



The
University
Of
Sheffield.

SCHOOL OF MATHEMATICS AND STATISTICS

**Spring Semester
2015–2016**

Dependent Data

3 hours

*Marks will be awarded for your best **five** answers.*

RESTRICTED OPEN BOOK EXAMINATION

Candidates may bring to the examination lecture notes and associated lecture material (but no textbooks) plus a calculator that conforms to University regulations.

There are 100 marks available on the paper.

**Please leave this exam paper on your desk
Do not remove it from the hall**

Registration number from U-Card (9 digits)
to be completed by student

--	--	--	--	--	--	--	--	--

Blank

- 1 This question considers a study reported by Sokal and Rohlf (1981) of air pollution in 41 US cities, between 1969 and 1971. The variables are as follows:

S02 Sulphur dioxide content of air in μg per cubic metre

Temp Average annual temperature in degrees Fahrenheit

Firms Number of manufacturers employing 20 or more workers

Pop Population in thousands in 1970

Wind Average annual wind speed in miles per hour

Rain Average annual rainfall in inches

Raindays Average number of rainy days per year

A principal components analysis was carried out on the data, and an R transcript is given below.

- (i) Explain why `cor=TRUE` is used in the `princomp` command. *(1 mark)*
- (ii) In the principal components analysis, it is decided to use only the first few components. Using an informal graphical technique, how many components would you choose? *(4 marks)*
- (iii) Phoenix is in the top left of the plot of the first two principal components. Explain why Phoenix gets such a low score on the first principal component.
There is an outlier on PC1, with a very high score. Can you suggest some characteristics of this city? *(7 marks)*
- (iv) Explain why Phoenix seems to get such a high score for the second principal component. *(3 marks)*
- (v) Interpret the third principal component, justifying your answer briefly. *(2 marks)*
- (vi) Comment on the fourth and fifth principal components. Is there anything in the plots of these principal components that strikes you? *(3 marks)*

