

Training Exam, Multivariate Analysis

Shaobo Jin

1. (4pt) Let

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \mathbf{X}_3 \end{bmatrix} \sim N \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \\ \boldsymbol{\mu}_3 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} & \boldsymbol{\Sigma}_{13} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} & \boldsymbol{\Sigma}_{23} \\ \boldsymbol{\Sigma}_{13}^T & \boldsymbol{\Sigma}_{23}^T & \boldsymbol{\Sigma}_{33} \end{bmatrix} \right),$$

and

$$\mathbf{A} = \begin{bmatrix} \mathbf{I}_{q \times q} & -\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1} \\ \mathbf{0} & \mathbf{I}_{r \times r} \end{bmatrix}.$$

(a) Find the joint distribution of $\mathbf{A} \left(\begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \right)$.

(b) Find the conditional distribution of $\begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_3 \end{bmatrix} \mid \mathbf{X}_2 = \mathbf{a}$.

2. (4p) Let X_{jk} be the response to the k th treatment on the j th unit, where $k = 1, 2, \dots, q$ and $j = 1, 2, \dots, n$. Let

$$\mathbf{X}_j = \begin{bmatrix} X_{j1} \\ X_{j2} \\ \vdots \\ X_{jq} \end{bmatrix}.$$

Suppose that $\mathbf{X}_1, \dots, \mathbf{X}_n$ is a random sample from an $N_q(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ population. We want to test a linear combination of $\boldsymbol{\mu}$ as

$$H_0 : \mathbf{C}\boldsymbol{\mu} = \mathbf{0} \quad \text{versus} \quad H_1 : \mathbf{C}\boldsymbol{\mu} \neq \mathbf{0},$$

where \mathbf{C} is a matrix of constants.

(a) Let

$$\begin{aligned} \bar{\mathbf{X}} &= \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j, \\ \mathbf{S} &= \frac{1}{n-1} \sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}})(\mathbf{X}_j - \bar{\mathbf{X}})^T. \end{aligned}$$

Find the distribution of $\mathbf{C}\bar{\mathbf{X}}$ and $(n-1)\mathbf{C}\mathbf{S}\mathbf{C}^T$.

- (b) Find the distribution of $T^2 = n(\mathbf{C}\bar{\mathbf{X}} - \mathbf{C}\boldsymbol{\mu})^T (\mathbf{CSC}^T)^{-1} (\mathbf{C}\bar{\mathbf{X}} - \mathbf{C}\boldsymbol{\mu})$.
3. (4p) Consider the OLS estimator of the multivariate regression coefficient $\hat{\boldsymbol{\beta}} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Y}$. The predicted value is $\hat{\mathbf{Y}} = \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Y}$, and the residual matrix is $\hat{\mathbf{E}} = \mathbf{Y} - \hat{\mathbf{Y}}$.
- (a) Show that the predicted value is $\hat{\mathbf{Y}}_{(i)}$ is perpendicular to the residual $\hat{\mathbf{E}}_{(k)}$, where $\hat{\mathbf{Y}}_{(i)}$ is the i th column of $\hat{\mathbf{Y}}$ and $\hat{\mathbf{E}}_{(k)}$ is the k th column of $\hat{\mathbf{E}}$.
- (b) Show that $\hat{\mathbf{E}}_{(k)}$ is perpendicular to the columns of \mathbf{Z} .
4. (15pt) Suppose that we have a dataset with the following variables *Station*, *PM25*, *SO2*, *NO2*, and *CO*. There are 12 stations in total.
- (a) Suppose that we have perform the following test in R

```
HotellingsT2Test(x = Station1, mu = c(60, 20, 60, 1000), test = "f")

##
## Hotelling's one sample T2-test
##
## data: Station1
## T.2 = 19.992, df1 = 4, df2 = 199, p-value = 7.472e-14
## alternative hypothesis: true location is not equal to c(60,20,60,1000)
```

