

11.7 Logistic Regression and Classification

Introduction

The classification functions already discussed are based on quantitative variables. Here we discuss an approach to classification where some or all of the variables are qualitative. This approach is called logistic regression. In its simplest setting, the response variable Y is restricted to two values. For example, Y may be recorded as “male” or “female” or “employed” and “not employed.”

Even though the response may be a two outcome qualitative variable, we can always code the two cases as 0 and 1. For instance, we can take male = 0 and female = 1. Then the probability p of 1 is a parameter of interest. It represents the proportion in the population who are coded 1. The mean of the distribution of 0's and 1's is also p since

$$\text{mean} = 0 \times (1 - p) + 1 \times p = p$$

The proportion of 0's is $1 - p$ which is sometimes denoted as q .

The variance of the distribution is

$$\text{variance} = 0^2 \times (1 - p) + 1^2 \times p - p^2 = p(1 - p)$$

It is clear the variance is not constant. For $p = .5$, it equals $.5 \times .5 = .25$ while for $p = .8$, it is $.8 \times .2 = .16$. The variance approaches 0 as p approaches either 0 or 1.

Let the response Y be either 0 or 1. If we were to model the probability of 1 with a single predictor linear model, we would write

$$p = E(Y|z) = \beta_0 + \beta_1 z$$

and then add an error term ϵ . But there are serious drawbacks to this model.

- The predicted values of the response Y could become greater than 1 or less than 0 because the linear expression for its expected value is unbounded.
- One of the assumptions of a regression analysis is that the variance of Y is constant across all values of the predictor variable Z . We have shown this is not the case. Of course, weighted least squares might improve the situation.

We need another approach to introduce predictor variables or covariates \mathbf{Z} into the model (see [26]). Throughout, if the covariates are not fixed by the investigator, the approach is to make the models for $p(z)$ conditional on the observed values of the covariates $\mathbf{Z} = \mathbf{z}$.

The Logit Model

Instead of modeling the probability p directly with a linear model, we first consider the *odds ratio*

$$\text{odds} = \frac{p}{1 - p}$$

which is the ratio of the probability of 1 to the probability of 0. Note, unlike probability, the odds ratio can be greater than 1. If a proportion .8 of persons will get

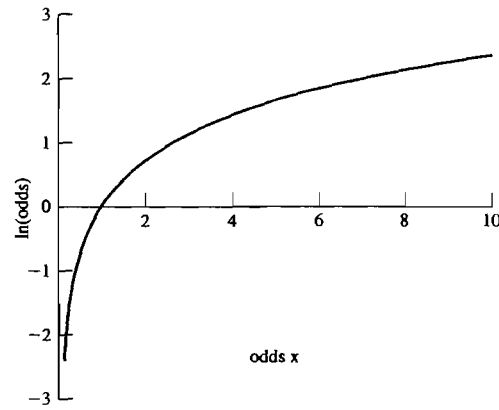


Figure 11.15 Natural log of odds ratio.

through customs without their luggage being checked, then $p = .8$ but the odds of not getting checked is $.8/.2 = 4$ or 4 to 1 of not being checked. There is a lack of symmetry here since the odds of being checked are $.2/.8 = 1/4$. Taking the natural logarithms, we find that $\ln(4) = 1.386$ and $\ln(1/4) = -1.386$ are exact opposites.

Consider the natural log function of the odds ratio that is displayed in Figure 11.15. When the odds x are 1, so outcomes 0 and 1 are equally likely, the natural log of x is zero. When the odds x are greater than one, the natural log increases slowly as x increases. However, when the odds x are less than one, the natural log decreases rapidly as x decreases toward zero.

In logistic regression for a binary variable, we model the natural log of the odds ratio, which is called *logit*(p). Thus

$$\text{logit}(p) = \ln(\text{odds}) = \ln\left(\frac{p}{1-p}\right) \quad (11-69)$$

The logit is a function of the probability p . In the simplest model, we assume that the logit graphs as a straight line in the predictor variable Z so

$$\text{logit}(p) = \ln(\text{odds}) = \ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 z \quad (11-70)$$

In other words, the log odds are linear in the predictor variable.

Because it is easier for most people to think in terms of probabilities, we can convert from the logit or log odds to the probability p . By first exponentiating

$$\ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 z$$

we obtain

$$\theta(z) = \frac{p(z)}{1-p(z)} = \exp(\beta_0 + \beta_1 z)$$

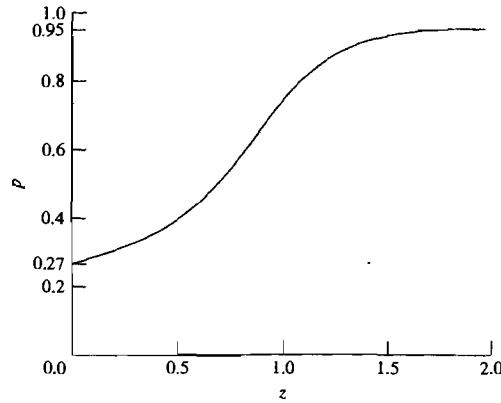


Figure 11.16 Logistic function with $\beta_0 = -1$ and $\beta_1 = 2$.

where $\exp = e = 2.718$ is the base of the natural logarithm. Next solving for $\theta(z)$, we obtain

$$p(z) = \frac{\exp(\beta_0 + \beta_1 z)}{1 + \exp(\beta_0 + \beta_1 z)} \quad (11-71)$$

which describes a *logistic curve*. The relation between p and the predictor z is not linear but has an *S-shaped* graph as illustrated in Figure 11.16 for the case $\beta_0 = -1$ and $\beta_1 = 2$. The value of β_0 gives the value $\exp(\beta_0)/(1 + \exp(\beta_0))$ for p when $z = 0$.

The parameter β_1 in the logistic curve determines how quickly p changes with z but its interpretation is not as simple as in ordinary linear regression because the relation is not linear, either in z or β_1 . However, we can exploit the linear relation for log odds.

To summarize, the logistic curve can be written as

$$p(z) = \frac{\exp(\beta_0 + \beta_1 z)}{1 + \exp(\beta_0 + \beta_1 z)} \quad \text{or} \quad p(z) = \frac{1}{1 + \exp(-\beta_0 - \beta_1 z)}$$

Logistic Regression Analysis

Consider the model with several predictor variables. Let $(z_{j1}, z_{j2}, \dots, z_{jr})$ be the values of the r predictors for the j -th observation. It is customary, as in normal theory linear regression, to set the first entry equal to 1 and $\mathbf{z}_j = [1, z_{j1}, z_{j2}, \dots, z_{jr}]'$. Conditional on these values, we assume that the observation Y_j is Bernoulli with success probability $p(\mathbf{z}_j)$, depending on the values of the covariates. Then

$$P(Y_j = y_j) = p^{y_j}(\mathbf{z}_j)(1 - p(\mathbf{z}_j))^{1-y_j} \quad \text{for } y_j = 0, 1$$

so

$$E(Y_j) = p(\mathbf{z}_j) \quad \text{and} \quad \text{Var}(Y_j) = p(\mathbf{z}_j)(1 - p(\mathbf{z}_j))$$

It is not the mean that follows a linear model but the natural log of the odds ratio. In particular, we assume the model

$$\ln\left(\frac{p(\mathbf{z})}{1-p(\mathbf{z})}\right) = \beta_0 + \beta_1 z_1 + \cdots + \beta_r z_r = \beta' \mathbf{z}_j \quad (11-72)$$

where $\beta = [\beta_0, \beta_1, \dots, \beta_r]'$.

Maximum Likelihood Estimation. Estimates of the β 's can be obtained by the method of maximum likelihood. The likelihood L is given by the joint probability distribution evaluated at the observed counts y_j . Hence

$$\begin{aligned} L(b_0, b_1, \dots, b_r) &= \prod_{j=1}^n p^{y_j}(\mathbf{z}_j) (1 - p(\mathbf{z}_j))^{1-y_j} \\ &= \frac{\prod_{j=1}^n e^{y_j(b_0 + b_1 z_{j1} + \dots + b_r z_{jr})}}{\prod_{j=1}^n (1 + e^{b_0 + b_1 z_{j1} + \dots + b_r z_{jr}})} \end{aligned} \quad (11-73)$$

The values of the parameters that maximize the likelihood cannot be expressed in a nice closed form solution as in the normal theory linear models case. Instead they must be determined numerically by starting with an initial guess and iterating to the maximum of the likelihood function. Technically, this procedure is called an iteratively re-weighted least squares method (see [26]).

We denote the numerically obtained values of the maximum likelihood estimates by the vector $\hat{\beta}$.

Confidence Intervals for Parameters. When the sample size is large, $\hat{\beta}$ is approximately normal with mean β , the prevailing values of the parameters and approximate covariance matrix

$$\widehat{\text{Cov}}(\hat{\beta}) \approx \left[\sum_{j=1}^n \hat{p}(\mathbf{z}_j) (1 - \hat{p}(\mathbf{z}_j)) \mathbf{z}_j \mathbf{z}_j' \right]^{-1} \quad (11-74)$$

The square roots of the diagonal elements of this matrix are the large sample estimated standard deviations or standard errors (SE) of the estimators $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_r$ respectively. The large sample 95% confidence interval for β_k is

$$\hat{\beta}_k \pm 1.96 SE(\hat{\beta}_k) \quad k = 0, 1, \dots, r \quad (11-75)$$

The confidence intervals can be used to judge the significance of the individual terms in the model for the logit. Large sample confidence intervals for the logit and for the population proportion $p(\mathbf{z}_j)$ can be constructed as well. See [17] for details.

Likelihood Ratio Tests. For the model with r predictor variables plus the constant, we denote the maximized likelihood by

$$L_{\max} = L(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_r)$$

If the null hypothesis is $H_0: \beta_k = 0$, numerical calculations again give the maximum likelihood estimate of the reduced model and, in turn, the maximized value of the likelihood

$$L_{\max, \text{Reduced}} = L(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{k-1}, \hat{\beta}_{k+1}, \dots, \hat{\beta}_r)$$

When doing logistic regression, it is common to test H_0 using minus twice the log-likelihood ratio

$$-2 \ln \left(\frac{L_{\max, \text{Reduced}}}{L_{\max}} \right) \quad (11-76)$$

which, in this context, is called the *deviance*. It is approximately distributed as chi-square with 1 degree of freedom when the reduced model has one fewer predictor variables. H_0 is rejected for a large value of the deviance.

An alternative test for the significance of an individual term in the model for the *logit* is due to Wald (see [17]). The Wald test of $H_0: \beta_k = 0$ uses the test statistic $Z = \hat{\beta}_k / SE(\hat{\beta}_k)$ or its chi-square version Z^2 with 1 degree of freedom. The likelihood ratio test is preferable to the Wald test as the level of this test is typically closer to the nominal α .

Generally, if the null hypothesis specifies a subset of, say, m parameters are simultaneously 0, the deviance is constructed for the implied reduced model and referred to a chi-squared distribution with m degrees of freedom.

When working with individual binary observations Y_i , the residuals

$$\frac{Y_i - \hat{p}(z_i)}{\sqrt{\hat{p}(z_i)(1 - \hat{p}(z_i))}}$$

each can assume only two possible values and are not particularly useful. It is better if they can be grouped into reasonable sets and a total residual calculated for each set. If there are, say, t residuals in each group, sum these residuals and then divide by \sqrt{t} to help keep the variances compatible.

We give additional details on logistic regression and model checking following and application to classification.

Classification

Let the response variable Y be 1 if the observational unit belongs to population 1 and 0 if it belongs to population 2. (The choice of 1 and 0 for response outcomes is arbitrary but convenient. In Example 11.17, we use 1 and 2 as outcomes.) Once a logistic regression function has been established, and using training sets for each of the two populations, we can proceed to classify. Priors and costs are difficult to incorporate into the analysis, so the classification rule becomes

Assign z to population 1 if the estimated odds ratio is greater than 1 or

$$\frac{\hat{p}(z)}{1 - \hat{p}(z)} = \exp(\hat{\beta}_0 + \hat{\beta}_1 z_1 + \dots + \hat{\beta}_r z_r) > 1$$

Equivalently, we have the simple linear discriminant rule

Assign \mathbf{z} to population 1 if the linear discriminant is greater than 0 or

$$\ln \frac{\hat{p}(\mathbf{z})}{1 - \hat{p}(\mathbf{z})} = \hat{\beta}_0 + \hat{\beta}_1 z_1 + \cdots + \hat{\beta}_r z_r > 0 \quad (11-77)$$

Example 11.17 (Logistic regression with the salmon data) We introduced the salmon data in Example 11.8 (see Table 11.2). In Example 11.8, we ignored the gender of the salmon when considering the problem of classifying salmon as Alaskan or Canadian based on growth ring measurements. Perhaps better classification is possible if gender is included in the analysis. Panel 11.2 contains the SAS output from a logistic regression analysis of the salmon data. Here the response Y is 1 if Alaskan salmon and 2 if Canadian salmon. The predictor variables (covariates) are gender (1 if female, 2 if male), freshwater growth and marine growth. From the SAS output under Testing the Global Null Hypothesis, the likelihood ratio test result (see 11-76) with the reduced model containing only a β_0 term) is significant at the $< .0001$ level. At least one covariate is required in the linear model for the logit. Examining the significance of individual terms under the heading Analysis of Maximum Likelihood Estimates, we see that the Wald test suggests gender is not significant (p -value = .7356). On the other hand, freshwater growth and marine are significant covariates. Gender can be dropped from the model. It is not a useful variable for classification. The logistic regression model can be re-estimated without gender and the resulting function used to classify the salmon as Alaskan or Canadian using rule (11-77).

