

**THE UNIVERSITY OF CHICAGO**  
**Booth School of Business**  
Business 41912, Spring Quarter 2018, Mr. Ruey S. Tsay

**Midterm**

**Chicago Booth Honor Code:**

*I pledge my honor that I have not violated the Honor Code during this examination.*

**Signature:**

**Name:**

**ID:**

**Notes:**

1. The exam has two parts. The first part consists of some simple questions whereas the second part focuses on data analysis. The total time is 180 minutes with the first 60 minutes (maximum) allocated to the first part. You can immediately start the second part of the exam once you hand in the solutions of the first part.
2. Turn off cell phones. No communication is allowed during the exam. Except for downloading data of part 2, no Internet connection is allowed.
3. You may consult textbooks during the second part of the exam.
4. Write your answers in a bluebook. Mark the solution clearly.
5. All tests are based on the 5% significance level.
6. The transpose of the matrix  $\mathbf{A}$  is  $\mathbf{A}'$ . Also, the  $(i, j)$ th element of the matrix  $\mathbf{A}$  is  $a_{ij}$ .

**Part One.** Basic concepts.

**Problem A.** (10 points) Suppose that  $\mathbf{X} = (X_1, X_2, X_3)'$  follows a 3-dimensional normal distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$  given by

$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} 4 & -3 & 0 \\ -3 & 6 & 0 \\ 0 & 0 & 5 \end{bmatrix}.$$

Answer the following questions:

1. Find the joint distribution of  $\mathbf{Z} = (X_1, X_2)'$ .
2. Are  $X_1$  and  $X_3$  independent? Why?
3. Are  $3X_1 - X_2$  and  $X_3$  independent? Why?
4. Are  $(X_1, X_3)$  and  $X_2$  independent? Why?
5. What is the conditional distribution of  $X_1|X_2 = 1$ ?

**Problem B.** (10 points) Suppose that  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are random samples from a  $p$ -dimensional multivariate distribution with mean  $\boldsymbol{\mu}$  and positive-definite covariance matrix  $\boldsymbol{\Sigma}$ . Let  $\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'$  be the sample covariance matrix and  $\bar{\mathbf{X}}$  be the sample mean. Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be another random sample from the same distribution. Denote the sample mean and sample covariance of this 2nd sample by  $\bar{\mathbf{Y}}$  and  $\mathbf{S}_2$ , respectively. The two random samples are independent.

1. Obtain the mean and covariance matrix of  $\mathbf{X}_1 - \mathbf{X}_2$ ?
2. Show that  $E(\bar{\mathbf{X}}) = \boldsymbol{\mu}$ .
3. Show that  $\text{var}(\bar{\mathbf{X}}) = \frac{1}{n}\boldsymbol{\Sigma}$ .
4. Show that  $E(\mathbf{S}) = \boldsymbol{\Sigma}$ .
5. What is the limiting distribution of  $n(\bar{\mathbf{X}} - \bar{\mathbf{Y}})'(\mathbf{S} + \mathbf{S}_2)^{-1}(\bar{\mathbf{X}} - \bar{\mathbf{Y}})$  as  $n \rightarrow \infty$ ?

**Problem C:** (12 points) Consider the multivariate multiple linear regression (mmr) model,

$$\mathbf{Y}_i' = \mathbf{Z}_i' \boldsymbol{\beta} + \boldsymbol{\epsilon}_i', \quad i = 1, \dots, n,$$

where  $\mathbf{Y}_i$  is an  $m$ -dimensional dependent vector,  $\mathbf{Z}_i = (1, Z_{1i}, \dots, Z_{pi})'$  is a  $(p+1)$  dimensional vector and  $\boldsymbol{\epsilon}_i = (\epsilon_{1i}, \dots, \epsilon_{mi})'$  is an  $m$ -dimensional random noise with mean zero and positive-definite covariance matrix  $\boldsymbol{\Sigma}$ . In matrix form, the model can be written as

$$\mathbf{Y}_{n \times m} = \mathbf{Z}_{n \times (p+1)} \boldsymbol{\beta}_{(p+1) \times m} + \boldsymbol{\epsilon}_{n \times m},$$

where  $\mathbf{Z}$  is a full rank design matrix and we assume that  $\boldsymbol{\epsilon}_i$  are uncorrelated. Assume that  $n > (p+1) + m$ . Let  $\widehat{\mathbf{Y}} = \mathbf{Z}\widehat{\boldsymbol{\beta}}$  be the fitted value of the mmr model and  $\widehat{\boldsymbol{\epsilon}} = \mathbf{Y} - \widehat{\mathbf{Y}}$  be the residual matrix, where  $\widehat{\boldsymbol{\beta}}$  is the ordinary least squares estimate of  $\boldsymbol{\beta}$ . Also, let  $\mathbf{H} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$  be the hat matrix of the regression.

