

EXAM IN MULTIVARIATE METHODS

March 25 2020

Time: 6 hours

The exam is for individual solving. It is an open-book exam, but you are not allowed to use the help of other students, friends, family, or similar. In case you need clarification, the teacher is available at Zoom:

Join Zoom Meeting from

<https://stockholmuniiversity.zoom.us/j/9255236581>

Meeting ID: 925 523 6581

The exam consists of five questions. To score maximum points on a question solutions need to be clear, detailed and well-motivated.

Question. 1 (6+2+3+3+2=16 Points)

For a data set with observations on two variables x_1 and x_2 the sample correlation matrix was found to be

$$R = \begin{bmatrix} 1 & 0.50 \\ 0.50 & 1 \end{bmatrix}.$$

$$\text{Var}(x_1) = 65.41 \text{ and } \text{Var}(x_2) = 1.27$$

- Using R, construct two principal components that are orthogonal to each other.
- What proportion of variance is accounted by these principal components?
- Compute the loadings of the variables.
- Find the covariance matrix (S).
- Find the generalized variance using S matrix.

Question. 2 (3+2+4+4+3=16 Points)

The sample correlation matrix given below arises from the scores of 220 boys in six school subjects: (1) French, (2) English, (3) History, (4) Arithmetic, (5) Algebra, and (6) Geometry.

$$R = \begin{matrix} & \begin{matrix} French \\ English \\ History \\ Arithmetic \\ Algebra \\ Geometry \end{matrix} & \begin{bmatrix} 1 & & & & & \\ 0.44 & 1 & & & & \\ 0.41 & 0.35 & 1 & & & \\ 0.29 & 0.35 & 0.16 & 1 & & \\ 0.33 & 0.32 & 0.19 & 0.59 & 1 & \\ 0.25 & 0.33 & 0.18 & 0.47 & 0.46 & 1 \end{bmatrix} \end{matrix}$$

A factor analysis was performed to analyze the correlation matrix from the scores of boys in six school subjects by the Principal component method where two factors were extracted. The two factors are assumed uncorrelated. The un-rotated two-factor solution is given below

Variable	F1	F2
French	0.66	0.44
English	0.69	0.29
History	0.52	0.64
Arithmetic	0.74	-0.42
Algebra	0.74	-0.37
Geometry	0.68	-0.35

Based on these reported results obtain:

- The communalities.
- The proportion of variance explained by each factor.
- The estimated/reproduced correlation matrix.
- The residual correlation matrix.
- RMSR.

Question. 3 (8+8=16 Points)

For the following data

Group	Y ₁	Y ₂
1	1	5
1	2	6
1	5	2
1	4	3
2	6	10
2	8	5
2	10	4
2	5	6
2	11	5

- a) Compute the **SSCP_b**, **SSCP_w** and **SSCP_t** matrices.
 b) Suppose n₁=5 and n₂=6 are observations in group-1 and group-2, respectively and

$$\text{Within-group covariance matrix for group-I} = S_1 = \begin{bmatrix} 9.70 & -3.45 \\ -3.45 & 2.70 \end{bmatrix}$$

$$\text{Within-group covariance matrix for group-II} = S_2 = \begin{bmatrix} 11.2 & -3.0 \\ -3.0 & 4.4 \end{bmatrix}$$

$$\bar{X}_1 = \begin{bmatrix} 4.2 \\ 3.8 \end{bmatrix} \text{ and } \bar{X}_2 = \begin{bmatrix} 7 \\ 6 \end{bmatrix}$$

Calculate Fisher's linear discriminant function for this data set.

Question. 4 (4+4+4+4=16 Points)

Observations on two variables were made for five subjects according to the following table.

Subject	Variable-1	Variable-2
1	1	4
2	2	5
3	5	3
4	4	2
5	6	8

- a) Construct a similarity matrix containing squared Euclidean distances

b) Use the similarity matrix in part (a) and perform a cluster analysis with the following method

- I. Farthest neighbor method.
- II. Nearest neighbor method.
- III. Average linkage method.

Question. 5 (2+3+3+4+4=16 Points)

A company that manufactures riding mowers wants to identify the best sales prospects for an intensive sales campaign. In particular, the manufacturer is interested in classifying households as prospective owners or nonowners on the basis of Income (in \$1000s) and Lot Size (in 1000 ft²). Data were collected and a logistic regression was fitted:

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-25.9382	11.4871	-2.258	0.0239 *
Income	0.1109	0.0543	2.042	0.0412 *
Lot Size	0.9638	0.4628	2.083	0.0373 *

The following table displays observations on 12 riding-mower owners and 12 nonowners as well as the estimated probability to be an owner based on the logistic regression.

Ownership	Income	Lot Size	\hat{P}	Ownership	Income	Lot Size	\hat{P}
Owner	60	18.4	0.17	Nonowner	75	19.6	0.78
Owner	85.5	16.8	0.43	Nonowner	52.8	20.8	0.49
Owner	64.8	21.6	0.89	Nonowner	64.8	17.2	0.10
Owner	61.5	20.8	0.72	Nonowner	43.2	20.4	0.18
Owner	87	23.6	0.99	Nonowner	84	17.6	0.58
Owner	110.1	19.2	0.99	Nonowner	49.2	17.6	0.03
Owner	108	17.6	0.95	Nonowner	59.4	16	0.02
Owner	82.8	22.4	0.99	Nonowner	66	18.4	0.29
Owner	69	20	0.72	Nonowner	47.4	16.4	0.01
Owner	93	20.8	0.98	Nonowner	33	18.8	0.02
Owner	51	22	0.71	Nonowner	51	14	0.01
Owner	81	20	0.91	Nonowner	63	14.8	0.01

- a) Interpret the parameter $\hat{\beta}_1$.
- b) What are the odds that a household with a \$60000 income and a lot size of 20 000 ft² is an owner?

- c) What is the classification of a household with a \$60000 income and a lot size of 20 000 ft²?
- d) What is the minimum income that a household with 16,000 ft² lot size should have before it is classified as an owner?
- e) Classify the observations given in the table and compute the sensitivity and specificity.