



The
University
Of
Sheffield.

SCHOOL OF MATHEMATICS AND STATISTICS

**Spring Semester
2016–2017**

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 (i) In a report from 1973, Weber reports on protein consumption in 25 European countries for nine food groups: Red meat, White meat, Eggs, Milk, Fish, Cereals, Starchy foods, Pulses and nuts, Fruits and vegetables. The R analysis of this data set follows on the next two pages.

(a) What is the correlation between red meat consumption and white meat consumption? *(2 marks)*

(b) Why is the option `cor=TRUE` included in the `princomp` command? What tends to happen if it is omitted? *(2 marks)*

(c) The session has been edited so that only 6 components are listed. But how many would R have computed? *(1 mark)*

(d) With the aid of an informal graphical technique, how many principal components would you retain in your study? Justify your answer. *(2 marks)*

(e) Suggest characteristics of countries with high scores on the first principal component. *(2 marks)*

(f) Albania has a very low consumption of fish. How might one expect to see that reflected in the PCA plots? Which country would you expect to have the highest consumption of fish? *(2 marks)*

(g) Give an interpretation of the third and fourth principal components, and comment on the protein consumption of Finland with particular reference to these components. *(3 marks)*

(ii) Find the principal components and the proportion of total variance explained by each when the covariance matrix is

$$S = \begin{pmatrix} 1 & \rho & 0 \\ \rho & 1 & \rho \\ 0 & \rho & 1 \end{pmatrix}.$$

(4 marks)

(iii) Explain briefly why multidimensional scaling might be regarded as a generalisation of principal components analysis. *(2 marks)*