1. Let  $\mathbf{u}$  be a  $n$ -dimensional vector of 1. Show that  $\mathbf{u}'\widehat{\boldsymbol{\epsilon}} = \mathbf{0}$ .
2. Show that  $\sum_{i=1}^n h_{ii} = r + 1$ .

3. (True or False)  $h_{ii} > 0$  for all  $i$ .
4. (True or False)  $\text{Cov}(\text{vec}(\hat{\boldsymbol{\beta}})) = \boldsymbol{\Sigma}^{-1} \otimes (\mathbf{Z}'\mathbf{Z})^{-1}$ .
5. (True or False) Residuals  $\hat{\boldsymbol{\epsilon}}_i$  and  $\hat{\boldsymbol{\epsilon}}_j$  are uncorrelated if  $i \neq j$ .
6. (True or False)  $E(\frac{1}{n}\hat{\boldsymbol{\epsilon}}'\hat{\boldsymbol{\epsilon}}) = \boldsymbol{\Sigma}$ .

**Problem D.** (10 points) Let  $\mathbf{x} = (x_1, \dots, x_p)$  be a  $p$ -dimensional random variable that has a continuous distribution with mean zero and positive-definite covariance matrix  $\boldsymbol{\Sigma}$ . Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be a random sample from  $\mathbf{x}$  and denote the data matrix by  $\mathbf{X}$ , which is a  $n \times p$  matrix.

1. Define the first principal component of  $\mathbf{x}$ .
2. Let  $(\lambda_i, \mathbf{e}_i)$  be the  $i$ th eigenvalue-eigenvector pair of  $\boldsymbol{\Sigma}$ , where  $\lambda_1 > \lambda_2 > \dots > \lambda_p$  and  $\mathbf{e}_i'\mathbf{e}_i = 1$ . What is the second principal component of  $\mathbf{x}$ ?
3. Let  $(\hat{\lambda}_i, \hat{\mathbf{e}}_i)$  be the  $i$ th eigenvalue-eigenvector pair of  $\mathbf{S}$ , where  $\hat{\lambda}_1 > \hat{\lambda}_2 > \dots > \hat{\lambda}_p$  and  $\mathbf{S}$  is the sample covariance matrix. True or False that  $\hat{\lambda}_i$  is a consistent estimate of  $\lambda_i$ ?
4. Under what condition on the data matrix  $\mathbf{X}$  that the singular value decomposition of  $\mathbf{X}$  provides a quick principal component analysis for  $\mathbf{x}$ ?
5. Describe a weakness of principal component analysis.

**Problem E.** (8 points) Answer briefly the following questions.

1. Describe two methods for model selection in the linear regression analysis.
2. One can use multivariate analysis of variance to compare the mean vectors of several populations. Describe two assumptions used in such an analysis.
3. Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be a realization of the  $p$ -dimensional time series model  $\mathbf{x}_t = \boldsymbol{\mu} + \mathbf{a}_t - \boldsymbol{\theta}\mathbf{a}_{t-1}$ , where  $\{\mathbf{a}_t\}$  is a sequence of independent and identically distributed Gaussian random vector with mean  $\mathbf{0}$  and positive-definite covariance matrix  $\boldsymbol{\Sigma}$ . Let  $\bar{\mathbf{x}}$  be the sample mean of the data. What is  $E(\bar{\mathbf{x}})$ ?
4. What is  $\text{Cov}(\bar{\mathbf{x}})$ , where  $\bar{\mathbf{x}}$  is defined in the prior question.

## Part Two: Data analysis

**Problem F.** (22 points) Consider the measurements of blood glucose levels on three occasions for 50 women. The dependent variables  $y$ 's represent fasting glucose measurements on the three occasions and the predictors  $x$ 's are glucose measurements 1 hour after sugar intake. The data are available on the course web: `Glucose.DAT`.