Turning to the classification problem, but retaining gender, we assign salmon j to population 1, Alaskan, if the linear classifier

$$\beta'z = 3.5054 + .2816 \text{ gender} + .1264 \text{ freshwater} + .0486 \text{ marine} \geq 0$$

The observations that are misclassified are

Row	Pop	Gender	Freshwater	Marine	Linear Classifier
2	1	1	131	355	3.093
12	1	2	123	372	1.537
13	1	1	123	372	1.255
30	1	2	118	381	0.467
51	2	1	129	420	-0.319
68	2	2	136	438	-0.028
71	2	2	90	385	-3.266

From these misclassifications, the confusion matrix is

		Predicted membership	
		π_1 : Alaskan	π_1 : Canadian
Actual	π_1 : Alaskan	46	4
	π_1 : Canadian	3	47

and the apparent error rate, expressed as a percentage is

$$\text{APER} = \frac{4 + 3}{50 + 50} \times 100 = 7\%$$

When performing a logistic classification, it would be preferable to have an estimate of the misclassification probabilities using the jackknife (holdout) approach but this is not currently available in the major statistical software packages.

We could have continued the analysis in Example 11.17 by dropping gender and using just the freshwater and marine growth measurements. However, when normal distributions with equal covariance matrices prevail, logistic classification can be quite inefficient compared to the normal theory linear classifier (see [7]).

Logistic Regression with Binomial Responses

We now consider a slightly more general case where several runs are made at the same values of the covariates \mathbf{z}_j and there are a total of m different sets where these predictor variables are constant. When n_j independent trials are conducted with the predictor variables \mathbf{z}_j , the response Y_j is modeled as a binomial distribution with probability $p(\mathbf{z}_j) = P(\text{Success} | \mathbf{z}_j)$.

Because the Y_j are assumed to be independent, the likelihood is the product

$$L(\beta_0, \beta_1, \dots, \beta_r) = \prod_{j=1}^m \binom{n_j}{y_j} p_j^{y_j} (1 - p_j)^{n_j - y_j} \quad (11-78)$$

where the probabilities $p(\mathbf{z}_j)$ follow the logit model (11-72)

PANEL 11.2 SAS ANALYSIS FOR SALMON DATA USING PROC LOGISTIC.

```
title 'Logistic Regression and Discrimination';
data salmon;
infile 'T11-2.dat';
input country gender freshwater marine;
proc logistic desc;
model country = gender freshwater marine / expb;
```

PROGRAM COMMANDS

OUTPUT

```
Logistic Regression and Discrimination

The LOGISTIC procedure

Model Information

Model                      binary logit

Response Profile

Ordered Value      country      Total
                    Frequency
1                  2          50
2                  1          50
```

(continues on next page)

PANEL 11.2 (continued)

Probability modeled is country = 2. Model Fit Statistics			
Criterion	Intercept Only	Intercept and Covariates	
AIC	140.629	46.674	
SC	143.235	57.094	
-2 Log L	138.629	38.674	

Testing Global Null Hypothesis: BETA = 0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	99.9557	3	<.0001
Wald	19.4435	3	0.0002

The LOGISTIC Procedure

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq	Exp (Est)
Intercept	1	3.5054	6.3935	0.3006	0.5835	33.293
gender	1	0.2816	0.8338	0.1140	0.7356	1.325
freshwater	1	0.1264	0.0357	12.5423	0.0004	1.135
marine	1	-0.0486	0.0146	11.1460	0.0008	0.953

The maximum likelihood estimates $\hat{\beta}$ must be obtained numerically because there is no closed form expression for their computation. When the total sample size is large, the approximate covariance matrix $\widehat{\text{Cov}}(\hat{\beta})$ is

$$\widehat{\text{Cov}}(\hat{\beta}) \approx \left[\sum_{j=1}^m n_j \hat{p}(\mathbf{z}_j)(1 - \hat{p}(\mathbf{z}_j)) \mathbf{z}_j \mathbf{z}_j' \right]^{-1} \quad (11-79)$$

and the i -th diagonal element is an estimate of the variance of $\hat{\beta}_{i+1}$. Its square root is an estimate of the large sample standard error $SE(\hat{\beta}_{i+1})$.

It can also be shown that a large sample estimate of the variance of the probability $\hat{p}(\mathbf{z}_j)$ is given by

$$\widehat{\text{Var}}(\hat{p}(\mathbf{z}_k)) \approx (\hat{p}(\mathbf{z}_k)(1 - \hat{p}(\mathbf{z}_k))^2 \mathbf{z}_k' \left[\sum_{j=1}^m n_j \hat{p}(\mathbf{z}_j)(1 - \hat{p}(\mathbf{z}_j)) \mathbf{z}_j \mathbf{z}_j' \right]^{-1} \mathbf{z}_k$$

Consideration of the interval plus and minus two estimated standard deviations from $\hat{p}(\mathbf{z}_j)$ may suggest observations that are difficult to classify.

Model Checking. Once any model is fit to the data, it is good practice to investigate the adequacy of the fit. The following questions must be addressed.

- Is there any systematic departure from the fitted logistic model?
- Are there any observations that are unusual in that they don't fit the overall pattern of the data (*outliers*)?
- Are there any observations that lead to important changes in the statistical analysis when they are included or excluded (*high influence*)?

If there is no parametric structure to the single trial probabilities $p(\mathbf{z}_j) = P(\text{Success} | \mathbf{z}_j)$, each would be estimated using the observed number of successes (1's) y_j in n_j trials. Under this nonparametric model, or saturated model, the contribution to the likelihood for the j -th case is

$$\binom{n_j}{y_j} p^{y_j}(\mathbf{z}_j) (1 - p(\mathbf{z}_j))^{n_j - y_j}$$

which is maximized by the choices $\hat{p}(\mathbf{z}_j) = y_j/n_j$ for $j = 1, 2, \dots, n$. Here $m = \sum n_j$. The resulting value for minus twice the maximized nonparametric (NP) likelihood is

$$-2 \ln L_{\max, NP} = -2 \sum_{j=1}^m \left[y_j \ln \left(\frac{y_j}{n_j} \right) + (n_j - y_j) \ln \left(1 - \frac{y_j}{n_j} \right) \right] + 2 \ln \left(\prod_{j=1}^m \binom{n_j}{y_j} \right) \quad (11-80)$$

The last term on the right hand side of (11-80) is common to all models.

We also define a deviance between the nonparametric model and a fitted model having a constant and $r-1$ predictors as minus twice the log-likelihood ratio or

$$G^2 = 2 \sum_{j=1}^m \left[y_j \ln \left(\frac{y_j}{\hat{y}_j} \right) + (n_j - y_j) \ln \left(\frac{n_j - y_j}{n_j - \hat{y}_j} \right) \right] \quad (11-81)$$

where $\hat{y}_j = n_j \hat{p}(\mathbf{z}_j)$ is the fitted number of successes. This is the specific deviance quantity that plays a role similar to that played by the residual (error) sum of squares in the linear models setting.

For large sample sizes, G^2 has approximately a chi square distribution with f degrees of freedom equal to the number of observations, m , minus the number of parameters β estimated.

Notice the deviance for the full model, G_{Full}^2 , and the deviance for a reduced model, $G_{Reduced}^2$, lead to a contribution for the extra predictor terms

$$G_{Reduced}^2 - G_{Full}^2 = -2 \ln \left(\frac{L_{\max, Reduced}}{L_{\max}} \right) \quad (11-82)$$

This difference is approximately χ^2 with degrees of freedom $df = df_{Reduced} - df_{Full}$. A large value for the difference implies the full model is required.

When m is large, there are too many probabilities to estimate under the nonparametric model and the chi-square approximation cannot be established by existing methods of proof. It is better to rely on likelihood ratio tests of logistic models where a few terms are dropped.

Residuals and Goodness-of-Fit Tests. Residuals can be inspected for patterns that suggest lack of fit of the logit model form and the choice of predictor variables (covariates). In logistic regression residuals are not as well defined as in the multiple regression models discussed in Chapter 7. Three different definitions of residuals are available.

Deviance residuals (d_j):

$$d_j = \pm \sqrt{2 \left[y_j \ln \left(\frac{y_j}{n_j \hat{p}(z_j)} \right) + (n_j - y_j) \ln \left(\frac{n_j - y_j}{n_j (1 - \hat{p}(z_j))} \right) \right]}$$

where the sign of d_j is the same as that of $y_j - n_j \hat{p}(z_j)$ and,

$$\begin{aligned} \text{if } y_j = 0, \text{ then } d_j &= -\sqrt{2n_j |\ln(1 - \hat{p}(z_j))|} \\ \text{if } y_j = n_j, \text{ then } d_j &= -\sqrt{2n_j |\ln \hat{p}(z_j)|} \end{aligned} \quad (11-83)$$

$$\text{Pearson residuals}(r_j): \quad r_j = \frac{y_j - n_j \hat{p}(z_j)}{\sqrt{n_j \hat{p}(z_j)(1 - \hat{p}(z_j))}} \quad (11-84)$$

$$\text{Standardized Pearson residuals}(r_{sj}): \quad r_{sj} = \frac{r_j}{\sqrt{1 - h_{jj}}} \quad (11-85)$$

where h_{jj} is the (j, j) th element in the “hat” matrix \mathbf{H} given by equation (11-87). Values larger than about 2.5 suggest lack of fit at the particular \mathbf{z}_j .

An overall test of goodness of fit—preferred especially for smaller sample sizes—is provided by Pearson’s chi square statistic

$$X^2 = \sum_{j=1}^m r_j^2 = \sum_{j=1}^n \frac{(y_j - n_j \hat{p}(z_j))^2}{n_j \hat{p}(z_j)(1 - \hat{p}(z_j))} \quad (11-86)$$

Notice that the chi square statistic, a single number summary of fit, is the sum of the squares of the Pearson residuals. Inspecting the Pearson residuals themselves allows us to examine the quality of fit over the entire pattern of covariates.

Another goodness-of-fit test due to Hosmer and Lemeshow [17] is only applicable when the proportion of observations with tied covariate patterns is small and all the predictor variables (covariates) are continuous.

Leverage Points and Influential Observations. The logistic regression equivalent of the hat matrix \mathbf{H} contains the estimated probabilities $\hat{p}_k(\mathbf{z}_j)$. The logistic regression version of *leverages* are the diagonal elements h_{jj} of this hat matrix.

$$\mathbf{H} = \mathbf{V}^{-1/2} \mathbf{Z} (\mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{V}^{-1/2} \quad (11-87)$$

where \mathbf{V}^{-1} is the diagonal matrix with (j, j) element $n_j \hat{p}(z_j)(1 - \hat{p}(z_j))$, $\mathbf{V}^{-1/2}$ is the diagonal matrix with (j, j) element $\sqrt{n_j \hat{p}(z_j)(1 - \hat{p}(z_j))}$.

Besides the leverages given in (11-87), other measures are available. We describe the most common called the *delta beta* or *deletion displacement*. It helps identify observations that, by themselves, have a strong influence on the regression

estimates. This change in regression coefficients, when all observations with the same covariate values as the j -th case \mathbf{z}_j are deleted, is quantified as

$$\Delta\beta_j = \frac{r_{sj}^2 h_{jj}}{1 - h_{jj}} \quad (11-88)$$

A plot of $\Delta\beta_j$ versus j can be inspected for influential cases.

11.8 Final Comments

Including Qualitative Variables

Our discussion in this chapter assumes that the discriminatory or classificatory variables, X_1, X_2, \dots, X_p have natural units of measurement. That is, each variable can, in principle, assume any real number, and these numbers can be recorded. Often, a *qualitative* or *categorical variable* may be a useful discriminator (classifier). For example, the presence or absence of a characteristic such as the color red may be a worthwhile classifier. This situation is frequently handled by creating a variable X whose numerical value is 1 if the object possesses the characteristic and zero if the object does not possess the characteristic. The variable is then treated like the measured variables in the usual discrimination and classification procedures.

Except for logistic classification, there is very little theory available to handle the case in which some variables are continuous and some qualitative. Computer simulation experiments (see [22]) indicate that Fisher's linear discriminant function can perform poorly or satisfactorily, depending upon the correlations between the qualitative and continuous variables. As Krzanowski [22] notes, "A low correlation in one population but a high correlation in the other, or a change in the sign of the correlations between the two populations could indicate conditions unfavorable to Fisher's linear discriminant function." This is a troublesome area and one that needs further study.

Classification Trees

An approach to classification completely different from the methods discussed in the previous sections of this chapter has been developed. (See [5].) It is very computer intensive and its implementation is only now becoming widespread. The new approach, called *classification and regression trees* (CART), is closely related to divisive clustering techniques. (See Chapter 12.)

Initially, all objects are considered as a single group. The group is split into two subgroups using, say, high values of a variable for one group and low values for the other. The two subgroups are then each split using the values of a second variable. The splitting process continues until a suitable stopping point is reached. The values of the splitting variables can be ordered or unordered categories. It is this feature that makes the CART procedure so general.