1 (continued)

```

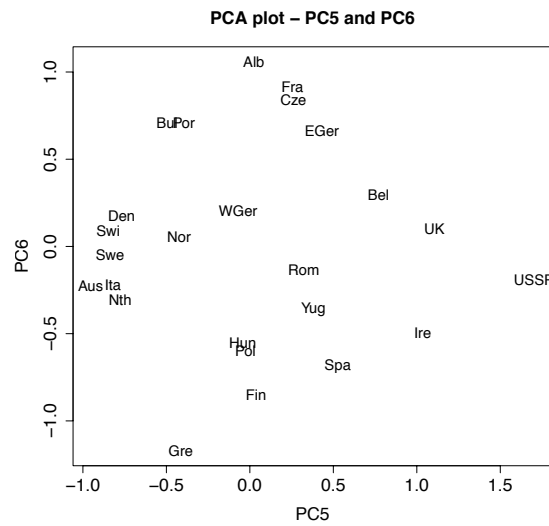
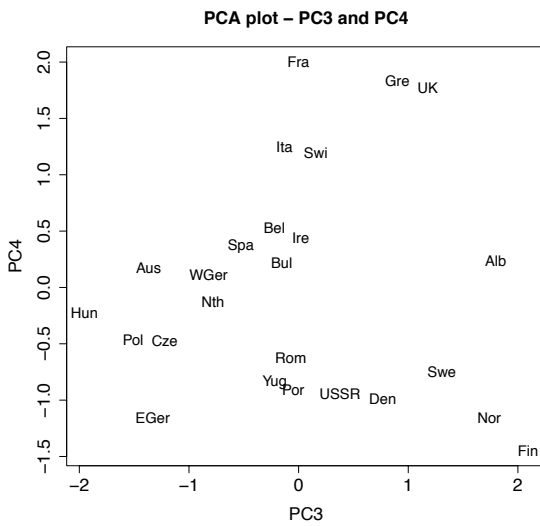
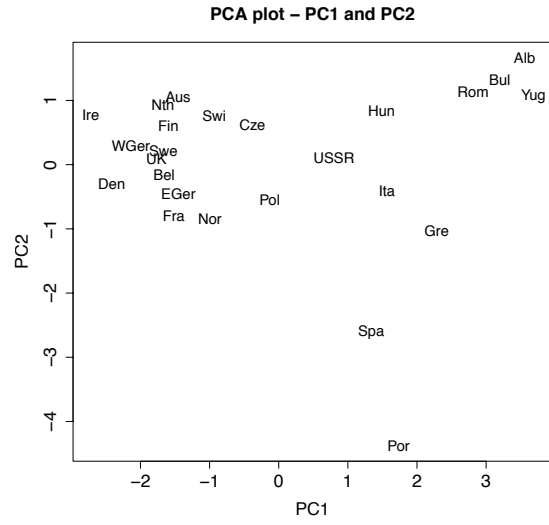
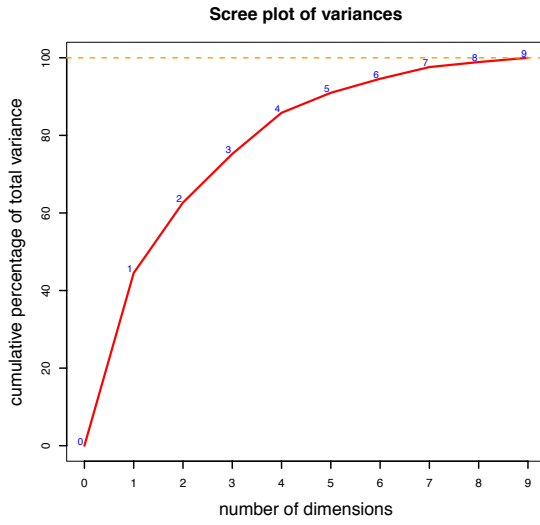
> attach(protein)
> protein[1:2,]
  RedMeat WhiteMeat Eggs Milk Fish Cereals Starch Nuts Fr.Veg
Alb    10.1      1.4 0.5 8.9 0.2      42   0.6 5.5  1.7
Aus     8.9     14.0 4.3 19.9 2.1      28   3.6 1.3  4.3
> apply(protein,2,mean)
  RedMeat WhiteMeat Eggs Milk Fish Cereals Starch Nuts Fr.Veg
      9.8      7.9  2.9 17.1  4.3    32.2    4.3  3.1  4.1
> var(protein)
      RedMeat WhiteMeat Eggs Milk Fish Cereals Starch Nuts Fr.Veg
RedMeat    11.20     1.89  2.191 12.0  0.69 -18.36  0.74 -2.32 -0.448
WhiteMeat   1.89    13.65  2.561  7.4 -2.94 -16.78  1.89 -4.66 -0.409
Eggs        2.19     2.56  1.249  4.6  0.25  -8.74  0.83 -1.24 -0.092
Milk       11.96     7.39  4.570 50.5  3.33 -46.22  2.58 -8.76 -5.234
Fish        0.69    -2.94  0.249  3.3 11.58 -19.58  2.25 -0.99  1.634
Cereals    -18.36   -16.78 -8.738 -46.2 -19.58 120.45 -9.56 14.19  0.922
Starch      0.74     1.89  0.826  2.6  2.25  -9.56  2.67 -1.54  0.249
Nuts       -2.32    -4.66 -1.242  -8.8 -0.99  14.19 -1.54  3.94  1.343
Fr.Veg     -0.45    -0.41 -0.092  -5.2  1.63  0.92  0.25  1.34  3.254

> pr.pca<-princomp(protein,cor=TRUE)
> loadings(pr.pca)
Loadings:
      Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6
RedMeat -0.303      0.298 0.646 0.322 0.460
WhiteMeat -0.311 0.237 -0.624      -0.300 0.121
Eggs      -0.427      -0.182 0.313      -0.361
Milk      -0.378 0.185 0.386      -0.200 -0.618
Fish      -0.136 -0.647 0.321 -0.216 -0.290 0.130
Cereals   0.438 0.233      0.238
Starch    -0.297 -0.353 -0.243 -0.337 0.736 -0.148
Nuts      0.420 -0.143      0.330 0.151 -0.447
Fr.Veg    0.110 -0.536 -0.408 0.462 -0.234 -0.119

> summary(pr.pca)
Importance of components:
      Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6
Standard deviation      2.00  1.28  1.06  0.98  0.681  0.570
Proportion of Variance  0.45  0.18  0.13  0.11  0.052  0.036
Cumulative Proportion  0.45  0.63  0.75  0.86  0.910  0.946

```

1 (continued)



2 Johnson and Wichern (2007) report on a study of 45 female hook-billed kites. The tail length and wing length are measured in centimetres for each bird. The mean tail length is 19.36cm, while the mean wing length is 27.98cm. The variance matrix is given by $S = \begin{pmatrix} 1.207 & 1.223 \\ 1.223 & 2.085 \end{pmatrix}$ (with the variable tail length appearing first, so that its variance is 1.207), with inverse $S^{-1} = \begin{pmatrix} 2.044 & -1.199 \\ -1.199 & 1.183 \end{pmatrix}$.

You may use the R output `qf(0.95, 2, 43)=3.214` and `qt(0.975, 44)=2.015`.

(i) A more extended study of male kites has shown that the mean tail length for male birds is 19.75cm, and the mean wing length is 28.45cm.