1. Find the maximum and minimum of Spearman's rho between  $y$  and  $x$ .
2. Find the maximum and minimum of Kendall's tau between  $y$  and  $x$ .
3. Let  $\Sigma_1$  and  $\Sigma_2$  be the covariance matrices of  $y$ 's and  $x$ 's, respectively. Test the hypothesis  $H_0 : \Sigma_1 = \Sigma_2$  versus  $H_a : \Sigma_1 \neq \Sigma_2$ . Draw a conclusion.
4. Let  $\mu_1$  and  $\mu_2$  be the mean vectors of  $y$ 's and  $x$ 's, respectively. Test the hypothesis  $H_0 : \mu_1 = \mu_2$  versus  $H_a : \mu_1 \neq \mu_2$ . Draw a conclusion.
5. Obtain the least squares estimate of  $\beta$  and the residual covariance matrix of the linear regression

$$\mathbf{Y}_i = \mathbf{X}_i\beta + \epsilon_i,$$

where  $\mathbf{Y} = (y_1, y_2, y_3)'$  and  $\mathbf{X} = (1, x_1, x_2, x_3)'$ .

6. Is the overall regression statistically significant? Why?
7. Test the significance of  $x_1$  given  $x_2$  and  $x_3$ .
8. Test the significance of  $x_3$  given  $x_1$  and  $x_2$ .
9. (3 points) Let  $\mathbf{x}_0 = (1, 100, 100, 100)'$ . Compute simultaneous 95% confidence intervals for the elements of  $E(\mathbf{Y}|\mathbf{X} = \mathbf{x}_0)$ . [You should update the "ama.R" file. The new version of `mmlr` provides the  $(\mathbf{Z}'\mathbf{Z})^{-1}$  in the output.]
10. (3 points) Calculate simultaneous 95% prediction intervals for the elements of  $E(\mathbf{Y}|\mathbf{X} = \mathbf{x}_0)$ .

**Problem G:** (8 points) An experiment involved a  $2 \times 4$  design with 4 replications was conducted to study the properties of bar steel. The two factors are rotational velocity [ $A_1$  (fast) and  $A_2$  (slow)] and lubricants [ $B_i$  for  $i = 1, 2, 3, 4$ ]. The measurements are  $y_1$  = ultimate torque and  $y_2$  = ultimate strain. The data are available on course web `Barsteel.DAT`. Use the Wilk's statistic to answer the following three questions.

1. Is there any interaction between the two factors? Why?
2. Is the effect of rotational velocity significant? Why?
3. Is the effect of lubricants significant? Why?

4. List the assumptions used in the analysis.

**Problem H:** (10 points) Consider the data set in `ProblemH.txt`. It consists of two scalar dependent variables [ $y$  (column 1) and  $y1$  (column 2)] and 300 regressors  $x_1, \dots, x_{300}$  (columns 3 to 302). You may use the following command to create column names for the data and design matrix for Question 3 below.

```
da <- read.table('ProblemH.txt')
y <- da[,1]; y1 <- da[,2]
x <- da[,3:302]
c1 <- paste('x',1:300,sep='')
colnames(x) <- c1
X <- model.matrix(y1~.(x1+x2+x3+x4+x5+x6+x7+x8+x9+x10)^2,data=data.frame(x))
X <- X[,-1] ## remove the intercept from the design matrix.
```

1. Apply LASSO to regress  $y$  on  $x$  and use 10-fold cross-validation to select the penalty parameter. What is the selected penalty parameter? With this choice of  $\lambda$ , identify the predictors of the LASSO regression with absolute coefficient greater than 0.5.
2. Apply elastic-net regress  $y$  on  $x$  with  $\alpha = 0.5$  (in package **glmnet**). Again, use 10-fold cross-validation to select the penalty parameter. What is the selected penalty parameter  $\lambda$ ? Identify the predictors of the resulting regression with absolute coefficient greater than 0.5.
3. (4 points) In addition to the individual predictors, we add all pairwise interactions of  $\{x_1, \dots, x_{10}\}$ . Apply elastic-net to regress  $y1$  on  $X$  with  $\alpha = 0.8$ . Identify the predictors with absolute coefficient greater than 0.5. Note that this problem concerns new dependent variable and new predictors.

**Problem I.** (10 points) The file `m-apca40stocks.txt` contains monthly returns of 40 stocks for two years. The first row of the file contains the codes for companies, and should be removed. Here we have  $p = 40$  and  $n = 24$ .

1. Perform the traditional principal component analysis. Identify the variable that has the maximum loading (in absolute value) of the first two principal components.
2. What are the variances of the first two principal components?
3. What percentage of total variability is explained by the first two principal components?
4. Perform a sparse principal component analysis. For instance, you may use the package **elasticnet** with command `spca` and subcommand `para=c(2,1)`. What percentages of total variability are explained by the first two sparse principal components?
5. Identify the second sparse principal component. What does it mean?