For example, suppose subjects are to be classified as

- π_1 : heart-attack prone
- π_2 : not heart-attack prone

on the basis of age, weight, and exercise activity. In this case, the CART procedure can be diagrammed as the tree shown in Figure 11.17. The branches of the tree actually

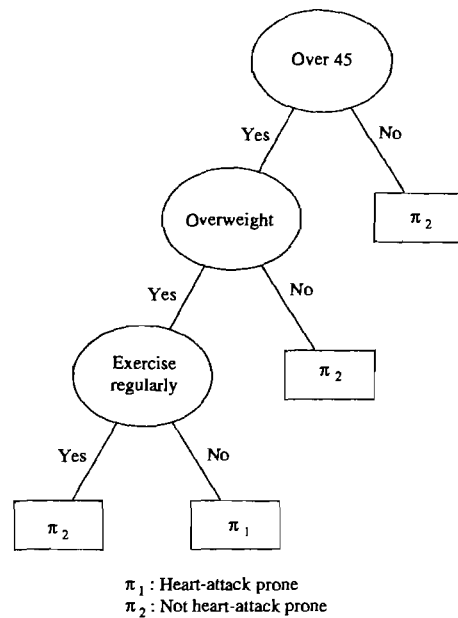


Figure 11.17 A classification tree.

correspond to divisions in the sample space. The region R_1 , defined as being over 45, being overweight, and undertaking no regular exercise, could be used to classify a subject as π_1 : heart-attack prone. The CART procedure would try splitting on different ages, as well as first splitting on weight or on the amount of exercise.

The classification tree that results from using the CART methodology with the Iris data (see Table 11.5), and variables X_3 = petal length (PetLength) and X_4 = petal width (PetWidth), is shown in Figure 11.18. The binary splitting rules are indicated in the figure. For example, the first split occurs at petal length = 2.45. Flowers with petal lengths ≤ 2.45 form one group (left), and those with petal lengths > 2.45 form the other group (right).

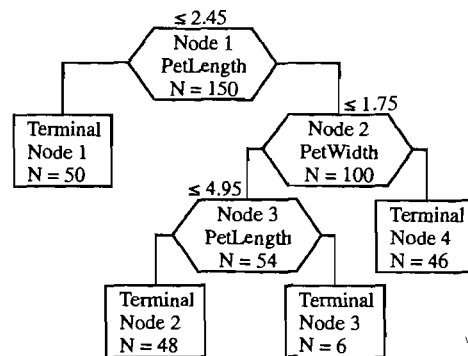


Figure 11.18 A classification tree for the Iris data.

The next split occurs with the right-hand side group (petal length > 2.45) at petal width = 1.75. Flowers with petal widths ≤ 1.75 are put in one group (left), and those with petal widths > 1.75 form the other group (right). The process continues until there is no gain with additional splitting. In this case, the process stops with four terminal nodes (TN).

The binary splits form terminal node rectangles (regions) in the positive quadrant of the X_3, X_4 sample space as shown in Figure 11.19. For example, TN #2 contains those flowers with $2.45 < \text{petal lengths} \leq 4.95$ and petal widths ≤ 1.75 —essentially the Iris Versicolor group.

Since the majority of the flowers in, for example, TN #3 are species *Virginica*, a new item in this group would be classified as *Virginica*. That is, TN #3 and TN #4 are both assigned to the *Virginica* population. We see that CART has correctly classified 50 of 50 of the *Setosa* flowers, 47 of 50 of the *Versicolor* flowers, and 49 of 50 of the *Virginica* flowers. The $\text{APER} = \frac{4}{150} = .027$. This result is comparable to the result obtained for the linear discriminant analysis using variables X_3 and X_4 discussed in Example 11.12.

The CART methodology is not tied to an underlying population probability distribution of characteristics. Nor is it tied to a particular optimality criterion. In practice, the procedure requires hundreds of objects and, often, many variables. The resulting tree is very complicated. Subjective judgments must be used to prune the tree so that it ends with groups of several objects rather than all single objects. Each terminal group is then assigned to the population holding the majority membership. A new object can then be classified according to its ultimate group.

Breiman, Friedman, Olshen, and Stone [5] have developed special-purpose software for implementing a CART analysis. Also, Loh (see [21] and [25]) has developed improved classification tree software called QUEST¹³ and CRUISE.¹⁴ Their programs use several intelligent rules for splitting and usually produces a tree that often separates groups well. CART has been very successful in data mining applications (see Supplement 12A).

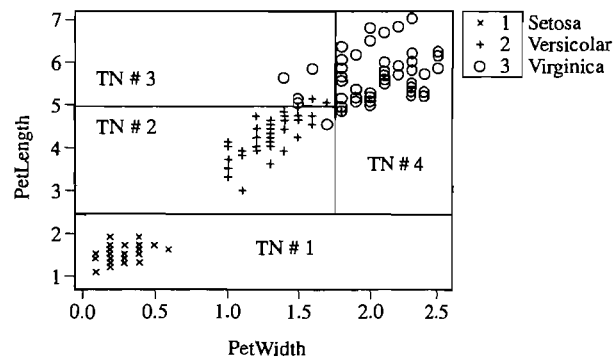


Figure 11.19 Classification tree terminal nodes (regions) in the petal width, petal length sample space.

¹³ Available for download at www.stat.wisc.edu/~loh/quest.html

¹⁴ Available for download at www.stat.wisc.edu/~loh/cruise.html

Neural Networks

A *neural network* (NN) is a computer-intensive, algorithmic procedure for transforming inputs into desired outputs using highly connected networks of relatively simple processing units (neurons or nodes). Neural networks are modeled after the neural activity in the human brain. The three essential features, then, of an NN are the basic computing units (neurons or nodes), the network architecture describing the connections between the computing units, and the training algorithm used to find values of the network parameters (weights) for performing a particular task.

The computing units are connected to one another in the sense that the output from one unit can serve as part of the input to another unit. Each computing unit transforms an input to an output using some prespecified function that is typically monotone, but otherwise arbitrary. This function depends on constants (parameters) whose values must be determined with a training set of inputs and outputs.

Network architecture is the organization of computing units and the types of connections permitted. In statistical applications, the computing units are arranged in a series of layers with connections between nodes in different layers, but not between nodes in the same layer. The layer receiving the initial inputs is called the input layer. The final layer is called the output layer. Any layers between the input and output layers are called hidden layers. A simple schematic representation of a multilayer NN is shown in Figure 11.20.

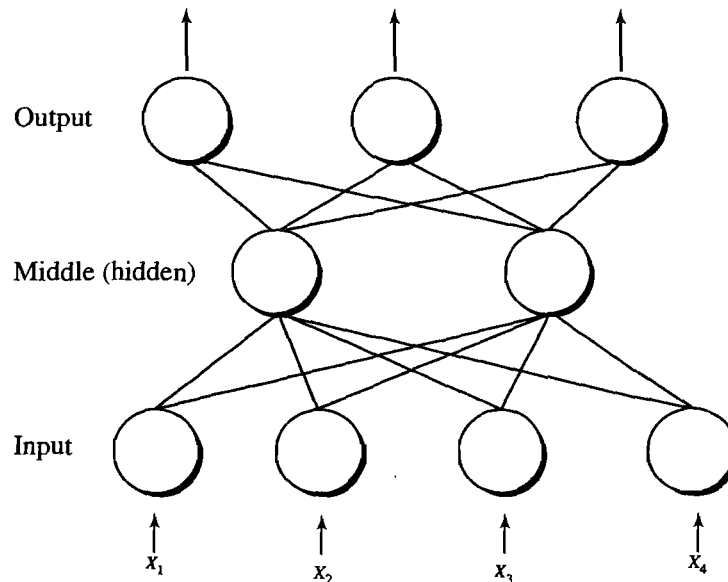


Figure 11.20 A neural network with one hidden layer.

Neural networks can be used for discrimination and classification. When they are so used, the input variables are the measured group characteristics X_1, X_2, \dots, X_p , and the output variables are categorical variables indicating group membership. Current practical experience indicates that properly constructed neural networks perform about as well as logistic regression and the discriminant functions we have discussed in this chapter. Reference [30] contains a good discussion of the use of neural networks in applied statistics.

Selection of Variables

In some applications of discriminant analysis, data are available on a large number of variables. Mucciardi and Gose [27] discuss a discriminant analysis based on 157 variables.¹⁵ In this case, it would obviously be desirable to select a relatively small subset of variables that would contain almost as much information as the original collection. This is the objective of *stepwise discriminant analysis*, and several popular commercial computer programs have such a capability.

If a stepwise discriminant analysis (or any variable selection method) is employed, the results should be interpreted with caution. (See [28].) There is no guarantee that the subset selected is “best,” regardless of the criterion used to make the selection. For example, subsets selected on the basis of minimizing the apparent error rate or maximizing “discriminatory power” may perform poorly in future samples. Problems associated with variable-selection procedures are magnified if there are large correlations among the variables or between linear combinations of the variables.

Choosing a subset of variables that seems to be optimal for a *given data set* is especially disturbing if classification is the objective. At the very least, the derived classification function should be evaluated with a validation sample. As Murray [28] suggests, a better idea might be to split the sample into a number of batches and determine the “best” subset for each batch. The number of times a given variable appears in the best subsets provides a measure of the worth of that variable for future classification.

Testing for Group Differences

We have pointed out, in connection with two group classification, that effective allocation is probably not possible unless the populations are well separated. The same is true for the many group situation. Classification is ordinarily not attempted, unless the population mean vectors differ significantly from one another. Assuming that the data are nearly multivariate normal, with a common covariance matrix, MANOVA can be performed to test for differences in the population mean vectors. Although apparent significant differences do not automatically imply effective classification, testing is a necessary first step. If no significant differences are found, constructing classification rules will probably be a waste of time.

¹⁵Imagine the problems of verifying the assumption of 157-variate normality and simultaneously estimating, for example, the 12,403 parameters of the 157×157 presumed common covariance matrix!

Graphics

Sophisticated computer graphics now allow one visually to examine multivariate data in two and three dimensions. Thus, groupings in the variable space for any choice of two or three variables can often be discerned by eye. In this way, potentially important classifying variables are often identified and outlying, or "atypical," observations revealed. Visual displays are important aids in discrimination and classification, and their use is likely to increase as the hardware and associated computer programs become readily available. Frequently, as much can be learned from a visual examination as by a complex numerical analysis.

Practical Considerations Regarding Multivariate Normality

The interplay between the choice of tentative assumptions and the form of the resulting classifier is important. Consider Figure 11.21, which shows the kidney-shaped density contours from two very nonnormal densities. In this case, the normal theory linear (or even quadratic) classification rule will be inadequate compared to another choice. That is, linear discrimination here is inappropriate.

Often discrimination is attempted with a large number of variables, some of which are of the presence-absence, or 0-1, type. In these situations and in others with restricted ranges for the variables, multivariate normality may not be a sensible assumption. As we have seen, classification based on Fisher's linear discriminants can be optimal from a minimum ECM or minimum TPM point of view only when multivariate normality holds. How are we to interpret these quantities when normality is clearly not viable?

In the absence of multivariate normality, Fisher's linear discriminants can be viewed as providing an approximation to the total sample information. The values of the first few discriminants themselves can be checked for normality and rule (11-67) employed. Since the discriminants are linear combinations of a large number of variables, they will often be nearly normal. Of course, one must keep in mind that the first few discriminants are an *incomplete* summary of the original sample information. Classification rules based on this restricted set may perform poorly, while optimal rules derived from all of the sample information may perform well.

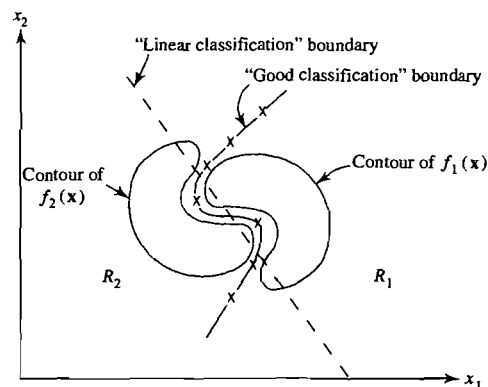


Figure 11.21 Two nonnormal populations for which linear discrimination is inappropriate.

EXERCISES

11.1. Consider the two data sets

$$\mathbf{X}_1 = \begin{bmatrix} 3 & 7 \\ 2 & 4 \\ 4 & 7 \end{bmatrix} \quad \text{and} \quad \mathbf{X}_2 = \begin{bmatrix} 6 & 9 \\ 5 & 7 \\ 4 & 8 \end{bmatrix}$$

for which

$$\bar{\mathbf{x}}_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \quad \bar{\mathbf{x}}_2 = \begin{bmatrix} 5 \\ 8 \end{bmatrix}$$

and

$$\mathbf{S}_{\text{pooled}} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

- (a) Calculate the linear discriminant function in (11-19).
 (b) Classify the observation $\mathbf{x}'_0 = [2 \ 7]$ as population π_1 or population π_2 , using Rule (11-18) with equal priors and equal costs.

- 11.2. (a) Develop a linear classification function for the data in Example 11.1 using (11-19).
 (b) Using the function in (a) and (11-20), construct the "confusion matrix" by classifying the given observations. Compare your classification results with those of Figure 11.1, where the classification regions were determined "by eye." (See Example 11.6.)
 (c) Given the results in (b), calculate the apparent error rate (APER).
 (d) State any assumptions you make to justify the use of the method in Parts a and b.

