



The
University
Of
Sheffield.

SCHOOL OF MATHEMATICS AND STATISTICS

**Autumn Semester
2017–2018**

Multivariate Data Analysis

2 hours

*Marks will be awarded for your best **three** 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 75 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 70 randomly selected kernels of the Kama variety of wheat were measured using a soft X-ray technique (Charytanowicz et al, 2010). Several characteristics were measured:

- x_1 a measure of compactness (got from the area and perimeter)
- x_2 roundedness (the width as a percentage of length)
- x_3 groove length as percentage of length

The means were given by $\bar{x} = (\bar{x}_1, \bar{x}_2, \bar{x}_3) = (88.01, 58.94, 92.35)$, with variance matrix $S = \begin{pmatrix} 2.621 & 4.107 & -0.090 \\ 4.107 & 7.901 & -0.670 \\ -0.090 & -0.670 & 5.817 \end{pmatrix}$. The data is believed to look approximately multivariate normal.

You may use the R output $qt(0.975, 69) = 1.994$, and the following values of the F -distribution (not all of which will be relevant):

$$qf(0.95, 2, 67) = 3.134 \quad qf(0.95, 2, 68) = 3.132 \quad qf(0.95, 2, 69) = 3.130.$$

- (i) Compute the correlation between x_1 and x_2 . **(1 mark)**
- (ii) There is interest in comparing the Kama variety with another, ancient variety. For the ancient variety, the means for the compactness and roundedness variables x_1 and x_2 are estimated as 88.36 and 58.44 respectively.
 - (a) Perform a t -test to test the hypothesis that the mean compactness for the Kama variety is equal to 88.36. **(2 marks)**
 - (b) Perform a t -test to test the hypothesis that the mean roundedness for the Kama variety is equal to 58.44. **(2 marks)**
 - (c) Test the hypothesis that the compactness and roundedness of the Kama variety are simultaneously the same as the ancient variety.

You may use the fact that $\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$.

(6 marks)
 - (d) Comment on the results of (a)–(c). **(3 marks)**
- (iii) Work out the mean and variance for the difference $x_1 - x_3$ for the Kama variety. Without doing any calculations, say how you would test that x_1 and x_3 had equal means. **(4 marks)**
- (iv) The same measurements were taken for 70 samples of the Rosa variety. For this variety, the mean for x_1 is 88.35, and the mean for x_2 is 59.86; the variance matrix for these two variables is $S_2 = \begin{pmatrix} 2.403 & 3.907 \\ 3.907 & 7.903 \end{pmatrix}$. Test the hypothesis that the means of these two variables are the same for the Kama and the Rosa varieties. You may assume that $F_{2,k}(0.95) > 2.5$ for all k . **(6 marks)**
- (v) A third variety, of Canadian wheat, was tested in the same way. How would you test whether all three varieties had the same means? **(1 mark)**

- 2 (i) If X_1, \dots, X_p denote the values of the original variables, and Y_1, \dots, Y_p denote the principal components, then $Y' = X'A$ where A is the matrix whose columns are the normalised eigenvectors in descending order of eigenvalue.

Show that $X = AY$, stating which properties of the matrix A you are using.

Explain why $A'\text{var}(X)A$ is diagonal, assuming any results from the course. *(4 marks)*

- (ii) The `state.x77` dataset contains data on the 50 US states with the following variables:

<code>Population</code>	population estimate as of July 1, 1975
<code>Income</code>	per capita income (1974)
<code>Illiteracy</code>	illiteracy (1970, percent of population)
<code>Life Exp</code>	life expectancy in years (1969–71)
<code>Murder</code>	murder rate per 100,000 population (1976)
<code>HS Grad</code>	percent high-school graduates (1970)
<code>Frost</code>	mean number of days below freezing (1931–1960)
<code>Area</code>	land area in square miles

A principal components analysis was carried out on the data, and an edited R transcript (using the `screeplot` function from the course) is given on the following two pages.

- (a) Can you use the information to say which variable is most positively correlated with `HS Grad`? Which is the most negatively correlated? *(2 marks)*
- (b) What proportion of the variance is explained by the first two principal components? What would seem like a reasonable number of components to work with? Justify your answer. *(4 marks)*
- (c) What are the characteristics of a state which gets a high score on PC1? *(3 marks)*
- (d) Explain why Alaska (AK) and California (CA) should score so highly on PC2. *(4 marks)*
- (e) Estimate the mean score of all the states on PC3. (Recall that the option `cor=TRUE` was used in the PCA command.) *(3 marks)*
- (f) There is an outlier for PC3. What are its likely characteristics? Which of the states listed in the first R command is most likely to be the outlier? *(3 marks)*
- (g) Investigators are interested in how the `Murder` variable depends on the other variables. What would you use for this instead of PCA? Give the R command. *(2 marks)*