1 (continued)

```

> attach(airpoll)
> airpoll[1:5,]
      S02 Temp Firms Pop Wind Rain Raindays
Phoenix    10 70.3  213 582  6.0  7.05      36
Little Rock 13 61.0   91 132  8.2 48.52     100
San Fransisco 12 56.7  453 716  8.7 20.66     67
Denver     17 51.9  454 515  9.0 12.95     86
Harford    56 49.1  412 158  9.0 43.37    127

> apply(airpoll,2,mean)
      S02    Temp    Firms    Pop    Wind    Rain Raindays
30.05  55.76  463.10  608.61   9.44  36.77  113.90

> airpoll.pca<-princomp(airpoll,cor=TRUE)
> summary(airpoll.pca)
Importance of components:
              Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7
Standard deviation      1.65  1.230  1.181  0.944  0.5889  0.3167  0.15973
Proportion of Variance  0.39  0.216  0.199  0.127  0.0495  0.0143  0.00364
Cumulative Proportion  0.39  0.606  0.805  0.932  0.9820  0.9964  1.00000

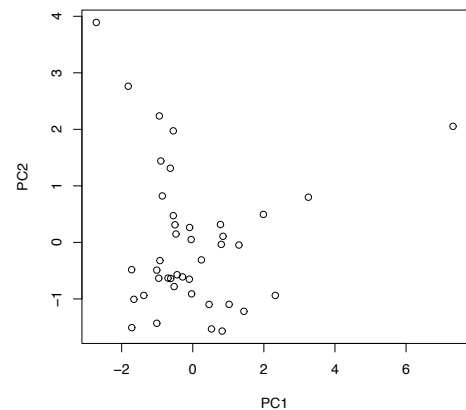
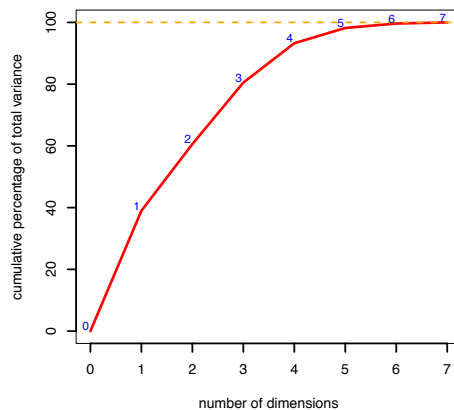
> print(airpoll.pca$loadings,digits=2)

Loadings:
              Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7
S02           0.49          0.40 -0.73 -0.18  0.15
Temp        -0.32          -0.68 -0.19 -0.16 -0.61
Firms        0.54  0.23  -0.27          0.16          -0.75
Pop          0.49  0.28  -0.34 -0.11  0.35          0.65
Wind         0.25          0.31 -0.86 -0.27 -0.15
Rain         -0.63 -0.49 -0.18 -0.16  0.55
Raindays   0.26 -0.68  0.11  0.11  0.44 -0.50

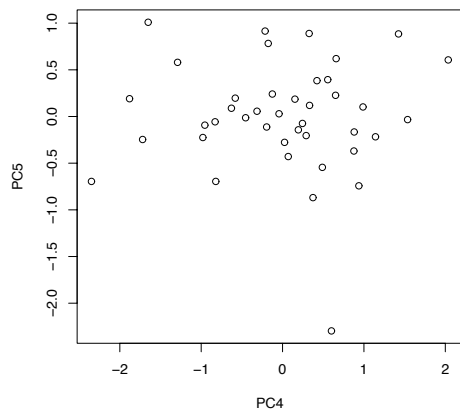
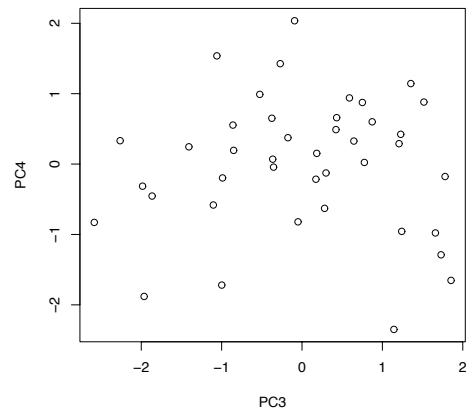
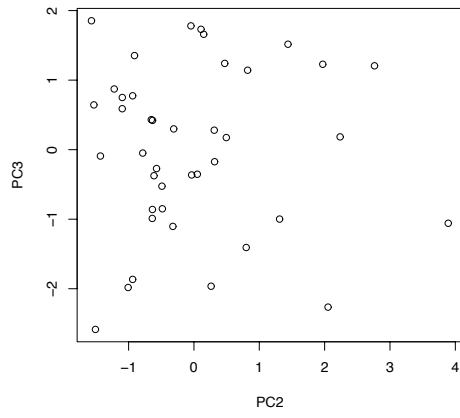
> screeplot(airpoll,cor=TRUE)
> airpoll.pc<-predict(airpoll.pca)
> airpoll.pc[1:5,]
              Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7
Phoenix        -2.7160  3.8914 -1.0584  1.5374 -0.0331 -0.654428  0.09061
Little Rock    -1.7177 -0.4824 -0.8500  0.1945 -0.1421  0.544992 -0.19867
San Fransisco  -0.9390  2.2372  0.1837  0.1528  0.1856  0.295074  0.00714
Denver         -0.5499  1.9719  1.2286  0.4227  0.3837 -0.130106 -0.15799
Harford        0.4604 -1.0973  0.5897  0.9387 -0.7427  0.400483 -0.26019

```

Scree plot of variances



1 (continued)



- 2 (i) Measurements of x_1 (stiffness) and x_2 (bending strength) for a new sample of $n_A = 20$ pieces of the highest grade of timber from a site A were taken. The means were given by $\bar{x}_A = (\bar{x}_1, \bar{x}_2) = (17.6, 81.8)$, while the variance matrix of the sample was $S = \begin{pmatrix} 2.1 & 2.2 \\ 2.2 & 7.4 \end{pmatrix}$, with inverse $S^{-1} = \begin{pmatrix} 0.7 & -0.2 \\ -0.2 & 0.2 \end{pmatrix}$.

The company wishes to compare the new sample with its established population, where the mean stiffness is $\bar{x} = (17.1, 83.0)$.

For the first two parts, R gives `qt(0.975, 19)=2.093`.