11.3. Prove Result 11.1.

Hint: Substituting the integral expressions for $P(2|1)$ and $P(1|2)$ given by (11-1) and (11-2), respectively, into (11-5) yields

$$\text{ECM} = c(2|1)p_1 \int_{R_2} f_1(\mathbf{x}) d\mathbf{x} + c(1|2)p_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x}$$

Noting that $\Omega = R_1 \cup R_2$, so that the total probability

$$1 = \int_{\Omega} f_1(\mathbf{x}) d\mathbf{x} = \int_{R_1} f_1(\mathbf{x}) d\mathbf{x} + \int_{R_2} f_1(\mathbf{x}) d\mathbf{x}$$

we can write

$$\text{ECM} = c(2|1)p_1 \left[1 - \int_{R_1} f_1(\mathbf{x}) d\mathbf{x} \right] + c(1|2)p_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x}$$

By the additive property of integrals (volumes),

$$\text{ECM} = \int_{R_1} [c(1|2)p_2 f_2(\mathbf{x}) - c(2|1)p_1 f_1(\mathbf{x})] d\mathbf{x} + c(2|1)p_1$$

Now, p_1 , p_2 , $c(1|2)$, and $c(2|1)$ are nonnegative. In addition, $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ are nonnegative for all \mathbf{x} and are the only quantities in ECM that depend on \mathbf{x} . Thus, ECM is minimized if R_1 includes those values \mathbf{x} for which the integrand

$$[c(1|2)p_2 f_2(\mathbf{x}) - c(2|1)p_1 f_1(\mathbf{x})] \leq 0$$

and excludes those \mathbf{x} for which this quantity is positive.

- 11.4. A researcher wants to determine a procedure for discriminating between two multivariate populations. The researcher has enough data available to estimate the density functions $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ associated with populations π_1 and π_2 , respectively. Let $c(2|1) = 50$ (this is the cost of assigning items as π_2 , given that π_1 is true) and $c(1|2) = 100$.

In addition, it is known that about 20% of all possible items (for which the measurements \mathbf{x} can be recorded) belong to π_2 .

- (a) Give the minimum ECM rule (in general form) for assigning a new item to one of the two populations.
 (b) Measurements recorded on a new item yield the density values $f_1(\mathbf{x}) = .3$ and $f_2(\mathbf{x}) = .5$. Given the preceding information, assign this item to population π_1 or population π_2 .

- 11.5. Show that

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \\ = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)$$

[see Equation (11-13).]

- 11.6. Consider the linear function $Y = \mathbf{a}'\mathbf{X}$. Let $E(\mathbf{X}) = \boldsymbol{\mu}_1$ and $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ if \mathbf{X} belongs to population π_1 . Let $E(\mathbf{X}) = \boldsymbol{\mu}_2$ and $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ if \mathbf{X} belongs to population π_2 . Let $m = \frac{1}{2}(\mu_{1Y} + \mu_{2Y}) = \frac{1}{2}(\mathbf{a}'\boldsymbol{\mu}_1 + \mathbf{a}'\boldsymbol{\mu}_2)$. Given that $\mathbf{a}' = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1}$, show each of the following.

(a) $E(\mathbf{a}'\mathbf{X}|\pi_1) - m = \mathbf{a}'\boldsymbol{\mu}_1 - m > 0$

(b) $E(\mathbf{a}'\mathbf{X}|\pi_2) - m = \mathbf{a}'\boldsymbol{\mu}_2 - m < 0$

Hint: Recall that $\boldsymbol{\Sigma}$ is of full rank and is positive definite, so $\boldsymbol{\Sigma}^{-1}$ exists and is positive definite.

- 11.7. Let $f_1(x) = (1 - |x|)$ for $|x| \leq 1$ and $f_2(x) = (1 - |x - .5|)$ for $-.5 \leq x \leq 1.5$.

- (a) Sketch the two densities.
 (b) Identify the classification regions when $p_1 = p_2$ and $c(1|2) = c(2|1)$.
 (c) Identify the classification regions when $p_1 = .2$ and $c(1|2) = c(2|1)$.

- 11.8. Refer to Exercise 11.7. Let $f_1(x)$ be the same as in that exercise, but take $f_2(x) = \frac{1}{4}(2 - |x - .5|)$ for $-1.5 \leq x \leq 2.5$.

- (a) Sketch the two densities.
 (b) Determine the classification regions when $p_1 = p_2$ and $c(1|2) = c(2|1)$.

- 11.9. For $g = 2$ groups, show that the ratio in (11-59) is proportional to the ratio

$$\frac{\left(\begin{array}{c} \text{squared distance} \\ \text{between means of } Y \end{array} \right)}{(\text{variance of } Y)} = \frac{(\mu_{1Y} - \mu_{2Y})^2}{\sigma_Y^2} = \frac{(\mathbf{a}'\boldsymbol{\mu}_1 - \mathbf{a}'\boldsymbol{\mu}_2)^2}{\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a}} \\ = \frac{\mathbf{a}'(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)'\mathbf{a}}{\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a}} = \frac{(\mathbf{a}'\boldsymbol{\delta})^2}{\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a}}$$

where $\boldsymbol{\delta} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ is the difference in mean vectors. This ratio is the population counterpart of (11-23). Show that the ratio is maximized by the linear combination

$$\mathbf{a} = c\boldsymbol{\Sigma}^{-1}\boldsymbol{\delta} = c\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

for any $c \neq 0$.

Hint: Note that $(\mu_i - \bar{\mu})(\mu_i - \bar{\mu})' = \frac{1}{4}(\mu_1 - \mu_2)(\mu_1 - \mu_2)'$ for $i = 1, 2$, where $\bar{\mu} = \frac{1}{2}(\mu_1 + \mu_2)$.

- 11.10.** Suppose that $n_1 = 11$ and $n_2 = 12$ observations are made on two random variables X_1 and X_2 , where X_1 and X_2 are assumed to have a bivariate normal distribution with a common covariance matrix Σ , but possibly different mean vectors μ_1 and μ_2 for the two samples. The sample mean vectors and pooled covariance matrix are

$$\bar{\mathbf{x}}_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}; \quad \bar{\mathbf{x}}_2 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$\mathbf{S}_{\text{pooled}} = \begin{bmatrix} 7.3 & -1.1 \\ -1.1 & 4.8 \end{bmatrix}$$

- (a) Test for the difference in population mean vectors using Hotelling's two-sample T^2 -statistic. Let $\alpha = .10$.
- (b) Construct Fisher's (sample) linear discriminant function. [See (11-19) and (11-25).]
- (c) Assign the observation $\mathbf{x}_0' = [0 \ 1]$ to either population π_1 or π_2 . Assume equal costs and equal prior probabilities.
- 11.11.** Suppose a univariate random variable X has a normal distribution with variance 4. If X is from population π_1 , its mean is 10; if it is from population π_2 , its mean is 14. Assume equal prior probabilities for the events $A1 = X$ is from population π_1 and $A2 = X$ is from population π_2 , and assume that the misclassification costs $c(2|1)$ and $c(1|2)$ are equal (for instance, \$10). We decide that we shall allocate (classify) X to population π_1 if $X \leq c$, for some c to be determined, and to population π_2 if $X > c$. Let $B1$ be the event X is classified into population π_1 and $B2$ be the event X is classified into population π_2 . Make a table showing the following: $P(B1|A2)$, $P(B2|A1)$, $P(A1 \text{ and } B2)$, $P(A2 \text{ and } B1)$, $P(\text{misclassification})$, and expected cost for various values of c . For what choice of c is expected cost minimized? The table should take the following form:

c	$P(B1 A2)$	$P(B2 A1)$	$P(A1 \text{ and } B2)$	$P(A2 \text{ and } B1)$	$P(\text{error})$	Expected cost
10						
\vdots						
14						

What is the value of the minimum expected cost?

- 11.12.** Repeat Exercise 11.11 if the prior probabilities of $A1$ and $A2$ are equal, but $c(2|1) = \$5$ and $c(1|2) = \$15$.
- 11.13.** Repeat Exercise 11.11 if the prior probabilities of $A1$ and $A2$ are $P(A1) = .25$ and $P(A2) = .75$ and the misclassification costs are as in Exercise 11.12.
- 11.14.** Consider the discriminant functions derived in Example 11.3. Normalize $\hat{\mathbf{a}}$ using (11-21) and (11-22). Compute the two midpoints \hat{m}_1^* and \hat{m}_2^* corresponding to the two choices of normalized vectors, say, $\hat{\mathbf{a}}_1^*$ and $\hat{\mathbf{a}}_2^*$. Classify $\mathbf{x}_0' = [-.210, -.044]$ with the function $\hat{y}_0^* = \hat{\mathbf{a}}^{*'} \mathbf{x}_0$ for the two cases. Are the results consistent with the classification obtained for the case of equal prior probabilities in Example 11.3? Should they be?
- 11.15.** Derive the expressions in (11-27) from (11-6) when $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ are multivariate normal densities with means μ_1, μ_2 and covariances Σ_1, Σ_2 , respectively.

11.16. Suppose \mathbf{x} comes from one of two populations:

π_1 : Normal with mean $\boldsymbol{\mu}_1$ and covariance matrix $\boldsymbol{\Sigma}_1$

π_2 : Normal with mean $\boldsymbol{\mu}_2$ and covariance matrix $\boldsymbol{\Sigma}_2$

If the respective density functions are denoted by $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$, find the expression for the quadratic discriminator

$$Q = \ln \left[\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \right]$$

If $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$, for instance, verify that Q becomes

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)$$

11.17. Suppose populations π_1 and π_2 are as follows:

	Population	
	π_1	π_2
Distribution	Normal	Normal
Mean $\boldsymbol{\mu}$	$[10, 15]'$	$[10, 25]'$
Covariance $\boldsymbol{\Sigma}$	$\begin{bmatrix} 18 & 12 \\ 12 & 32 \end{bmatrix}$	$\begin{bmatrix} 20 & -7 \\ -7 & 5 \end{bmatrix}$

Assume equal prior probabilities and misclassifications costs of $c(2|1) = \$10$ and $c(1|2) = \$73.89$. Find the posterior probabilities of populations π_1 and π_2 , $P(\pi_1|\mathbf{x})$ and $P(\pi_2|\mathbf{x})$, the value of the quadratic discriminator Q in Exercise 11.16, and the classification for each value of \mathbf{x} in the following table:

\mathbf{x}	$P(\pi_1 \mathbf{x})$	$P(\pi_2 \mathbf{x})$	Q	Classification
$[10, 15]'$				
$[12, 17]'$				
\vdots				
$[30, 35]'$				

(Note: Use an increment of 2 in each coordinate—11 points in all.)

Show each of the following on a graph of the x_1, x_2 plane.

- The mean of each population
- The ellipse of minimal area with probability .95 of containing \mathbf{x} for each population
- The region R_1 (for population π_1) and the region $\Omega - R_1 = R_2$ (for population π_2)
- The 11 points classified in the table

11.18. If \mathbf{B} is defined as $c(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)'$ for some constant c , verify that $\mathbf{e} = c\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ is in fact an (unscaled) eigenvector of $\boldsymbol{\Sigma}^{-1}\mathbf{B}$, where $\boldsymbol{\Sigma}$ is a covariance matrix.

11.19. (a) Using the original data sets \mathbf{X}_1 and \mathbf{X}_2 given in Example 11.7, calculate $\bar{\mathbf{x}}_i$, \mathbf{S}_i , $i = 1, 2$, and $\mathbf{S}_{\text{pooled}}$, verifying the results provided for these quantities in the example.

- (b) Using the calculations in Part a, compute Fisher's linear discriminant function, and use it to classify the sample observations according to Rule (11-25). Verify that the confusion matrix given in Example 11.7 is correct.
- (c) Classify the sample observations on the basis of smallest squared distance $D_i^2(\mathbf{x})$ of the observations from the group means $\bar{\mathbf{x}}_1$ and $\bar{\mathbf{x}}_2$. [See (11-54).] Compare the results with those in Part b. Comment.

11.20. The matrix identity (see Bartlett [3])

$$\mathbf{S}_{H,\text{pooled}}^{-1} = \frac{n-3}{n-2} \left(\mathbf{S}_{\text{pooled}}^{-1} + \frac{c_k}{1 - c_k(\mathbf{x}_H - \bar{\mathbf{x}}_k)' \mathbf{S}_{\text{pooled}}^{-1} (\mathbf{x}_H - \bar{\mathbf{x}}_k)} \cdot \mathbf{S}_{\text{pooled}}^{-1} (\mathbf{x}_H - \bar{\mathbf{x}}_k) (\mathbf{x}_H - \bar{\mathbf{x}}_k)' \mathbf{S}_{\text{pooled}}^{-1} \right)$$

where

$$c_k = \frac{n_k}{(n_k - 1)(n - 2)}$$

allows the calculation of $\mathbf{S}_{H,\text{pooled}}^{-1}$ from $\mathbf{S}_{\text{pooled}}^{-1}$. Verify this identity using the data from Example 11.7. Specifically, set $n = n_1 + n_2$, $k = 1$, and $\mathbf{x}'_H = [2, 12]$. Calculate $\mathbf{S}_{H,\text{pooled}}^{-1}$ using the full data $\mathbf{S}_{\text{pooled}}^{-1}$ and $\bar{\mathbf{x}}_1$, and compare the result with $\mathbf{S}_{H,\text{pooled}}^{-1}$ in Example 11.7.