(a) Compute a 95% confidence interval for the mean tail length of female kites. Can we reject the hypothesis that the mean tail length of female kites is 19.75cm at the 5% level of significance? **(2 marks)**

(b) Similarly, can we reject the hypothesis that the mean wing length for female kites is 28.45cm at the 5% level of significance? **(2 marks)**

(c) Test the multivariate hypothesis that the mean of the pairs of measurement $\bar{x} = (19.36, 27.98)$ for female kites is equal to $\mu_0 = (19.75, 28.45)$ at the 95% level. **(6 marks)**

(d) Discuss briefly the results of parts (a)–(c). Illustrate your answer with a sketch of the multivariate confidence region. **(3 marks)**

(ii) Suppose that x_1, \dots, x_n are independent observations of a bivariate normal distribution. We work with the principal components of the generated data, so that we suppose that the sample variance matrix is of the form $S = \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}$, with uncorrelated variables. Using these variables, we write \bar{x} for the sample mean, and μ and Σ for the parameters of the distribution, so that our data comes from an $N_2(\mu, \Sigma)$ distribution.

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

We wish to test the hypothesis that $\det \Sigma = \delta$, some fixed positive number. You may assume that, under $H_0 : \det \Sigma = \delta$, we have $\hat{\mu} = \bar{x}$, and that

$$\hat{\Sigma} = \sqrt{\frac{\delta}{ab}} \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}.$$

Hence give the likelihood ratio test for the null hypothesis that $\det \Sigma = \delta$, using Wilks's Theorem. **(7 marks)**

3 Campbell and Mahon (1974) took various measurements of certain species of rock crab of genus *Leptograpsus*. Observations 1–50 were male, and observations 51–100 were female. Of interest in this question is whether we can predict the sex of the crab from the two dimension measurements, given in millimetres, as:

FL frontal lobe
 RW rear width

Each line of the data set consists of 3 fields: firstly, the sex of the crab, and then the remaining 2 dimension measurements in the order above.

It seems that the dimensions of the two groups appear to be distributed as bivariate normal distributions with a similar variance matrix, which we take to be the same as the variance matrix in the R output which follows on the next page.

(i) Compute Fisher’s linear discriminant function for this data set. *(5 marks)*

(ii) Using the output from the previous part, classify the sex of a crab whose measurements are FL = 13 and RW = 10. *(2 marks)*

(iii) By working out the probability of misclassification under the assumption that the two groups are distributed as bivariate normal distributions with the variance matrix in the R output, does Fisher’s linear discriminant function classify the data set better or worse than expected? You are given that $\text{pnorm}(-1)=0.159$. *(8 marks)*

(iv) Suppose that two bivariate normal distributions have means $\mu_1 = \begin{pmatrix} -1 \\ 0 \end{pmatrix}$ and $\mu_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, and that both have variance matrix $\begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$. Find the decision boundary (i.e., the set of points for which the probability densities agree) as a function of ρ . *(5 marks)*

3 (continued)

```

> crab[1,]
  sex FL RW
1  M 8.1 6.7
> print(S<-var(crab[,-1]))
  FL RW
FL 9.12 6.18
RW 6.18 5.20
> print(Sinv<-solve(S))
  FL RW
FL 0.566 -0.673
RW -0.673 0.993
> print(crab.lda<-lda(sex~FL+RW))
Group means:
  FL RW
F 13.3 12.1
M 14.8 11.7

Coefficients of linear discriminants:
  LD1
FL 1.21
RW -1.52
> crab.pred<-predict(crab.lda,data.frame(cbind(FL,RW)))
> crab.pred$class
 [1] M M M M M M F M M F F F M M M M M M M M M M M
 [26] M M M M M M M M M M M M M M M M F M M M M M M
 [51] F F F F M F F F F F F F F F F F F F F F F F F
 [76] F F F F F F F F F F F F F F F F F F F F F F F
Levels: F M
> table(sex,crab.pred$class)
sex F M
  F 49 1
  M 5 45

```


4 (i) Explain briefly why the definition of the sample autocorrelation function is not appropriate for non-stationary time series. (2 marks)

(ii) (a) For the following time series data

time (t)	response (y_t)
1	1
2	3
3	4
4	10
5	2
6	2
7	5
8	11

calculate 4-span moving averages for time points $t = 3, 4, 5, 6$. (4 marks)

(b) Consider the time series

$$y_t = \alpha + \beta t,$$

for some known α and β and for $t = 1, 2, 3, \dots$. Show that the 4-span moving average of $\{y_t\}$ at time t is exactly equal to the value of y_t . (4 marks)

(c) Now consider the time series

$$x_t = y_t + \epsilon_t,$$

where y_t is the time series in part (b) above and ϵ_t is white noise with variance 2. Show that the de-trended time series defined by

$$d_t = x_t - \hat{y}_t,$$

where \hat{y}_t is the 4-span moving average of y_t , is a weakly stationary time series. (2 marks)

(iii) Let $\{Z_t\}$ be a sequence of independent identically distributed random variables, so that Z_t follows a normal distribution with zero mean and variance 1.