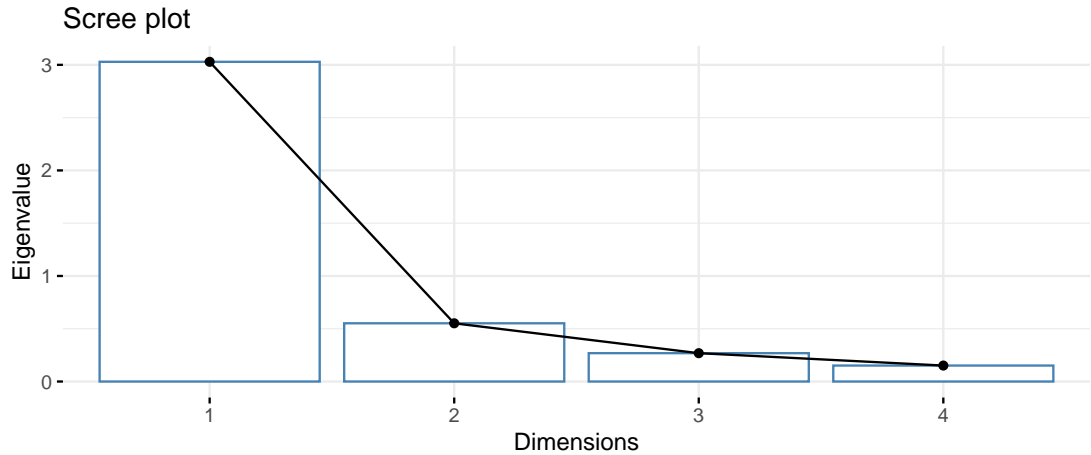
Which test problem has been considered here?

- (b) Given the degrees of freedoms of the test produced by the function HotellingsT2Test(), how many observations do you have in the first station?
- (c) A statistician has computed some confidence intervals for the mean values of these variables of the first station. The R code of computing those confidence interval are given below.

```
Xbar <- colMeans(Station1)
S <- cov(Station1)
c(Xbar[1] - sqrt(S[1, 1] / n1) * qt(0.975, n1 - 1),
  Xbar[1] + sqrt(S[1, 1] / n1) * qt(0.975, n1 - 1))
c(Xbar[2] - sqrt(S[2, 2] / n1) * qt(0.975, n1 - 1),
  Xbar[2] + sqrt(S[2, 2] / n1) * qt(0.975, n1 - 1))
c(Xbar[3] - sqrt(S[3, 3] / n1) * qt(0.975, n1 - 1),
  Xbar[3] + sqrt(S[3, 3] / n1) * qt(0.975, n1 - 1))
c(Xbar[4] - sqrt(S[4, 4] / n1) * qt(0.975, n1 - 1),
  Xbar[4] + sqrt(S[4, 4] / n1) * qt(0.975, n1 - 1))
```

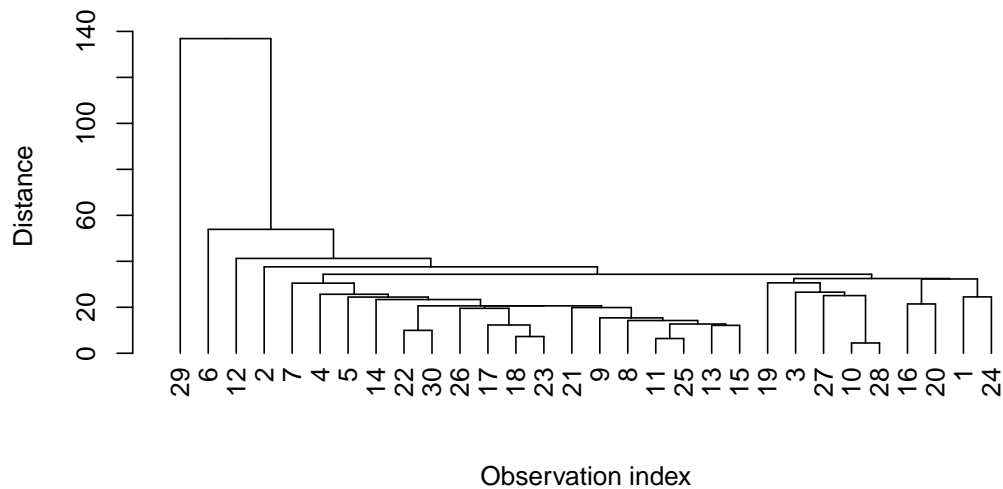
Let $n1$ be the sample size of the first station. Can these confidence intervals give you the correct confidence level? State also the reason.