- (a) Perform a t -test to test the hypothesis that the mean stiffness in the new sample is equal to 17.1. (1 mark)
- (b) Perform a t -test to test the hypothesis that the mean bending strength in the new sample is equal to 83.0. (1 mark)
- (c) Test the hypothesis that $\bar{x}_A = \bar{x}$. You may use the R output `qf(0.95, 2, 18)=3.555`. (5 marks)
- (d) Discuss briefly the results of parts (a)–(c), bearing in mind that \bar{x}_A seems to be not too far from \bar{x} , illustrating your answer with a sketch of the multivariate confidence region. (You may like to compute the correlation between the two variables to help with your answer.) (2 marks)
- (ii) Suppose that x_1, \dots, x_n are independent observations of a bivariate normal distribution $x \sim N_2(\mu, \Sigma)$, where neither μ nor Σ are known. Write \bar{x} for the sample mean, and $S = \begin{pmatrix} a & c \\ c & b \end{pmatrix}$ for the sample variance.

Recall that the log-likelihood is given by

$$\ell(\mu, \Sigma) = -\frac{1}{2}(n-1)\text{tr}(\Sigma^{-1}S) - \frac{1}{2}n\text{tr}(\Sigma^{-1}(\bar{x}-\mu)(\bar{x}-\mu)') - n\log(2\pi) - \frac{1}{2}n\log|\Sigma|,$$

as $p = 2$, and you may assume that the unrestricted MLEs are given by $\hat{\mu} = \bar{x}$ and $\hat{\Sigma} = \frac{n-1}{n}S$.

- (a) Suppose that $\Sigma = \begin{pmatrix} \alpha & 0 \\ 0 & \beta \end{pmatrix}$ is diagonal, and $\mu = \bar{x}$. By explicitly computing $\Sigma^{-1}S$, show that

$$\ell(\mu, \Sigma) = -\frac{1}{2}(n-1)(\alpha^{-1}a + \beta^{-1}b) - n\log(2\pi) - \frac{1}{2}n\log(\alpha\beta).$$

By computing $\frac{\partial \ell}{\partial \alpha}$, deduce that under the constraint that Σ is diagonal, its MLE is

$$\hat{\Sigma} = \frac{n-1}{n}\text{diag}(S) = \frac{n-1}{n} \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}. \quad (5 \text{ marks})$$

- (b) Develop a likelihood-ratio test to test the null hypothesis that Σ is diagonal, using Wilks's Theorem. (6 marks)

- 3 A study was performed in Malaysia (1985) relating to the rate of settling of red blood cells out of suspension in blood plasma. This rate increases if the level of certain proteins in the blood plasma rise, and two plasma proteins, whose levels are denoted F and G in what follows, were considered. Subjects were classed as healthy or unhealthy depending on the value of this rate of settling.

The mean values of F and G of 26 healthy subjects were 4.66 and 61.81 respectively. The corresponding values of 6 subjects classified as not healthy were 5.97 and 66.88. The pooled within group variances of the two measurements were 1.022 and 63.067, with covariance -0.316 .

A key aim of the study was to determine the strength of any relationship between the levels of the two plasma proteins and the rate of settling.

- (i) Estimate Fisher's linear discriminant function for classifying a subject as healthy or not on the basis of measurements of the levels of plasma proteins F and G . *(7 marks)*
- (ii) Informal investigations suggest that the data for each group is reasonably well approximated by a bivariate normal distribution, and, further, that the variance matrices for both groups appear to be very similar, so that they may be assumed to be the same. Using your answer to part (i) to classify subjects as healthy or unhealthy, estimate the probability of misclassifying a randomly selected healthy subject as unhealthy. You may leave your answer in the form $\Phi(z)$ for some number z which you should determine. *(5 marks)*
- (iii) Suppose that for a particular subject, the G sample is contaminated, and an assessment is made of the subject purely based on the level of the plasma protein F . What value of F should be used as a lower limit to ensure that the probability of missing an unhealthy subject is the same as that using the rule determined in part (i)? *(5 marks)*
- (iv) What proportion of healthy subjects will be falsely diagnosed as unhealthy by the rule in (iii)? (Again, you may give your answer in a form involving $\Phi(z)$ for some z .) *(3 marks)*

- 4 (i) A model is to be fitted to a time series of length 81. Values of the sample autocorrelation function (ACF) and sample partial ACF (PACF) are tabulated below.