11.21. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s > 0$ denote the $s \leq \min(g-1, p)$ nonzero eigenvalues of $\Sigma^{-1} \mathbf{B}_\mu$ and $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_s$ the corresponding eigenvectors (scaled so that $\mathbf{e}' \Sigma \mathbf{e} = 1$). Show that the vector of coefficients \mathbf{a} that maximizes the ratio

$$\frac{\mathbf{a}' \mathbf{B}_\mu \mathbf{a}}{\mathbf{a}' \Sigma \mathbf{a}} = \frac{\mathbf{a}' \left[\sum_{i=1}^g (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})' \right] \mathbf{a}}{\mathbf{a}' \Sigma \mathbf{a}}$$

is given by $\mathbf{a}_1 = \mathbf{e}_1$. The linear combination $\mathbf{a}'_1 \mathbf{X}$ is called the *first discriminant*. Show that the value $\mathbf{a}_2 = \mathbf{e}_2$ maximizes the ratio subject to $\text{Cov}(\mathbf{a}'_1 \mathbf{X}, \mathbf{a}'_2 \mathbf{X}) = 0$. The linear combination $\mathbf{a}'_2 \mathbf{X}$ is called the *second discriminant*. Continuing, $\mathbf{a}_k = \mathbf{e}_k$ maximizes the ratio subject to $0 = \text{Cov}(\mathbf{a}'_k \mathbf{X}, \mathbf{a}'_i \mathbf{X})$, $i < k$, and $\mathbf{a}'_k \mathbf{X}$ is called the *kth discriminant*. Also, $\text{Var}(\mathbf{a}'_i \mathbf{X}) = 1$, $i = 1, \dots, s$. [See (11-62) for the sample equivalent.]

Hint: We first convert the maximization problem to one already solved. By the spectral decomposition in (2-20), $\Sigma = \mathbf{P}' \Lambda \mathbf{P}$ where Λ is a diagonal matrix with positive elements λ_i . Let $\Lambda^{1/2}$ denote the diagonal matrix with elements $\sqrt{\lambda_i}$. By (2-22), the symmetric square-root matrix $\Sigma^{1/2} = \mathbf{P}' \Lambda^{1/2} \mathbf{P}$ and its inverse $\Sigma^{-1/2} = \mathbf{P}' \Lambda^{-1/2} \mathbf{P}$ satisfy $\Sigma^{1/2} \Sigma^{1/2} = \Sigma$, $\Sigma^{1/2} \Sigma^{-1/2} = \mathbf{I} = \Sigma^{-1/2} \Sigma^{1/2}$ and $\Sigma^{-1/2} \Sigma^{-1/2} = \Sigma^{-1}$. Next, set

$$\mathbf{u} = \Sigma^{1/2} \mathbf{a}$$

so $\mathbf{u}' \mathbf{u} = \mathbf{a}' \Sigma^{1/2} \Sigma^{1/2} \mathbf{a} = \mathbf{a}' \Sigma \mathbf{a}$ and $\mathbf{u}' \Sigma^{-1/2} \mathbf{B}_\mu \Sigma^{-1/2} \mathbf{u} = \mathbf{a}' \Sigma^{1/2} \Sigma^{-1/2} \mathbf{B}_\mu \Sigma^{-1/2} \Sigma^{1/2} \mathbf{a} = \mathbf{a}' \mathbf{B}_\mu \mathbf{a}$. Consequently, the problem reduces to maximizing

$$\frac{\mathbf{u}' \Sigma^{-1/2} \mathbf{B}_\mu \Sigma^{-1/2} \mathbf{u}}{\mathbf{u}' \mathbf{u}}$$

over \mathbf{u} . From (2-51), the maximum of this ratio is λ_1 , the largest eigenvalue of $\Sigma^{-1/2} \mathbf{B}_\mu \Sigma^{-1/2}$. This maximum occurs when $\mathbf{u} = \mathbf{e}_1$, the normalized eigenvector

associated with λ_1 . Because $\mathbf{e}_1 = \mathbf{u} = \Sigma^{1/2}\mathbf{a}_1$, or $\mathbf{a}_1 = \Sigma^{-1/2}\mathbf{e}_1$, $\text{Var}(\mathbf{a}_1'\mathbf{X}) = \mathbf{a}_1'\Sigma\mathbf{a}_1 = \mathbf{e}_1'\Sigma^{-1/2}\Sigma\Sigma^{-1/2}\mathbf{e}_1 = \mathbf{e}_1'\Sigma^{-1/2}\Sigma^{1/2}\Sigma^{1/2}\Sigma^{-1/2}\mathbf{e}_1 = \mathbf{e}_1'\mathbf{e}_1 = 1$. By (2-52), $\mathbf{u} \perp \mathbf{e}_1$ maximizes the preceding ratio when $\mathbf{u} = \mathbf{e}_2$, the normalized eigenvector corresponding to λ_2 . For this choice, $\mathbf{a}_2 = \Sigma^{-1/2}\mathbf{e}_2$, and $\text{Cov}(\mathbf{a}_2'\mathbf{X}, \mathbf{a}_1'\mathbf{X}) = \mathbf{a}_2'\Sigma\mathbf{a}_1 = \mathbf{e}_2'\Sigma^{-1/2}\Sigma\Sigma^{-1/2}\mathbf{e}_1 = \mathbf{e}_2'\mathbf{e}_1 = 0$, since $\mathbf{e}_2 \perp \mathbf{e}_1$. Similarly, $\text{Var}(\mathbf{a}_2'\mathbf{X}) = \mathbf{a}_2'\Sigma\mathbf{a}_2 = \mathbf{e}_2'\mathbf{e}_2 = 1$. Continue in this fashion for the remaining discriminants. Note that if λ and \mathbf{e} are an eigenvalue-eigenvector pair of $\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}$, then

$$\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}\mathbf{e} = \lambda\mathbf{e}$$

and multiplication on the left by $\Sigma^{-1/2}$ gives

$$\Sigma^{-1/2}\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}\mathbf{e} = \lambda\Sigma^{-1/2}\mathbf{e} \quad \text{or} \quad \Sigma^{-1}\mathbf{B}_\mu(\Sigma^{-1/2}\mathbf{e}) = \lambda(\Sigma^{-1/2}\mathbf{e})$$

Thus, $\Sigma^{-1}\mathbf{B}_\mu$ has the same eigenvalues as $\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}$, but the corresponding eigenvector is proportional to $\Sigma^{-1/2}\mathbf{e} = \mathbf{a}$, as asserted.

11.22. Show that $\Delta_S^2 = \lambda_1 + \lambda_2 + \cdots + \lambda_p = \lambda_1 + \lambda_2 + \cdots + \lambda_s$, where $\lambda_1, \lambda_2, \dots, \lambda_s$ are the nonzero eigenvalues of $\Sigma^{-1}\mathbf{B}_\mu$ (or $\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}$) and Δ_S^2 is given by (11-68). Also, show that $\lambda_1 + \lambda_2 + \cdots + \lambda_r$ is the resulting separation when only the first r discriminants, Y_1, Y_2, \dots, Y_r are used.

Hint: Let \mathbf{P} be the orthogonal matrix whose i th row \mathbf{e}_i' is the eigenvector of $\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}$ corresponding to the i th largest eigenvalue, $i = 1, 2, \dots, p$. Consider

$$\underset{(p \times 1)}{\mathbf{Y}} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_s \\ \vdots \\ Y_p \end{bmatrix} = \begin{bmatrix} \mathbf{e}_1'\Sigma^{-1/2}\mathbf{X} \\ \vdots \\ \mathbf{e}_s'\Sigma^{-1/2}\mathbf{X} \\ \vdots \\ \mathbf{e}_p'\Sigma^{-1/2}\mathbf{X} \end{bmatrix} = \mathbf{P}\Sigma^{-1/2}\mathbf{X}$$

Now, $\mu_{iY} = E(\mathbf{Y} | \pi_i) = \mathbf{P}\Sigma^{-1/2}\boldsymbol{\mu}_i$ and $\bar{\boldsymbol{\mu}}_Y = \mathbf{P}\Sigma^{-1/2}\bar{\boldsymbol{\mu}}$, so

$$\begin{aligned} (\boldsymbol{\mu}_{iY} - \bar{\boldsymbol{\mu}}_Y)'(\boldsymbol{\mu}_{iY} - \bar{\boldsymbol{\mu}}_Y) &= (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})'\Sigma^{-1/2}\mathbf{P}'\mathbf{P}\Sigma^{-1/2}(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}}) \\ &= (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})'\Sigma^{-1}(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}}) \end{aligned}$$

Therefore, $\Delta_S^2 = \sum_{i=1}^g (\boldsymbol{\mu}_{iY} - \bar{\boldsymbol{\mu}}_Y)'(\boldsymbol{\mu}_{iY} - \bar{\boldsymbol{\mu}}_Y)$. Using Y_1 , we have

$$\begin{aligned} \sum_{i=1}^g (\mu_{iY_1} - \bar{\mu}_{Y_1})^2 &= \sum_{i=1}^g \mathbf{e}_1'\Sigma^{-1/2}(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})'\Sigma^{-1/2}\mathbf{e}_1 \\ &= \mathbf{e}_1'\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}\mathbf{e}_1 = \lambda_1 \end{aligned}$$

because \mathbf{e}_1 has eigenvalue λ_1 . Similarly, Y_2 produces

$$\sum_{i=1}^g (\mu_{iY_2} - \bar{\mu}_{Y_2})^2 = \mathbf{e}_2'\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}\mathbf{e}_2 = \lambda_2$$

and Y_p produces

$$\sum_{i=1}^g (\mu_{iY_p} - \bar{\mu}_{Y_p})^2 = \mathbf{e}_p'\Sigma^{-1/2}\mathbf{B}_\mu\Sigma^{-1/2}\mathbf{e}_p = \lambda_p$$

Thus,

$$\begin{aligned}\Delta_S^2 &= \sum_{i=1}^g (\mu_{iY} - \bar{\mu}_Y)' (\mu_{iY} - \bar{\mu}_Y) \\ &= \sum_{i=1}^g (\mu_{iY_1} - \bar{\mu}_{Y_1})^2 + \sum_{i=1}^g (\mu_{iY_2} - \bar{\mu}_{Y_2})^2 + \cdots + \sum_{i=1}^g (\mu_{iY_p} - \bar{\mu}_{Y_p})^2 \\ &= \lambda_1 + \lambda_2 + \cdots + \lambda_p = \lambda_1 + \lambda_2 + \cdots + \lambda_s\end{aligned}$$

since $\lambda_{s+1} = \cdots = \lambda_p = 0$. If only the first r discriminants are used, their contribution to Δ_S^2 is $\lambda_1 + \lambda_2 + \cdots + \lambda_r$.

The following exercises require the use of a computer.

11.23. Consider the data given in Exercise 1.14.

- Check the marginal distributions of the x_i 's in both the multiple-sclerosis (MS) group and non-multiple-sclerosis (NMS) group for normality by graphing the corresponding observations as normal probability plots. Suggest appropriate data transformations if the normality assumption is suspect.
- Assume that $\Sigma_1 = \Sigma_2 = \Sigma$. Construct Fisher's linear discriminant function. Do all the variables in the discriminant function appear to be important? Discuss your answer. Develop a classification rule assuming equal prior probabilities and equal costs of misclassification.
- Using the results in (b), calculate the apparent error rate. If computing resources allow, calculate an estimate of the expected actual error rate using Lachenbruch's holdout procedure. Compare the two error rates.

11.24. Annual financial data are collected for bankrupt firms approximately 2 years prior to their bankruptcy and for financially sound firms at about the same time. The data on four variables, $X_1 = \text{CF/TD} = (\text{cash flow})/(\text{total debt})$, $X_2 = \text{NI/TA} = (\text{net income})/(\text{total assets})$, $X_3 = \text{CA/CL} = (\text{current assets})/(\text{current liabilities})$, and $X_4 = \text{CA/NS} = (\text{current assets})/(\text{net sales})$, are given in Table 11.4.

- Using a different symbol for each group, plot the data for the pairs of observations (x_1, x_2) , (x_1, x_3) and (x_1, x_4) . Does it appear as if the data are approximately bivariate normal for any of these pairs of variables?
- Using the $n_1 = 21$ pairs of observations (x_1, x_2) for bankrupt firms and the $n_2 = 25$ pairs of observations (x_1, x_2) for nonbankrupt firms, calculate the sample mean vectors \bar{x}_1 and \bar{x}_2 and the sample covariance matrices S_1 and S_2 .
- Using the results in (b) and assuming that both random samples are from bivariate normal populations, construct the classification rule (11-29) with $p_1 = p_2$ and $c(1|2) = c(2|1)$.
- Evaluate the performance of the classification rule developed in (c) by computing the apparent error rate (APER) from (11-34) and the estimated expected actual error rate \hat{E} (AER) from (11-36).
- Repeat Parts c and d, assuming that $p_1 = .05$, $p_2 = .95$, and $c(1|2) = c(2|1)$. Is this choice of prior probabilities reasonable? Explain.
- Using the results in (b), form the pooled covariance matrix S_{pooled} , and construct Fisher's sample linear discriminant function in (11-19). Use this function to classify the sample observations and evaluate the APER. Is Fisher's linear discriminant function a sensible choice for a classifier in this case? Explain.
- Repeat Parts b–e using the observation pairs (x_1, x_3) and (x_1, x_4) . Do some variables appear to be better classifiers than others? Explain.
- Repeat Parts b–e using observations on all four variables (X_1, X_2, X_3, X_4) .