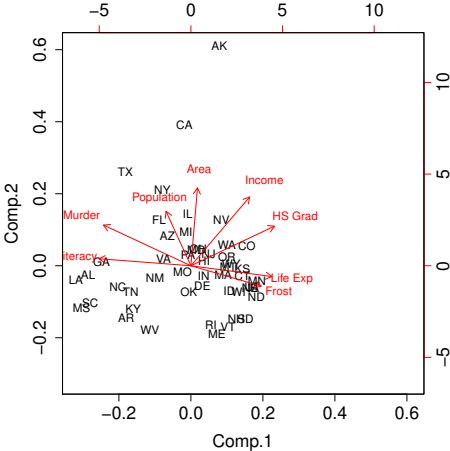
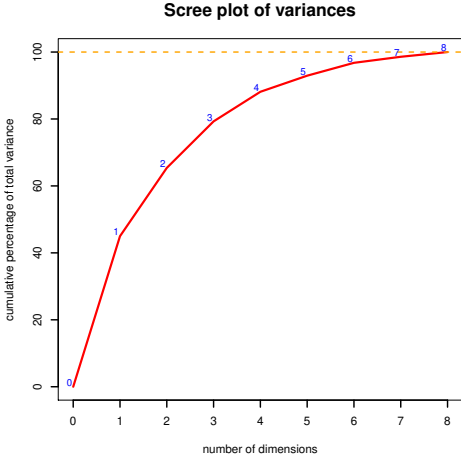
2 (continued)

```

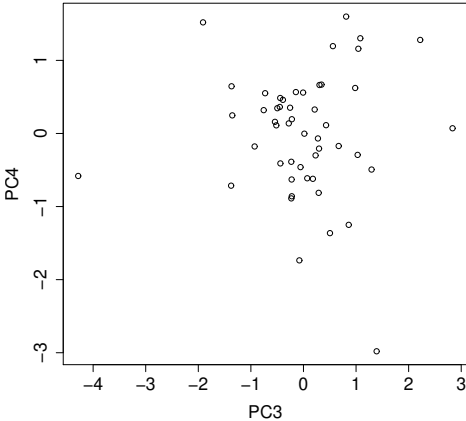
> state.x77[c(1:5,11),]
  Population Income Illiteracy Life Exp Murder HS Grad Frost Area
AL      3615   3624      2.1   69.05   15.1   41.3   20 50708
AK       365   6315      1.5   69.31   11.3   66.7  152 566432
AZ      2212   4530      1.8   70.55    7.8   58.1   15 113417
AR       210   3378      1.9   70.66   10.1   39.9   65 51945
CA     21198   5114      1.1   71.71   10.3   62.6   20 156361
> apply(state.x77,2,mean)
Population Income Illiteracy Life Exp Murder HS Grad Frost Area
  4246.42 4435.80      1.17   70.88   7.38   53.11 104.46 70735.88
> apply(state.x77,2,sd)
Population Income Illiteracy Life Exp Murder HS Grad Frost Area
  4464.49 614.47      0.61   1.34   3.69   8.08  51.98 85327.30
> state.pca<-princomp(state.x77,cor=TRUE)
> summary(state.pca)
Importance of components:
                Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8
Standard deviation    1.90  1.277  1.054 0.8411 0.6202 0.5545 0.3801 0.3364
Proportion of Variance 0.45  0.204  0.139 0.0884 0.0481 0.0384 0.0181 0.0141
Cumulative Proportion 0.45  0.654  0.793 0.8813 0.9294 0.9678 0.9859 1.0000
> loadings(state.pca)
Loadings:
      Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8
Population -0.126  0.411  0.656  0.409 -0.406          0.219
Income      0.299  0.519  0.100          0.638  0.462
Illiteracy -0.468          -0.353          0.387 -0.620  0.339
Life Exp    0.412          0.360 -0.443 -0.327  0.219 -0.256 -0.527
Murder     -0.444  0.307 -0.108  0.166  0.128 -0.325 -0.295 -0.678
HS Grad     0.425  0.299          -0.232          -0.645 -0.393  0.307
Frost       0.357 -0.154 -0.387  0.619 -0.217  0.213 -0.472
Area        0.588 -0.510 -0.201 -0.499  0.148  0.286
> screeplot(state.x77,cor=TRUE)
> state.pc<-predict(state.pca)
> biplot(state.pca)
> plot(state.pc[,3:4],xlab="PC3",ylab="PC4")

```

2 (continued)



