}, \underline{C})$: assume \underline{C} is positive definite

$\downarrow \quad \downarrow$
mean vector covariance matrix

$$f(\underline{x}) = f(x_1, \dots, x_p) = \frac{1}{(2\pi)^{p/2} |\underline{C}|^{1/2}} \exp\left[-\frac{1}{2} (\underline{x} - \underline{\mu})^T \underline{C}^{-1} (\underline{x} - \underline{\mu})\right]$$

determinants

Example: $p=2$

$$\underline{C} = \begin{pmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix} \quad |\underline{C}| = \sigma_1^2 \sigma_2^2 - \rho^2 \sigma_1^2 \sigma_2^2 = \sigma_1^2 \sigma_2^2 (1 - \rho^2)$$

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

$$\text{so } f(x_1, x_2) = \frac{1}{2\pi \sigma_1 \sigma_2 (1 - \rho^2)^{1/2}} \exp\left[-\frac{1}{2(1 - \rho^2)} \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} - 2\rho \frac{(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1 \sigma_2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right]\right]$$

Conditional dist'n of X_1 given $X_2 = x_2$

$$f(x_1 | x_2) = \frac{f(x_1, x_2)}{f(x_2)}$$

$\Rightarrow X_1 | X_2 = x_2 \sim N(\mu_1 + \rho \frac{(x_2 - \mu_2)}{\sigma_2}, \sigma_1^2 (1 - \rho^2))$ conditional dist'n is still a Normal dist'n

Properties: $\underline{X} \sim N_p(\mu, C)$

$$\textcircled{1}: \underline{\alpha}^T \underline{X} \sim N(\underline{\alpha}^T \mu, \underline{\alpha}^T C \underline{\alpha})$$

$$\textcircled{2} \text{ If } \underline{X} = \begin{pmatrix} \underline{X}_1 \\ \underline{X}_2 \end{pmatrix} \xleftarrow{\text{length } r} \xleftarrow{\text{length } (p-r)}$$

then the conditional dist'n \underline{X}_1 given $\underline{X}_2 = \underline{x}_2$ is also multivariate normal
(mean & covariance revealed later!)

$$\textcircled{3} (\underline{X} - \mu)^T C^{-1} (\underline{X} - \mu) \sim \chi^2_{(p)} \quad \text{chi-square with df=p}$$

Application:

Checking multivariate normality

Data: $\underline{X}_1, \dots, \underline{X}_n \rightarrow$ is this multivariate normal?

- estimate μ, C by \bar{X} & S (as before)

Two approaches

① Mahalanobis distance (see textbook)

$$\text{Define } d_i^2 = (\underline{X}_i - \bar{X})^T S^{-1} (\underline{X}_i - \bar{X}) \quad \begin{matrix} \downarrow \text{est. of } C \\ \downarrow \text{est. of } \mu \end{matrix}$$

If $\underline{X}_1, \dots, \underline{X}_n$ are multivariate normal then d_1^2, \dots, d_n^2 should be approaching $\chi^2_{(p)}$

- order d_1^2, \dots, d_n^2 from smallest to largest

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

- compare $d_{(1)}^2, \dots, d_{(n)}^2$ to quantiles of $\chi^2_{(p)}$ dist'n.

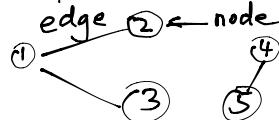
② Ad hoc: Use fact that all 1-dim projections of multivar normal are normal

- take (random) projections $\underline{a}_1, \dots, \underline{a}_m$ and look at qqplots of $\{\underline{a}_k^T \underline{X}_i\}$ for $k=1, \dots, m$

Application: Graphical models

Idea: Describe dependence structure of \underline{X} using an (undirected) graph.

Example: $p=5$



$$\underline{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_5 \end{pmatrix} \sim N_5(\mu, C)$$

X_1 depends on X_2 & X_3

X_2 & X_3 are conditionally independent given X_1 ,

X_4 & X_5 are dependent

(X_4, X_5) & (X_1, X_2, X_3) are independent

Question:

Given C , how to extract this graph structure?

$$\text{Partition: } \underline{X} = \begin{pmatrix} \underline{X}_1 \\ \underline{X}_2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

$$C = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix} \quad \begin{matrix} C_{11} \leftarrow r \times r \\ C_{12} \leftarrow r \times (p-r) \\ C_{21} \leftarrow (p-r) \times r \\ C_{22} \leftarrow (p-r) \times (p-r) \end{matrix}$$

$$\text{and } K = C^{-1} = \begin{pmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{pmatrix} \quad \begin{matrix} K_{11} \leftarrow (p-r) \times (p-r) \\ K_{12} \leftarrow (p-r) \times r \\ K_{21} \leftarrow r \times (p-r) \\ K_{22} \leftarrow r \times r \end{matrix}$$

↳ Called "concentration matrix"

$$X_1 | X_2 = x_2 \sim N_r(\mu_{1|2}, C_{1|2})$$

where $\mu_{1|2} = \mu_1 + C_{12} C_{22}^{-1} (x_2 - \mu_2)$ conditional mean
 $= \mu_1 - K_{11}^{-1} K_{12} (x_2 - \mu_2)$

$$C_{1|2} = C_{11} - C_{12} C_{22}^{-1} C_{21}$$
 conditional covariance
 $= K_{11}^{-1}$

Take $r=2$, look at dist'n of 2 variables conditional on remaining $p-2$.

covariance = K_{11}^{-1}
 $\rightarrow 2 \times 2$ matrix

What do we need for 2 variables to be independent given the remaining $p-2$?

\Rightarrow off-diagonal elements of $K_{11}^{-1} = 0$. \Leftrightarrow off-diagonal elements of $K_{11} = 0$



How to draw graph?

- all dependence information necessary to draw graph contained in $K = C^{-1}$.

$$K = \begin{pmatrix} k_{11} & \dots & k_{1p} \\ \vdots & \ddots & \vdots \\ k_{p1} & \dots & k_{pp} \end{pmatrix}$$

- If $k_{ij} \neq 0$ then there is an edge between vars i & j

- If $k_{ij} = 0$... no edge

Example: Exam marks in 5 subjects

- mechanics
 - vectors
 - algebra
 - analysis
 - statistics
- 88 students

$$\hat{R} = \begin{pmatrix} 1 & 0.55 & 0.55 & 0.41 & 0.39 \\ & 1 & 0.61 & 0.49 & 0.44 \\ & & 1 & 0.71 & 0.66 \\ & & & 1 & 0.61 \\ & & & & 1 \end{pmatrix}$$

Look at $\hat{R} = \hat{R}^{-1}$