Table 11.4 Bankruptcy Data

Row	$x_1 = \frac{CF}{TD}$	$x_2 = \frac{NI}{TA}$	$x_3 = \frac{CA}{CL}$	$x_4 = \frac{CA}{NS}$	Population $\pi_i, i = 1, 2$
1	-.45	-.41	1.09	.45	0
2	-.56	-.31	1.51	.16	0
3	.06	.02	1.01	.40	0
4	-.07	-.09	1.45	.26	0
5	-.10	-.09	1.56	.67	0
6	-.14	-.07	.71	.28	0
7	.04	.01	1.50	.71	0
8	-.06	-.06	1.37	.40	0
9	.07	-.01	1.37	.34	0
10	-.13	-.14	1.42	.44	0
11	-.23	-.30	.33	.18	0
12	.07	.02	1.31	.25	0
13	.01	.00	2.15	.70	0
14	-.28	-.23	1.19	.66	0
15	.15	.05	1.88	.27	0
16	.37	.11	1.99	.38	0
17	-.08	-.08	1.51	.42	0
18	.05	.03	1.68	.95	0
19	.01	-.00	1.26	.60	0
20	.12	.11	1.14	.17	0
21	-.28	-.27	1.27	.51	0
1	.51	.10	2.49	.54	1
2	.08	.02	2.01	.53	1
3	.38	.11	3.27	.35	1
4	.19	.05	2.25	.33	1
5	.32	.07	4.24	.63	1
6	.31	.05	4.45	.69	1
7	.12	.05	2.52	.69	1
8	-.02	.02	2.05	.35	1
9	.22	.08	2.35	.40	1
10	.17	.07	1.80	.52	1
11	.15	.05	2.17	.55	1
12	-.10	-.01	2.50	.58	1
13	.14	-.03	.46	.26	1
14	.14	.07	2.61	.52	1
15	.15	.06	2.23	.56	1
16	.16	.05	2.31	.20	1
17	.29	.06	1.84	.38	1
18	.54	.11	2.33	.48	1
19	-.33	-.09	3.01	.47	1
20	.48	.09	1.24	.18	1
21	.56	.11	4.29	.45	1
22	.20	.08	1.99	.30	1
23	.47	.14	2.92	.45	1
24	.17	.04	2.45	.14	1
25	.58	.04	5.06	.13	1

Legend: $\pi_1 = 0$: bankrupt firms; $\pi_2 = 1$: nonbankrupt firms.

Source: 1968, 1969, 1970, 1971, 1972 Moody's Industrial Manuals.

- 11.25.** The annual financial data listed in Table 11.4 have been analyzed by Johnson [19] with a view toward detecting influential observations in a discriminant analysis. Consider variables $X_1 = \text{CF/TD}$ and $X_3 = \text{CA/CL}$.
- Using the data on variables X_1 and X_3 , construct Fisher's linear discriminant function. Use this function to classify the sample observations and evaluate the APER. [See (11-25) and (11-34).] Plot the data and the discriminant line in the (x_1, x_3) coordinate system.
 - Johnson [19] has argued that the multivariate observations in rows 16 for bankrupt firms and 13 for sound firms are influential. Using the X_1, X_3 data, calculate Fisher's linear discriminant function with *only* data point 16 for bankrupt firms deleted. Repeat this procedure with *only* data point 13 for sound firms deleted. Plot the respective discriminant lines on the scatter in part a, and calculate the APERs, ignoring the deleted point in each case. Does deleting either of these multivariate observations make a difference? (Note that neither of the potentially influential data points is particularly "distant" from the center of its respective scatter.)
- 11.26.** Using the data in Table 11.4, define a binary response variable Z that assumes the value 0 if a firm is bankrupt and 1 if a firm is not bankrupt. Let $X = \text{CA/CL}$, and consider the straight-line regression of Z on X .
- Although a binary response variable does not meet the standard regression assumptions, consider using least squares to determine the fitted straight line for the X, Z data. Plot the fitted values for bankrupt firms as a dot diagram on the interval $[0, 1]$. Repeat this procedure for nonbankrupt firms and overlay the two dot diagrams. A reasonable discrimination rule is to predict that a firm will go bankrupt if its fitted value is closer to 0 than to 1. That is, the fitted value is less than .5. Similarly, a firm is predicted to be sound if its fitted value is greater than .5. Use this decision rule to classify the sample firms. Calculate the APER.
 - Repeat the analysis in Part a using all four variables, X_1, \dots, X_4 . Is there any change in the APER? Do data points 16 for bankrupt firms and 13 for nonbankrupt firms stand out as influential?
 - Perform a logistic regression using all four variables.
- 11.27.** The data in Table 11.5 contain observations on $X_2 = \text{sepal width}$ and $X_4 = \text{petal width}$ for samples from three species of iris. There are $n_1 = n_2 = n_3 = 50$ observations in each sample.
- Plot the data in the (x_2, x_4) variable space. Do the observations for the three groups appear to be bivariate normal?

Table 11.5 Data on Irises											
π_1 : <i>Iris setosa</i>				π_2 : <i>Iris versicolor</i>				π_3 : <i>Iris virginica</i>			
Sepal length x_1	Sepal width x_2	Petal length x_3	Petal width x_4	Sepal length x_1	Sepal width x_2	Petal length x_3	Petal width x_4	Sepal length x_1	Sepal width x_2	Petal length x_3	Petal width x_4
5.1	3.5	1.4	0.2	7.0	3.2	4.7	1.4	6.3	3.3	6.0	2.5
4.9	3.0	1.4	0.2	6.4	3.2	4.5	1.5	5.8	2.7	5.1	1.9
4.7	3.2	1.3	0.2	6.9	3.1	4.9	1.5	7.1	3.0	5.9	2.1
4.6	3.1	1.5	0.2	5.5	2.3	4.0	1.3	6.3	2.9	5.6	1.8
5.0	3.6	1.4	0.2	6.5	2.8	4.6	1.5	6.5	3.0	5.8	2.2
5.4	3.9	1.7	0.4	5.7	2.8	4.5	1.3	7.6	3.0	6.6	2.1

(continues on next page)

Table 11.5 (continued)

π_1 : <i>Iris setosa</i>				π_2 : <i>Iris versicolor</i>				π_3 : <i>Iris virginica</i>			
Sepal length x_1	Sepal width x_2	Petal length x_3	Petal width x_4	Sepal length x_1	Sepal width x_2	Petal length x_3	Petal width x_4	Sepal length x_1	Sepal width x_2	Petal length x_3	Petal width x_4
4.6	3.4	1.4	0.3	6.3	3.3	4.7	1.6	4.9	2.5	4.5	1.7
5.0	3.4	1.5	0.2	4.9	2.4	3.3	1.0	7.3	2.9	6.3	1.8
4.4	2.9	1.4	0.2	6.6	2.9	4.6	1.3	6.7	2.5	5.8	1.8
4.9	3.1	1.5	0.1	5.2	2.7	3.9	1.4	7.2	3.6	6.1	2.5
5.4	3.7	1.5	0.2	5.0	2.0	3.5	1.0	6.5	3.2	5.1	2.0
4.8	3.4	1.6	0.2	5.9	3.0	4.2	1.5	6.4	2.7	5.3	1.9
4.8	3.0	1.4	0.1	6.0	2.2	4.0	1.0	6.8	3.0	5.5	2.1
4.3	3.0	1.1	0.1	6.1	2.9	4.7	1.4	5.7	2.5	5.0	2.0
5.8	4.0	1.2	0.2	5.6	2.9	3.6	1.3	5.8	2.8	5.1	2.4
5.7	4.4	1.5	0.4	6.7	3.1	4.4	1.4	6.4	3.2	5.3	2.3
5.4	3.9	1.3	0.4	5.6	3.0	4.5	1.5	6.5	3.0	5.5	1.8
5.1	3.5	1.4	0.3	5.8	2.7	4.1	1.0	7.7	3.8	6.7	2.2
5.7	3.8	1.7	0.3	6.2	2.2	4.5	1.5	7.7	2.6	6.9	2.3
5.1	3.8	1.5	0.3	5.6	2.5	3.9	1.1	6.0	2.2	5.0	1.5
5.4	3.4	1.7	0.2	5.9	3.2	4.8	1.8	6.9	3.2	5.7	2.3
5.1	3.7	1.5	0.4	6.1	2.8	4.0	1.3	5.6	2.8	4.9	2.0
4.6	3.6	1.0	0.2	6.3	2.5	4.9	1.5	7.7	2.8	6.7	2.0
5.1	3.3	1.7	0.5	6.1	2.8	4.7	1.2	6.3	2.7	4.9	1.8
4.8	3.4	1.9	0.2	6.4	2.9	4.3	1.3	6.7	3.3	5.7	2.1
5.0	3.0	1.6	0.2	6.6	3.0	4.4	1.4	7.2	3.2	6.0	1.8
5.0	3.4	1.6	0.4	6.8	2.8	4.8	1.4	6.2	2.8	4.8	1.8
5.2	3.5	1.5	0.2	6.7	3.0	5.0	1.7	6.1	3.0	4.9	1.8
5.2	3.4	1.4	0.2	6.0	2.9	4.5	1.5	6.4	2.8	5.6	2.1
4.7	3.2	1.6	0.2	5.7	2.6	3.5	1.0	7.2	3.0	5.8	1.6
4.8	3.1	1.6	0.2	5.5	2.4	3.8	1.1	7.4	2.8	6.1	1.9
5.4	3.4	1.5	0.4	5.5	2.4	3.7	1.0	7.9	3.8	6.4	2.0
5.2	4.1	1.5	0.1	5.8	2.7	3.9	1.2	6.4	2.8	5.6	2.2
5.5	4.2	1.4	0.2	6.0	2.7	5.1	1.6	6.3	2.8	5.1	1.5
4.9	3.1	1.5	0.2	5.4	3.0	4.5	1.5	6.1	2.6	5.6	1.4
5.0	3.2	1.2	0.2	6.0	3.4	4.5	1.6	7.7	3.0	6.1	2.3
5.5	3.5	1.3	0.2	6.7	3.1	4.7	1.5	6.3	3.4	5.6	2.4
4.9	3.6	1.4	0.1	6.3	2.3	4.4	1.3	6.4	3.1	5.5	1.8
4.4	3.0	1.3	0.2	5.6	3.0	4.1	1.3	6.0	3.0	4.8	1.8
5.1	3.4	1.5	0.2	5.5	2.5	4.0	1.3	6.9	3.1	5.4	2.1
5.0	3.5	1.3	0.3	5.5	2.6	4.4	1.2	6.7	3.1	5.6	2.4
4.5	2.3	1.3	0.3	6.1	3.0	4.6	1.4	6.9	3.1	5.1	2.3
4.4	3.2	1.3	0.2	5.8	2.6	4.0	1.2	5.8	2.7	5.1	1.9
5.0	3.5	1.6	0.6	5.0	2.3	3.3	1.0	6.8	3.2	5.9	2.3
5.1	3.8	1.9	0.4	5.6	2.7	4.2	1.3	6.7	3.3	5.7	2.5
4.8	3.0	1.4	0.3	5.7	3.0	4.2	1.2	6.7	3.0	5.2	2.3
5.1	3.8	1.6	0.2	5.7	2.9	4.2	1.3	6.3	2.5	5.0	1.9
4.6	3.2	1.4	0.2	6.2	2.9	4.3	1.3	6.5	3.0	5.2	2.0
5.3	3.7	1.5	0.2	5.1	2.5	3.0	1.1	6.2	3.4	5.4	2.3
5.0	3.3	1.4	0.2	5.7	2.8	4.1	1.3	5.9	3.0	5.1	1.8

Source: Anderson [1].

- (b) Assume that the samples are from bivariate normal populations with a common covariance matrix. Test the hypothesis $H_0: \mu_1 = \mu_2 = \mu_3$ versus H_1 : at least one μ_i is different from the others at the $\alpha = .05$ significance level. Is the assumption of a common covariance matrix reasonable in this case? Explain.
- (c) Assuming that the populations are bivariate normal, construct the quadratic discriminate scores $\hat{d}_i^Q(\mathbf{x})$ given by (11-47) with $p_1 = p_2 = p_3 = \frac{1}{3}$. Using Rule (11-48), classify the new observation $\mathbf{x}'_0 = [3.5 \ 1.75]$ into population π_1, π_2 , or π_3 .
- (d) Assume that the covariance matrices Σ_i are the same for all three bivariate normal populations. Construct the linear discriminate score $\hat{d}_i(\mathbf{x})$ given by (11-51), and use it to assign $\mathbf{x}'_0 = [3.5 \ 1.75]$ to one of the populations $\pi_i, i = 1, 2, 3$ according to (11-52). Take $p_1 = p_2 = p_3 = \frac{1}{3}$. Compare the results in Parts c and d. Which approach do you prefer? Explain.
- (e) Assuming equal covariance matrices and bivariate normal populations, and supposing that $p_1 = p_2 = p_3 = \frac{1}{3}$, allocate $\mathbf{x}'_0 = [3.5 \ 1.75]$ to π_1, π_2 , or π_3 using Rule (11-56). Compare the result with that in Part d. Delineate the classification regions \hat{R}_1, \hat{R}_2 , and \hat{R}_3 on your graph from Part a determined by the linear functions $\hat{d}_{ki}(\mathbf{x}_0)$ in (11-56).
- (f) Using the linear discriminant scores from Part d, classify the sample observations. Calculate the APER and $\hat{E}(\text{AER})$. (To calculate the latter, you should use Lachenbruch's holdout procedure. [See (11-57).])