- (d) Write down the MANOVA model equation that tests whether different stations have the same $PM_{2.5}$, SO_2 , NO_2 , and CO . What is the null hypothesis in terms of your MANOVA model?
- (e) Specify the assumptions in order to apply MANOVA.
- (f) A statistician has applied PCA to the correlation matrix of the variables $PM_{2.5}$, SO_2 , NO_2 , and CO . The following scree plot is obtained.



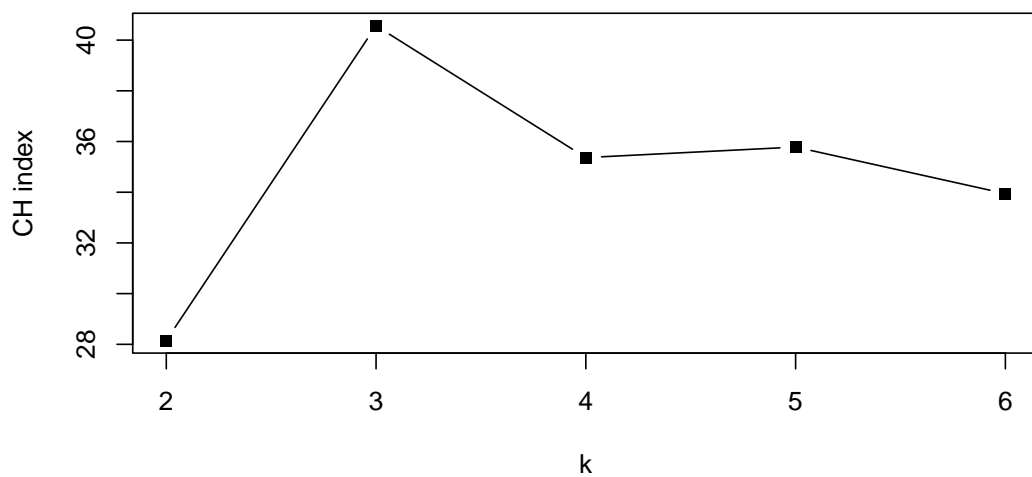
How many principal components would you like to choose? State also the reason.

- (g) Let the eigenvalues of the correlation matrix be 3, 0.6, 0.3, and 0.1. How much total variation is explained by the first principal component?
 - (h) The statistician performed orthogonal factor rotation. Are the communalities after rotation the same as the communalities before rotation?
 - (i) Formulate what types of research questions principal component analysis and factor analysis can answer from this data set, respectively.
5. (5pt) Suppose that we have a dataset with the following variables. The response variable is *heart disease* with two classes: with heart disease (1) or without heart disease (0). The covariates are *age* (A), *resting blood pressure* (B), and *cholesterol level* (C), and *maximum heart rate* (H).
- (a) Suppose that the response variable *heart disease* is actually missing. A hierarchical clustering analysis with single linkage has been performed. The resulting dendrogram is shown below.



If three clusters are desired, how would you cluster the observations based on the dendrogram?

- (b) We also want to perform a k-means clustering analysis. To determine the number of clusters, we have obtained the following plot of the CH index.



How many clusters would like to choose?

- (c) A k-means clustering has been done with the following R code. Can you identify any potential issues? Name at least two, and motivate your answer.

```
kmeans(Data, centers = 3, nstart = 1)
```

6. (4pt) Let $Z \sim \text{Bernoulli}(\phi)$ be a random variable that indicates the class, where $Z = 1$ or 0 :

$$\mathbf{X} \mid Z = 1 \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}), \quad \mathbf{X} \mid Z = 0 \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}).$$

Here $\boldsymbol{\Sigma}$ is known, but ϕ , $\boldsymbol{\mu}_1$, and $\boldsymbol{\mu}_2$ are unknown. Derive a naive Bayes classifier.

7. (4pt) Suppose that we have observed a univariate random sample x_1, x_2, \dots, x_n from K populations, but we do not observe which population each observation comes from. We assume that

$$\begin{aligned} X \mid Z = k &\sim N_q(\mu_k, \sigma_k^2), \\ P(Z = k) &= p_k, \end{aligned}$$

where σ_k^2 are known. Find the expression of p_k and μ_k that maximizes the conditional expectation computed from the E-step of an EM algorithm.