Lag (h)	1	2	3	4	5
ACF (r_h)	*	0.7	0.05	0.02	0.01
PACF ($\hat{a}_h^{(h)}$)	0.9	**	0.4	-0.15	0.10

- (a) Find the omitted values (* and **). *(3 marks)*
- (b) Check whether the time series is stationary. *(1 mark)*
- (c) Test whether the time series is consistent with white noise. *(2 marks)*
- (d) Test whether the time series is consistent with MA(1) and MA(2) moving average models. *(5 marks)*
- (e) Test whether the time series is consistent with AR(1), AR(2), AR(3) and AR(4) autoregressive models. *(3 marks)*
- (f) Giving your reason, state which of the models in (d) and (e) you prefer for these data. *(2 marks)*
- (ii) Consider the time series model

$$y_t = 2 \sin\left(\frac{\pi}{t}\right) + \epsilon_t, \quad t = 1, 2, \dots,$$

where ϵ_t follows a white noise process with variance 10.

- (a) Show that y_t is non-stationary process. *(2 marks)*
- (b) Define an appropriate transformation of y_t to result in a stationary time series model. Justify your choice. *(2 marks)*

- 5 Consider that y_t is generated by an ARMA(1,1) model

$$y_t = \alpha y_{t-1} + \epsilon_t + \beta \epsilon_{t-1},$$

where α, β are the AR and MA coefficients and ϵ_t is a Gaussian white noise with variance σ^2 .

- (i) Write down the likelihood and the log-likelihood functions of the parameters α, β and σ^2 , based on a collection of observations $y_{1:n} = (y_1, y_2, \dots, y_n)$.
(6 marks)
- (ii) Using conditional least squares,
- (a) derive the partial derivatives of the conditional log-likelihood with respect to α and β ;
(6 marks)
- (b) using part (a) show that the maximum likelihood estimates of α and β are

$$\hat{\alpha} = \frac{\sum_{t=2}^n y_t y_{t-1} \sum_{t=2}^n \epsilon_{t-1}^2 - \sum_{t=2}^n y_t \epsilon_{t-1} \sum_{t=2}^n y_{t-1} \epsilon_{t-1}}{\sum_{t=2}^n y_{t-1}^2 \sum_{t=2}^n \epsilon_{t-1}^2 - (\sum_{t=2}^n y_{t-1} \epsilon_{t-1})^2}$$

$$\hat{\beta} = \frac{\sum_{t=2}^n y_{t-1}^2 \sum_{t=2}^n y_t \epsilon_{t-1} - \sum_{t=2}^n y_{t-1} \epsilon_{t-1} \sum_{t=2}^n y_t y_{t-1}}{\sum_{t=2}^n y_{t-1}^2 \sum_{t=2}^n \epsilon_{t-1}^2 - (\sum_{t=2}^n y_{t-1} \epsilon_{t-1})^2}.$$

(8 marks)

- 6 Consider the ARMA(1,1) model for the time series y_t :

$$y_t = 0.2y_{t-1} + \epsilon_t + \epsilon_{t-1}, \quad (1)$$

where ϵ_t follows Gaussian white noise with variance $\sigma^2 = 1$.

Define the state vector

$$\beta_t = \begin{bmatrix} y_t \\ \epsilon_{t+1} \\ \epsilon_t \end{bmatrix}.$$

Using this state vector write down model (1) in state space form, i.e.

$$\begin{aligned} y_t &= x\beta_t + \eta_t \\ \beta_t &= F\beta_{t-1} + \zeta_t \end{aligned}$$

and determine x , F , η_t , ζ_t and the variances of η_t and ζ_t . **(3 marks)**

- (i) If the posterior distribution of the state β_1 , given $y_1 = 2$ is

$$\beta_1 | \{y_1 = 2\} \sim N \left\{ \begin{bmatrix} 3 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right\},$$

then find the one-step forecast distribution of y_2 . **(6 marks)**

- (ii) Using the result in (i) obtain a 95% predictive interval for y_2 . **(2 marks)**

- (iii) At time $t = 2$ the observation y_2 is observed to be 3. Perform the Kalman filter iteration for $t = 2$ and obtain the posterior distribution of

$$\beta_2 | \{y_2 = 3\}.$$

(9 marks)

End of Question Paper