11.28. Darroch and Mosimann [6] have argued that the three species of iris indicated in Table 11.5 can be discriminated on the basis of "shape" or scale-free information alone. Let $Y_1 = X_1/X_2$ be sepal shape and $Y_2 = X_3/X_4$ be petal shape.

- (a) Plot the data in the $(\log Y_1, \log Y_2)$ variable space. Do the observations for the three groups appear to be bivariate normal?
- (b) Assuming equal covariance matrices and bivariate normal populations, and supposing that $p_1 = p_2 = p_3 = \frac{1}{3}$, construct the linear discriminant scores $\hat{d}_i(\mathbf{x})$ given by (11-51) using both variables $\log Y_1, \log Y_2$ and each variable individually. Calculate the APERs.
- (c) Using the linear discriminant functions from Part b, calculate the holdout estimates of the expected AERs, and fill in the following summary table:

Variable(s)	Misclassification rate
$\log Y_1$	
$\log Y_2$	
$\log Y_1, \log Y_2$	

Compare the preceding misclassification rates with those in the summary tables in Example 11.12. Does it appear as if information on shape alone is an effective discriminator for these species of iris?

- (d) Compare the corresponding error rates in Parts b and c. Given the scatter plot in Part a, would you expect these rates to differ much? Explain.

11.29. The GPA and GMAT data alluded to in Example 11.11 are listed in Table 11.6.

- (a) Using these data, calculate $\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2, \bar{\mathbf{x}}_3, \bar{\mathbf{x}}$, and $\mathbf{S}_{\text{pooled}}$ and thus verify the results for these quantities given in Example 11.11.

Table 11.6 Admission Data for Graduate School of Business								
π_1 : Admit			π_2 : Do not admit			π_3 : Borderline		
Applicant no.	GPA (x_1)	GMAT (x_2)	Applicant no.	GPA (x_1)	GMAT (x_2)	Applicant no.	GPA (x_1)	GMAT (x_2)
1	2.96	596	32	2.54	446	60	2.86	494
2	3.14	473	33	2.43	425	61	2.85	496
3	3.22	482	34	2.20	474	62	3.14	419
4	3.29	527	35	2.36	531	63	3.28	371
5	3.69	505	36	2.57	542	64	2.89	447
6	3.46	693	37	2.35	406	65	3.15	313
7	3.03	626	38	2.51	412	66	3.50	402
8	3.19	663	39	2.51	458	67	2.89	485
9	3.63	447	40	2.36	399	68	2.80	444
10	3.59	588	41	2.36	482	69	3.13	416
11	3.30	563	42	2.66	420	70	3.01	471
12	3.40	553	43	2.68	414	71	2.79	490
13	3.50	572	44	2.48	533	72	2.89	431
14	3.78	591	45	2.46	509	73	2.91	446
15	3.44	692	46	2.63	504	74	2.75	546
16	3.48	528	47	2.44	336	75	2.73	467
17	3.47	552	48	2.13	408	76	3.12	463
18	3.35	520	49	2.41	469	77	3.08	440
19	3.39	543	50	2.55	538	78	3.03	419
20	3.28	523	51	2.31	505	79	3.00	509
21	3.21	530	52	2.41	489	80	3.03	438
22	3.58	564	53	2.19	411	81	3.05	399
23	3.33	565	54	2.35	321	82	2.85	483
24	3.40	431	55	2.60	394	83	3.01	453
25	3.38	605	56	2.55	528	84	3.03	414
26	3.26	664	57	2.72	399	85	3.04	446
27	3.60	609	58	2.85	381			
28	3.37	559	59	2.90	384			
29	3.80	521						
30	3.76	646						
31	3.24	467						

- (b) Calculate \mathbf{W}^{-1} and \mathbf{B} and the eigenvalues and eigenvectors of $\mathbf{W}^{-1}\mathbf{B}$. Use the linear discriminants derived from these eigenvectors to classify the new observation $\mathbf{x}_0 = [3.21 \quad 497]$ into one of the populations π_1 : admit; π_2 : not admit; and π_3 : borderline. Does the classification agree with that in Example 11.11? Should it? Explain.

11.30. Gerrild and Lantz [13] chemically analyzed crude-oil samples from three zones of sandstone:

- π_1 : Wilhelm
 π_2 : Sub-Mulinia
 π_3 : Upper

The values of the trace elements

- X_1 = vanadium (in percent ash)
 X_2 = iron (in percent ash)
 X_3 = beryllium (in percent ash)

and two measures of hydrocarbons,

X_4 = saturated hydrocarbons (in percent area)

X_5 = aromatic hydrocarbons (in percent area)

are presented for 56 cases in Table 11.7. The last two measurements are determined from areas under a gas-liquid chromatography curve.

- Obtain the estimated minimum TPM rule, assuming normality. Comment on the adequacy of the assumption of normality.
- Determine the estimate of $E(AER)$ using Lachenbruch's holdout procedure. Also give the confusion matrix.
- Consider various transformations of the data to normality (see Example 11.14), and repeat Parts a and b.

	x_1	x_2	x_3	x_4	x_5
π_1	3.9	51.0	0.20	7.06	12.19
	2.7	49.0	0.07	7.14	12.23
	2.8	36.0	0.30	7.00	11.30
	3.1	45.0	0.08	7.20	13.01
	3.5	46.0	0.10	7.81	12.63
	3.9	43.0	0.07	6.25	10.42
	2.7	35.0	0.00	5.11	9.00
π_2	5.0	47.0	0.07	7.06	6.10
	3.4	32.0	0.20	5.82	4.69
	1.2	12.0	0.00	5.54	3.15
	8.4	17.0	0.07	6.31	4.55
	4.2	36.0	0.50	9.25	4.95
	4.2	35.0	0.50	5.69	2.22
	3.9	41.0	0.10	5.63	2.94
	3.9	36.0	0.07	6.19	2.27
	7.3	32.0	0.30	8.02	12.92
	4.4	46.0	0.07	7.54	5.76
π_3	3.0	30.0	0.00	5.12	10.77
	6.3	13.0	0.50	4.24	8.27
	1.7	5.6	1.00	5.69	4.64
	7.3	24.0	0.00	4.34	2.99
	7.8	18.0	0.50	3.92	6.09
	7.8	25.0	0.70	5.39	6.20
	7.8	26.0	1.00	5.02	2.50
	9.5	17.0	0.05	3.52	5.71
	7.7	14.0	0.30	4.65	8.63
	11.0	20.0	0.50	4.27	8.40
	8.0	14.0	0.30	4.32	7.87
	8.4	18.0	0.20	4.38	7.98

(continues on next page)

Table 11.7 (continued)					
	x_1	x_2	x_3	x_4	x_5
	10.0	18.0	0.10	3.06	7.67
	7.3	15.0	0.05	3.76	6.84
	9.5	22.0	0.30	3.98	5.02
	8.4	15.0	0.20	5.02	10.12
	8.4	17.0	0.20	4.42	8.25
	9.5	25.0	0.50	4.44	5.95
	7.2	22.0	1.00	4.70	3.49
	4.0	12.0	0.50	5.71	6.32
	6.7	52.0	0.50	4.80	3.20
	9.0	27.0	0.30	3.69	3.30
	7.8	29.0	1.50	6.72	5.75
	4.5	41.0	0.50	3.33	2.27
	6.2	34.0	0.70	7.56	6.93
	5.6	20.0	0.50	5.07	6.70
	9.0	17.0	0.20	4.39	8.33
	8.4	20.0	0.10	3.74	3.77
	9.5	19.0	0.50	3.72	7.37
	9.0	20.0	0.50	5.97	11.17
	6.2	16.0	0.05	4.23	4.18
	7.3	20.0	0.50	4.39	3.50
	3.6	15.0	0.70	7.00	4.82
	6.2	34.0	0.07	4.84	2.37
	7.3	22.0	0.00	4.13	2.70
	4.1	29.0	0.70	5.78	7.76
	5.4	29.0	0.20	4.64	2.65
	5.0	34.0	0.70	4.21	6.50
	6.2	27.0	0.30	3.97	2.97

11.31. Refer to the data on salmon in Table 11.2.

- Plot the bivariate data for the two groups of salmon. Are the sizes and orientation of the scatters roughly the same? Do bivariate normal distributions with a common covariance matrix appear to be viable population models for the Alaskan and Canadian salmon?
- Using a linear discriminant function for two normal populations with equal priors and equal costs [see (11-19)], construct dot diagrams of the discriminant scores for the two groups. Does it appear as if the growth ring diameters separate for the two groups reasonably well? Explain.
- Repeat the analysis in Example 11.8 for the male and female salmon separately. Is it easier to discriminate Alaskan male salmon from Canadian male salmon than it is to discriminate the females in the two groups? Is gender (male or female) likely to be a useful discriminatory variable?

11.32. Data on hemophilia A carriers, similar to those used in Example 11.3, are listed in Table 11.8 on page 664. (See [15].) Using these data,

- Investigate the assumption of bivariate normality for the two groups.

Table 11.8 Hemophilia Data

Noncarriers (π_1)			Obligatory carriers (π_2)		
Group	\log_{10} (AHF activity)	\log_{10} (AHF antigen)	Group	\log_{10} (AHF activity)	\log_{10} (AHF antigen)
1	-.0056	-.1657	2	-.3478	.1151
1	-.1698	-.1585	2	-.3618	-.2008
1	-.3469	-.1879	2	-.4986	-.0860
1	-.0894	.0064	2	-.5015	-.2984
1	-.1679	.0713	2	-.1326	.0097
1	-.0836	.0106	2	-.6911	-.3390
1	-.1979	-.0005	2	-.3608	.1237
1	-.0762	.0392	2	-.4535	-.1682
1	-.1913	-.2123	2	-.3479	-.1721
1	-.1092	-.1190	2	-.3539	.0722
1	-.5268	-.4773	2	-.4719	-.1079
1	-.0842	.0248	2	-.3610	-.0399
1	-.0225	-.0580	2	-.3226	.1670
1	.0084	.0782	2	-.4319	-.0687
1	-.1827	-.1138	2	-.2734	-.0020
1	.1237	.2140	2	-.5573	.0548
1	-.4702	-.3099	2	-.3755	-.1865
1	-.1519	-.0686	2	-.4950	-.0153
1	.0006	-.1153	2	-.5107	-.2483
1	-.2015	-.0498	2	-.1652	.2132
1	-.1932	-.2293	2	-.2447	-.0407
1	.1507	.0933	2	-.4232	-.0998
1	-.1259	-.0669	2	-.2375	.2876
1	-.1551	-.1232	2	-.2205	.0046
1	-.1952	-.1007	2	-.2154	-.0219
1	.0291	.0442	2	-.3447	.0097
1	-.2228	-.1710	2	-.2540	-.0573
1	-.0997	-.0733	2	-.3778	-.2682
1	-.1972	-.0607	2	-.4046	-.1162
1	-.0867	-.0560	2	-.0639	.1569
			2	-.3351	-.1368
			2	-.0149	.1539
			2	-.0312	.1400
			2	-.1740	-.0776
			2	-.1416	.1642
			2	-.1508	.1137
			2	-.0964	.0531
			2	-.2642	.0867
			2	-.0234	.0804
			2	-.3352	.0875
			2	-.1878	.2510
			2	-.1744	.1892
			2	-.4055	-.2418
			2	-.2444	.1614
			2	-.4784	.0282

Source: See [15].

- (b) Obtain the sample linear discriminant function, assuming equal prior probabilities, and estimate the error rate using the holdout procedure.
- (c) Classify the following 10 new cases using the discriminant function in Part b.
- (d) Repeat Parts a–c, assuming that the prior probability of obligatory carriers (group 2) is $\frac{1}{4}$ and that of noncarriers (group 1) is $\frac{3}{4}$.

New Cases Requiring Classification		
Case	$\log_{10}(\text{AHF activity})$	$\log_{10}(\text{AHF antigen})$
1	-.112	-.279
2	-.059	-.068
3	.064	.012
4	-.043	-.052
5	-.050	-.098
6	-.094	-.113
7	-.123	-.143
8	-.011	-.037
9	-.210	-.090
10	-.126	-.019

11.33. Consider the data on bulls in Table 1.10.

- (a) Using the variables YrHgt, FtFrBody, PrctFFB, Frame, BkFat, SaleHt, and SaleWt, calculate Fisher's linear discriminants, and classify the bulls as Angus, Hereford, or Simmental. Calculate an estimate of $E(\text{AER})$ using the holdout procedure. Classify a bull with characteristics YrHgt = 50, FtFrBody = 1000, PrctFFB = 73, Frame = 7, BkFat = .17, SaleHt = 54, and SaleWt = 1525 as one of the three breeds. Plot the discriminant scores for the bulls in the two-dimensional discriminant space using different plotting symbols to identify the three groups.
- (b) Is there a subset of the original seven variables that is almost as good for discriminating among the three breeds? Explore this possibility by computing the estimated $E(\text{AER})$ for various subsets.

