

**MTH 552: STATISTICAL & AI TECHNIQUES IN DATA MINING**  
**Mid semester Examination: Full Marks 60**

- [1] (a) From a dataset of track records of 100 athletes in 5 track and field events; 100m, 200m, 400m, 5000m, 10000m; calculations based on sample correlation matrix gives the following:

$$\hat{\lambda}_1 = 3.75; \hat{f}_1' = (0.469, 0.532, 0.465, 0.387, 0.361)$$

$$\hat{\lambda}_2 = 1.01; \hat{f}_2' = (-0.368, -0.236, -0.315, 0.585, 0.606)$$

$\hat{\lambda}_1$  is the largest eigen value of the sample correlation matrix,  $\hat{f}_1$  is the corresponding orthonormalized eigen vector;  $\hat{\lambda}_2$  is the second largest eigen value of the sample correlation matrix,  $\hat{f}_2$  is the corresponding orthonormalized eigen vector.

The standardized observation vectors of two athletes of interest are

Athlete1:  $(0.01, 0.06, 0.8, 0.65, 0.75)'$

Athlete2:  $(0.8, 0.7, 0.3, 0.12, 0.01)'$

- (i) How would you interpret the 1<sup>st</sup> 2 principal components?
  - (ii) What proportion of total (standardized) sample variation does the 1<sup>st</sup> principal component explain?
  - (iii) With the given information, can you suggest a ranking of the 2 athletes?
- (b) Let  $\underline{X} = (\dot{X}_1, \dots, \dot{X}_p)'$  be a random vector with  $E(\underline{X}) = \underline{\mu} = (\mu_1, \dots, \mu_p)'$  and  $Cov(\underline{X}) = \Sigma = ((\sigma_{ij}))$ ,  $(\Sigma > 0)$  and  $\underline{Y} = (Y_1, \dots, Y_p)'$  denote the vector of principal components derived from standardized variables  $\underline{Z} = (Z_1, \dots, Z_p)'$ ; for  $i = 1, \dots, p$ ,  $Z_i = (X_i - \mu_i) / \sqrt{\sigma_{ii}}$ . Find the covariance matrix,  $Cov(\underline{X}, \underline{Y})$ .

**12 (6+6) Marks**

- [2] The distance matrix corresponding to 6 multidimensional cases  $C_1, C_2, C_3, C_4, C_5, C_6$  is given by

$$D = \begin{pmatrix} 0 & 4 & 13 & 24 & 12 & 8 \\ & \ddots & 0 & 10 & 22 & 11 & 10 \\ & & & 0 & 7 & 3 & 9 \\ & & & & 0 & 6 & 18 \\ & & & & & 0 & 8.5 \\ & & & & & & 0 \end{pmatrix}$$

- (a) Construct the dendrogram tree corresponding to an agglomerative complete linkage hierarchical clustering algorithm.
- (b) Identify the clusters at a merger level 8.
- (c) Find the merger level at which we get 4 clusters and list the objects in the clusters.

**10 Marks**

- [3] Consider the divergence distance measure between two multidimensional ( $p$ -dimensional) populations ( $\pi_1$  and  $\pi_2$ )

$$J_D = \int \dots \int (f(x|\pi_1) - f(x|\pi_2)) \log(f(x|\pi_1)/f(x|\pi_2)) dx$$

$f(x|\pi_i)$  denote the joint density under population  $\pi_i$ ,  $i=1,2$ . Prove or disprove the following statement "If the  $p$ -components of the underlying random vector are independent then  $J_D = \sum_{i=1}^p J_D^i$ ,  $J_D^i$  is the divergence distance for the  $i^{th}$  component of the random vector".

6 Marks

- [4] Let  $(X_1, \dots, X_n)$  be random sample from a population having a mixture exponential model density

$$p(x) = \sum_{j=1}^g \pi_j p(x|\theta_j),$$

$$\text{where, } p(x|\theta_j) = \begin{cases} \theta_j e^{-\theta_j x}, & x > 0 \\ 0, & \text{o/w} \end{cases}, j=1, \dots, g.$$

- (a) Formulate the maximum likelihood estimation of the parameters involved in the E-M algorithm framework.  
(b) Derive the E-M algorithm update equations for  $\pi_j$  and  $\theta_j$  and hence outline the density estimation procedure based on a given set of observations.

12 Marks

- [5] Let  $x_1 = (1, 2)'$ ,  $x_2 = (3, 2)'$ ,  $x_3 = (1, 6)'$  &  $x_4 = (5, 4)'$  be observed feature vectors of 4 cases. 3 different clustering algorithms gave the following 3 partitions:

Algorithm I Partition: (Case 1, Case 3), (Case 2, Case 4)

Algorithm II Partition: (Case 1, Case 2), (Case 3, Case 4)

Algorithm III Partition: (Case 1, Case 4), (Case 2, Case 3)

Let  $S_W = \frac{1}{n} \sum_{j=1}^g \sum_{i=1}^n Z_{ji} (x_i - \bar{m}_j)(x_i - \bar{m}_j)'$  be the pooled within cluster scatter matrix for a fixed number,  $g$ , of clusters obtained from  $n$  cases.  $Z_{ji} = 1$ , if  $x_i \in \text{cluster } j$ ; 0, otherwise.  $\bar{m}_j$  is the mean of cluster  $j$ . Which of the above partition(s) would you prefer if clustering criterion based on  $\text{trace}(S_W)$  is to be used?

10 Marks

- [6] Let  $(20, 10, 16, 2, 3, 4, 4, 8, 1, 12, 11, 19, 18, 21, 5, 11, 11, 12, 19, 2)$  be a sample from an univariate population with unknown probability density function  $f(x)$ .

- (a) Find non-parametric, rectangular kernel based, density estimate at  $x = 7, 11, 24, 30$  with kernel bandwidth,  $h$ , equal to 4 for the rectangular kernel.  
(b) Find non-parametric  $k$ -nearest neighbor (with a symmetric region having center at the point  $x$ ) density estimates at the same points as in (a) using  $k = 4$ .

10 Marks