Principal Components: PC3, PC4



- 3 (i) Suppose that we use linear discriminant analysis to classify *univariate* data between two groups. Write $p_1(x)$ for the probability that an observation with reading x is classified in group 1, and $p_2(x) = 1 - p_1(x)$ that the observation is classified into group 2. We suppose that group 1 is generated from a univariate normal distribution with mean μ_1 and variance σ^2 , and that group 2 is generated from a univariate normal distribution with mean μ_2 and variance σ^2 . Write

$$f_i(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_i)^2\right)$$

for the density of the group i distribution, and let

$$p_1(x) = \frac{f_1(x)}{f_1(x) + f_2(x)}, \quad p_2(x) = \frac{f_2(x)}{f_1(x) + f_2(x)},$$

which could be interpreted as the probabilities of an observation x coming from the two groups. Show that

$$\log\left(\frac{p_1(x)}{1 - p_1(x)}\right) = c_0 + c_1x$$

for some constants c_0 and c_1 which you should give explicitly. **(5 marks)**

Which other discrimination/classification technique has a similar formula? **(1 mark)**

- (ii) Suppose that group 1 is generated from a bivariate normal distribution with mean $\mu_1 = (-1, -1)'$ and variance $\Sigma_1 = \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix}$, and group 2 is generated from a bivariate normal distribution with mean $\mu_2 = (1, 1)'$ and variance $\Sigma_2 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix}$.

Show that the decision boundary between the two groups is in fact a pair of straight lines, and draw a picture of the plane explaining which regions would be classified as in group 1, and which as group 2.

(7 marks)

- (iii) What is the main difference in modelling between linear and quadratic discriminant analysis? **(2 marks)**

3 (continued)

(iv) Lohweg et al (2012) recorded data on two measurements (x_1, x_2) (the *variance* and *skewness*) for 100 genuine banknotes and 100 forged banknotes. For the genuine banknotes, the means are $(2.297, 3.764)$, and for the forged notes, the means are $(-1.903, -1.118)$. The inverse pooled variance matrix is $S^{-1} = \begin{pmatrix} 0.28 & 0.01 \\ 0.01 & 0.04 \end{pmatrix}$.

(a) Give Fisher's linear discriminant function for this data. **(5 marks)**

(b) Hence classify a note with measurements $(1, 0)$. **(1 mark)**

(c) It would appear that the data for each group is reasonably well approximated by a bivariate normal distribution, and that the variance matrices are very similar. Estimate the probability of misclassifying a genuine banknote as forged. You may leave your answer in the form $\Phi(z)$ for some number z which you should determine.

(4 marks)

- 4 (i) Consider the following simple data frame, consisting of 4 observations Ob1, Ob2, Ob3 and Ob4 of 3 variables, X1, X2, X3 and a response variable Y:

	X1	X2	X3	Y
Ob1	0	1	0	1
Ob2	0	1	2	4
Ob3	1	0	1	7
Ob4	1	1	2	3

- (a) Construct a regression tree for this data. *(7 marks)*
- (b) Hence predict the response for a variable with inputs $X1 = 1$, $X2 = 1$ and $X3 = 0$. *(1 mark)*
- (c) Suppose now that the response column is replaced by a class, where Ob1 and Ob2 are both in class A, and Ob3 and Ob4 are both in class B. Compute the Gini index for the split $X2 < 0.5$. *(2 marks)*

- (ii) Suppose that we take n observations of an identically distributed random variable with variance σ^2 .

- (a) Assuming our n observations are independent, what is the variance of the mean? *(1 mark)*
- (b) If we knew that all observations would be identical, what would be the variance of the mean? *(1 mark)*
- (c) Now suppose that any two of our n observations have a correlation of $\rho > 0$. Show that the variance of the mean is $\rho\sigma^2 + \frac{1-\rho}{n}\sigma^2$. How does this behave as $n \rightarrow \infty$?

[You may use the result that the variance of the mean is given by

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{Cov}(X_i, X_j).] \quad (4 \text{ marks})$$

- (d) Explain the significance of this calculation to the construction of random forests. *(3 marks)*

- (iii) Suppose that we have 4 observations, for which we compute dissimilarities

as $\begin{pmatrix} 0.3 & & & \\ 0.4 & 0.5 & & \\ 0.7 & 0.8 & 0.45 & \\ & & & \end{pmatrix}$.

- (a) Complete this to a 4×4 -matrix of distances between the observations. *(1 mark)*
- (b) On your graph paper, draw the dendrogram that results from hierarchical clustering using single linkage, plotting the heights at which clusters fuse. Do the same for clustering with complete linkage. *(5 marks)*

End of Question Paper