11.34. Table 11.9 on pages 666–667 contains data on breakfast cereals produced by three different American manufacturers: General Mills (G), Kellogg (K), and Quaker (Q). Assuming multivariate normal data with a common covariance matrix, equal costs, and equal priors, classify the cereal brands according to manufacturer. Compute the estimated $E(\text{AER})$ using the holdout procedure. Interpret the coefficients of the discriminant functions. Does it appear as if some manufacturers are associated with more “nutritional” cereals (high protein, low fat, high fiber, low sugar, and so forth) than others? Plot the cereals in the two-dimensional discriminant space, using different plotting symbols to identify the three manufacturers.

11.35. Table 11.10 on page 668 contains measurements on the gender, age, tail length (mm), and snout to vent length (mm) for Concho Water Snakes.

Define the variables

$$\begin{aligned} X_1 &= \text{Gender} \\ X_2 &= \text{Age} \\ X_3 &= \text{TailLength} \\ X_4 &= \text{SntoVnLength} \end{aligned}$$

Brand	Manufacturer	Calories	Protein	Fat	Sodium	Fiber	Carbohydrates	Sugar	Potassium	Group
1 Apple_Cinnamon_Cheerios	G	110	2	2	180	1.5	10.5	10	70	1
2 Cheerios	G	110	6	2	290	2.0	17.0	1	105	1
3 Cocoa_Puffs	G	110	1	1	180	0.0	12.0	13	55	1
4 Count_Chocula	G	110	1	1	180	0.0	12.0	13	65	1
5 Golden_Grahams	G	110	1	1	280	0.0	15.0	9	45	1
6 Honey_Nut_Cheerios	G	110	3	1	250	1.5	11.5	10	90	1
7 Kix	G	110	2	1	260	0.0	21.0	3	40	1
8 Lucky_Charm	G	110	2	1	180	0.0	12.0	12	55	1
9 Multi_Grain_Cheerios	G	100	2	1	220	2.0	15.0	6	90	1
10 Oatmeal_Raisin_Crisp	G	130	3	2	170	1.5	13.5	10	120	1
11 Raisin_Nut_Bran	G	100	3	2	140	2.5	10.5	8	140	1
12 Total_Corn_Flakes	G	110	2	1	200	0.0	21.0	3	35	1
13 Total_Raisin_Bran	G	140	3	1	190	4.0	15.0	14	230	1
14 Total_Whole_Grain	G	100	3	1	200	3.0	16.0	3	110	1
15 Trix	G	110	1	1	140	0.0	13.0	12	25	1
16 Wheaties	G	100	3	1	200	3.0	17.0	3	110	1
17 Wheaties_Honey_Gold	G	110	2	1	200	1.0	16.0	8	60	1
18 All_Bran	K	70	4	1	260	9.0	7.0	5	320	2
19 Apple_Jacks	K	110	2	0	125	1.0	11.0	14	30	2
20 Corn_Flakes	K	100	2	0	290	1.0	21.0	2	35	2
21 Corn_Pops	K	110	1	0	90	1.0	13.0	12	20	2

continued

22 Cracklin' Oat Bran	K	110	3	3	140	4.0	10.0	7	160	2
23 Crispix	K	110	2	0	220	1.0	21.0	3	30	2
24 Froot Loops	K	110	2	1	125	1.0	11.0	13	30	2
25 Frosted Flakes	K	110	1	0	200	1.0	14.0	11	25	2
26 Frosted Mini Wheats	K	100	3	0	0	3.0	14.0	7	100	2
27 Fruitful Bran	K	120	3	0	240	5.0	14.0	12	190	2
28 Just Right Crunchy Nuggets	K	110	2	1	170	1.0	17.0	6	60	2
29 Mueslix Crispy Blend	K	160	3	2	150	3.0	17.0	13	160	2
30 Nut&Honey Crunch	K	120	2	1	190	0.0	15.0	9	40	2
31 Nutri-grain Almond-Raisin	K	140	3	2	220	3.0	21.0	7	130	2
32 Nutri-grain Wheat	K	90	3	0	170	3.0	18.0	2	90	2
33 Product 19	K	100	3	0	320	1.0	20.0	3	45	2
34 Raisin Bran	K	120	3	1	210	5.0	14.0	12	240	2
35 Rice Krispies	K	110	2	0	290	0.0	22.0	3	35	2
36 Smacks	K	110	2	1	70	1.0	9.0	15	40	2
37 Special K	K	110	6	0	230	1.0	16.0	3	55	2
38 Cap'n Crunch	Q	120	1	2	220	0.0	12.0	12	35	3
39 Honey Graham Ohs	Q	120	1	2	220	1.0	12.0	11	45	3
40 Life	Q	100	4	2	150	2.0	12.0	6	95	3
41 Puffed Rice	Q	50	1	0	0	0.0	13.0	0	15	3
42 Puffed Wheat	Q	50	2	0	0	1.0	10.0	0	50	3
43 Quaker Oatmeal	Q	100	5	2	0	2.7	1.0	1	110	3

Source: Data courtesy of Chad Dacus.

Table 11.10 Concho Water Snake Data

	Gender	Age	TailLength	Snto VnLength		Gender	Age	TailLength	Snto VnLength
1	Female	2	127	441	1	Male	2	126	457
2	Female	2	171	455	2	Male	2	128	466
3	Female	2	171	462	3	Male	2	151	466
4	Female	2	164	446	4	Male	2	115	361
5	Female	2	165	463	5	Male	2	138	473
6	Female	2	127	393	6	Male	2	145	477
7	Female	2	162	451	7	Male	3	145	507
8	Female	2	133	376	8	Male	3	145	493
9	Female	2	173	475	9	Male	3	158	558
10	Female	2	145	398	10	Male	3	152	495
11	Female	2	154	435	11	Male	3	159	521
12	Female	3	165	491	12	Male	3	138	487
13	Female	3	178	485	13	Male	3	166	565
14	Female	3	169	477	14	Male	3	168	585
15	Female	3	186	530	15	Male	3	160	550
16	Female	3	170	478	16	Male	4	181	652
17	Female	3	182	511	17	Male	4	185	587
18	Female	3	172	475	18	Male	4	172	606
19	Female	3	182	487	19	Male	4	180	591
20	Female	3	172	454	20	Male	4	205	683
21	Female	3	183	502	21	Male	4	175	625
22	Female	3	170	483	22	Male	4	182	612
23	Female	3	171	477	23	Male	4	185	618
24	Female	3	181	493	24	Male	4	181	613
25	Female	3	167	490	25	Male	4	167	600
26	Female	3	175	493	26	Male	4	167	602
27	Female	3	139	477	27	Male	4	160	596
28	Female	3	183	501	28	Male	4	165	611
29	Female	4	198	537	29	Male	4	173	603
30	Female	4	190	566					
31	Female	4	192	569					
32	Female	4	211	574					
33	Female	4	206	570					
34	Female	4	206	573					
35	Female	4	165	531					
36	Female	4	189	528					
37	Female	4	195	536					

Source: Data courtesy of Raymond J. Carroll.

- (a) Plot the data as a scatter plot with tail length (x_3) as the horizontal axis and snout to vent length (x_4) as the vertical axis. Use different plotting symbols for female and male snakes, and different symbols for different ages. Does it appear as if tail length and snout to vent length might usefully discriminate the genders of snakes? The different ages of snakes?
- (b) Assuming multivariate normal data with a common covariance matrix, equal priors, and equal costs, classify the Concho Water Snakes according to gender. Compute the estimated $E(\text{AER})$ using the holdout procedure.

- (c) Repeat part (b) using age as the groups rather than gender.
- (d) Repeat part (b) using only snout to vent length to classify the snakes according to age. Compare the results with those in part (c). Can effective classification be achieved with only a single variable in this case? Explain.
- 11.36.** Refer to Example 11.17. Using logistic regression, refit the salmon data in Table 11.2 with only the covariates freshwater growth and marine growth. Check for the significance of the model and the significance of each individual covariate. Set $\alpha = .05$. Use the fitted function to classify each of the observations in Table 11.2 as Alaskan salmon or Canadian salmon using rule (11-77). Compute the apparent error rate, APER, and compare this error rate with the error rate from the linear classification function discussed in Example 11.8.

References

1. Anderson, E. "The Irises of the Gaspé Peninsula." *Bulletin of the American Iris Society*, **59** (1939), 2-5.
2. Anderson, T. W. *An Introduction to Multivariate Statistical Analysis* (3rd ed.). New York: John Wiley, 2003.
3. Bartlett, M. S. "An Inverse Matrix Adjustment Arising in Discriminant Analysis." *Annals of Mathematical Statistics*, **22** (1951), 107-111.
4. Bouma, B. N., et al. "Evaluation of the Detection Rate of Hemophilia Carriers." *Statistical Methods for Clinical Decision Making*, **7**, no. 2 (1975), 339-350.
5. Breiman, L., J. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Belmont, CA: Wadsworth, Inc., 1984.
6. Darroch, J. N., and J. E. Mosimann. "Canonical and Principal Components of Shape." *Biometrika*, **72**, no. 1 (1985), 241-252.
7. Efron, B. "The Efficiency of Logistic Regression Compared to Normal Discriminant Analysis." *Journal of the American Statistical Association*, **81** (1975), 321-327.
8. Eisenbeis, R. A. "Pitfalls in the Application of Discriminant Analysis in Business, Finance and Economics." *Journal of Finance*, **32**, no. 3 (1977), 875-900.
9. Fisher, R. A. "The Use of Multiple Measurements in Taxonomic Problems." *Annals of Eugenics*, **7** (1936), 179-188.
10. Fisher, R. A. "The Statistical Utilization of Multiple Measurements." *Annals of Eugenics*, **8** (1938), 376-386.
11. Ganesalingam, S. "Classification and Mixture Approaches to Clustering via Maximum Likelihood." *Applied Statistics*, **38**, no. 3 (1989), 455-466.
12. Geisser, S. "Discrimination, Allocatory and Separatory, Linear Aspects." In *Classification and Clustering*, edited by J. Van Ryzin, pp. 301-330. New York: Academic Press, 1977.
13. Gerrild, P. M., and R. J. Lantz. "Chemical Analysis of 75 Crude Oil Samples from Pliocene Sand Units, Elk Hills Oil Field, California." *U.S. Geological Survey Open-File Report*, 1969.
14. Gnanadesikan, R. *Methods for Statistical Data Analysis of Multivariate Observations* (2nd ed.). New York: Wiley-Interscience, 1997.
15. Habbema, J. D. F., J. Hermans, and K. Van Den Broek. "A Stepwise Discriminant Analysis Program Using Density Estimation." In *Compstat 1974, Proc. Computational Statistics*, pp. 101-110. Vienna: Physica, 1974.

16. Hills, M. "Allocation Rules and Their Error Rates." *Journal of the Royal Statistical Society (B)*, **28** (1966), 1–31.
17. Hosmer, D. W. and S. Lemeshow. *Applied Logistic Regression* (2nd ed.). New York: Wiley-Interscience, 2000.
18. Hudlet, R., and R. A. Johnson. "Linear Discrimination and Some Further Results on Best Lower Dimensional Representations." In *Classification and Clustering*, edited by J. Van Ryzin, pp. 371–394. New York: Academic Press, 1977.
19. Johnson, W. "The Detection of Influential Observations for Allocation, Separation, and the Determination of Probabilities in a Bayesian Framework." *Journal of Business and Economic Statistics*, **5**, no. 3 (1987), 369–381.
20. Kendall, M. G. *Multivariate Analysis*. New York: Hafner Press, 1975.
21. Kim, H. and Loh, W. Y., "Classification Trees with Unbiased Multiway Splits," *Journal of the American Statistical Association*, **96**, (2001), 589–604.
22. Krzanowski, W. J. "The Performance of Fisher's Linear Discriminant Function under Non-Optimal Conditions." *Technometrics*, **19**, no. 2 (1977), 191–200.
23. Lachenbruch, P. A. *Discriminant Analysis*. New York: Hafner Press, 1975.
24. Lachenbruch, P. A., and M. R. Mickey. "Estimation of Error Rates in Discriminant Analysis." *Technometrics*, **10**, no. 1 (1968), 1–11.
25. Loh, W. Y. and Shih, Y. S., "Split Selection Methods for Classification Trees," *Statistica Sinica*, **7**, (1997), 815–840.
26. McCullagh, P., and J. A. Nelder. *Generalized Linear Models* (2nd ed.). London: Chapman and Hall, 1989.
27. Mucciardi, A. N., and E. E. Gose. "A Comparison of Seven Techniques for Choosing Subsets of Pattern Recognition Properties." *IEEE Trans. Computers*, **C20** (1971), 1023–1031.
28. Murray, G. D. "A Cautionary Note on Selection of Variables in Discriminant Analysis." *Applied Statistics*, **26**, no. 3 (1977), 246–250.
29. Rencher, A. C. "Interpretation of Canonical Discriminant Functions, Canonical Variates and Principal Components." *The American Statistician*, **46** (1992), 217–225.
30. Stern, H. S. "Neural Networks in Applied Statistics." *Technometrics*, **38**, (1996), 205–214.
31. Wald, A. "On a Statistical Problem Arising in the Classification of an Individual into One of Two Groups." *Annals of Mathematical Statistics*, **15** (1944), 145–162.
32. Welch, B. L. "Note on Discriminant Functions." *Biometrika*, **31** (1939), 218–220.