

STA 825 Final Exam – Wednesday, March 18
SHOW ALL WORK

Name: _____

1. A study was conducted to evaluate the effects of different washing treatments on the quality of beef after storage. There were three washing treatments: normal, chlorine, and lactic acid. On each of four consecutive days, three samples of beef were prepared and each sample was randomly assigned to one of the three treatments such that each treatment was observed on one sample per day. The experiment had to be conducted over several days because it was not possible to prepare and treat more than three samples per day. After treatment, all samples were stored for 10 days. After storage, each of the samples was evaluated on six criteria: beefy aroma, bloody/serummy aroma, metallic aroma, grassy/barnyard aroma, sour aroma, and spoiled aroma.

a. Why is this a multivariate rather than a univariate data situation?

b. What is the design here? (Hint: we only talked about three designs, and its one of those three.)

c. Write down an appropriate MANOVA model for analyzing these data. Include a very brief description of each of the terms of your model.

d. In terms of the parameters of your model, write down the null hypothesis for testing that there is no difference among the three washing methods.

- e. After testing the null hypothesis from (d), we will typically want to obtain simultaneous confidence intervals for specific comparisons among the treatment means on particular response variables. Give one reason that the Bonferroni method may be preferred in some cases over the Roy method for obtaining simultaneous confidence intervals. Give one reason that the Roy method may be preferred in some cases over the Bonferroni method.

2. In the early 1900's, several investigators were interested in predicting behavioral and social outcomes among people based on physical characteristics. Macdonnell (1902) reports a correlation matrix for the following seven physical variables measured on 3000 British criminals: (1) head length, (2) head breadth, (3) face breadth, (4) left finger length, (5) left forearm length, (6) left foot length, and (7) height. Assume that all original variables were measured in centimeters.

The attached SAS program and output (labelled "Final Exam Problem #2") performs a principal components analysis based Macdonnell's correlation matrix. Answer the following questions:

a. One of the goals of principal components analysis is to reduce the dimension of the original data. How would you choose the number of principal components to retain for subsequent analyses? In this example, how many principal components would you retain?

b. Briefly interpret the first two principal components in this example. That is, what aspect of the original variables is captured by the first principal component? the second?

c. Explain why it is not appropriate for this example to perform a principal components analysis on the covariance matrix rather than the correlation matrix.

d. Suppose that in addition to the seven variables described above, an eighth variable, computed as head length minus head breadth, had been included in the analysis to capture head shape. What would be the variance of the last principal component in such an analysis and why?

3. Annual financial data are available on firms. Four financial variables including $x_1 = (\text{cash flow})/(\text{total debt})$, $x_2 = (\text{net income})/(\text{total assets})$, $x_3 = (\text{current assets})/(\text{current liabilities})$, and $x_4 = (\text{current assets})/(\text{net sales})$, were collected for 21 firms that subsequently went bankrupt and 25 financially sound firms at about the same point in time. A discriminant analysis for these data is performed in the attached SAS program and output (labelled “Final Exam Problem #3”). The discriminant analysis is based only on x_2 and x_3 . Answer the following questions.
- a. In the SAS program, a hypothesis test is performed using PROC GLM of the hypothesis that the mean vectors are the same in the two groups (financially sound and financially troubled firms). Why is this test performed prior to performing a discriminant analysis, and what does the result of the hypothesis test say about how the discriminant analysis will perform?

b. In this example, equal priors and costs of misclassification were assumed. Explain why it may be more appropriate to use different costs of misclassification in this analysis? Why might it be more appropriate to use different prior probabilities?

c. Write down the linear discriminant rule for classifying firms as financially sound or financially troubled.

- d. There are two estimated error rates in the output. What error rate would you use and why, if you were describing how you expected the linear discriminant rule from (d) to perform in practice (when actually classifying firms of unknown status as either financial sound or troubled)?

4. The following table lists measurements on 5 nutritional variables for 12 breakfast cereals.

A cluster analysis for these 12 cereal brands is performed in the attached SAS program and output (labelled “Final Exam Problem #4”). Answer the following questions.

- a. Identify the type of clustering algorithm that was used in the SAS program. (Don't describe the algorithm, just tell me the type of algorithm it is — more than one adjective is necessary).

- b. What is the cluster structure for 5 clusters given by the results of PROC CLUSTER in this example. That is, assuming there are 5 clusters, tell me which cereals belong together.

c. Based on the pseudo- T^2 statistics printed in the SAS output, what is the appropriate number of clusters for this example.

d. Attached to the output is a page of star plots for the 12 cereals. Explain why just a dot appears for Puffed Rice.