Define the process

$$y_t = \begin{cases} Z_t, & \text{if } t \text{ is even} \\ \frac{Z_t^2 - 1}{\sqrt{2}}, & \text{if } t \text{ is odd.} \end{cases}$$

(a) Show that $\{y_t\}$ is a white noise process $WN(0, 1)$. (6 marks)

(b) Show $\{y_t\}$ is not an identically distributed sequence. (2 marks)

HINT: If a random variable X follows the chi-square distribution with ν degrees of freedom $X \sim \chi_\nu^2$, then $E(X) = \nu$ and $\text{Var}(X) = 2\nu$.

5 Suppose that observations y_1, y_2, \dots, y_n are generated from the autoregressive (AR) model of order one:

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

where α is the AR parameter and ϵ_t is a Gaussian white noise with variance σ^2 .

(i) Write down the likelihood function $L(\alpha, \sigma^2; y_{1:n})$ and the log-likelihood function $\ell(\alpha, \sigma^2; y_{1:n})$ of the parameters α and σ^2 , based on observation $y_{1:n} = \{y_1, y_2, \dots, y_n\}$.
(3 marks)

(ii) Using the approximation $\log(1 + x) \approx x$, for some x , with $|x| < 1$, show that the log-likelihood of part (i) can be approximated as

$$\ell(\alpha, \sigma^2; y_{1:n}) \approx -\frac{n}{2} \log(2\pi\sigma^2) - \frac{y_1^2}{2\sigma^2} + \left(\frac{y_1^2 - \sigma^2}{2\sigma^2} \right) \alpha^2 - \frac{1}{2\sigma^2} \sum_{t=2}^n (y_t - \alpha y_{t-1})^2.$$

(7 marks)

(iii) Using (ii) and adopting unconditional least squares, show that the approximate likelihood estimates of α and σ^2 satisfy

$$\hat{\alpha} = \frac{\sum_{t=2}^n y_t y_{t-1}}{\sum_{t=3}^n y_{t-1}^2 + \hat{\sigma}^2} \quad \text{and} \quad \hat{\sigma}^2 = \frac{(1 - \hat{\alpha}^2) y_1^2 + \sum_{t=2}^n (y_t - \hat{\alpha} y_{t-1})^2}{n}.$$

(8 marks)

(iv) If $\sigma^2 = y_1^2$, show that the maximum likelihood of α based on unconditional least squares is approximately equal to the maximum likelihood of α using conditional least squares.
(2 marks)

6 An airliner on a long-haul journey is set to fly at a pre-specified steady speed of v km/hr, but atmospheric conditions sometimes speed it up and sometimes slow it down. A simple model for its true position x_t along the flight path at time t hours supposes that

$$x_{t+1} = x_t + v + \eta_t,$$

where $\{\eta_t\}$ is a zero mean normal white noise sequence with variance σ_η^2 . The position x_t can be observed only with error, the observed position y_t at time t being

$$y_t = x_t + \epsilon_t,$$

where ϵ_t is a zero mean normal white noise with variance σ_ϵ^2 , uncorrelated with the other variables.

(i) Show that if $\beta_t = (x_t, v)^\top$, the system can be described by the equations

$$y_t = g^\top \beta_t + \epsilon_t, \tag{1}$$

$$\beta_t = F\beta_{t-1} + \zeta_t, \tag{2}$$

where g is a constant vector, F is a constant matrix and ζ_t is a random vector. Give the values of g and F and write down the mean and covariance matrix of ζ_t . What are equations (2) and (3) called in the context of state space linear modelling? **(4 marks)**

(ii) Suppose that hourly observations of the plane's position are available up to time $t - 1$ and it is believed that the true position x_{t-1} at time $t - 1$ has posterior normal distribution with mean m_{t-1} and variance v_{t-1} . Show that, given this information, the prior mean of the position x_t at time t is $m_{t-1} + v$ and that the prior variance is $v_{t-1} + \sigma_\eta^2$. **(4 marks)**

(iii) Find the posterior distribution of the true position x_t at time t when an observation y_t of the plane's position at time t becomes available. **(8 marks)**

(iv) Suppose that $v = 720$, $\sigma_\epsilon = 10$ and $\sigma_\eta = 30$. If, one hour into the flight, an observation gives the position as $y_1 = 750$ km and it was known with certainty that at time $t = 0$ the true position had been $x_0 = 0$, find the posterior mean of the plane's position at $t = 1$, and the associated posterior variance. Hence give an approximate 95% credible interval for the true position x_1 after one hour. **(4 marks)**

End of Question Paper