# Chapter 7
# Introduction to Discriminant Procedures

## Chapter Table of Contents

# Chapter 7
# Introduction to
## Discriminant Procedures

## Overview

The SAS procedures for discriminant analysis treat data with one classification variable and several quantitative variables. The purpose of discriminant analysis can be to find one or more of the following:

- a mathematical rule, or *discriminant function*, for guessing to which class an observation belongs, based on knowledge of the quantitative variables only
- a set of linear combinations of the quantitative variables that best reveals the differences among the classes
- a subset of the quantitative variables that best reveals the differences among the classes

The SAS discriminant procedures are as follows:

DISCRIM     computes various discriminant functions for classifying observations. Linear or quadratic discriminant functions can be used for data with approximately multivariate normal within-class distributions. Nonparametric methods can be used without making any assumptions about these distributions.

CANDISC     performs a canonical analysis to find linear combinations of the quantitative variables that best summarize the differences among the classes.

STEPDISC    uses forward selection, backward elimination, or stepwise selection to try to find a subset of quantitative variables that best reveals differences among the classes.

## Background

The term *discriminant analysis* (Fisher 1936; Cooley and Lohnes 1971; Tatsuoka 1971; Kshirsagar 1972; Lachenbruch 1975, 1979; Gnanadesikan 1977; Klecka 1980; Hand 1981,1982; Silverman, 1986) refers to several different types of analysis. Classificatory discriminant analysis is used to classify observations into two or more known groups on the basis of one or more quantitative variables. Classification can be done by either a parametric method or a nonparametric method in the DISCRIM procedure. A parametric method is appropriate only for approximately normal within-class distributions. The method generates either a linear discriminant function (the

within-class covariance matrices are assumed to be equal) or a quadratic discriminant function (the within-class covariance matrices are assumed to be unequal).

When the distribution within each group is not assumed to have any specific distribution or is assumed to have a distribution different from the multivariate normal distribution, nonparametric methods can be used to derive classification criteria. These methods include the kernel method and nearest-neighbor methods. The kernel method uses uniform, normal, Epanechnikov, biweight, or triweight kernels in estimating the group-specific density at each observation. The within-group covariance matrices or the pooled covariance matrix can be used to scale the data.

The performance of a discriminant function can be evaluated by estimating error rates (probabilities of misclassification). Error count estimates and posterior probability error rate estimates can be evaluated with PROC DISCRIM. When the input data set is an ordinary SAS data set, the error rates can also be estimated by cross validation.

In multivariate statistical applications, the data collected are largely from distributions different from the normal distribution. Various forms of nonnormality can arise, such as qualitative variables or variables with underlying continuous but nonnormal distributions. If the multivariate normality assumption is violated, the use of parametric discriminant analysis may not be appropriate. When a parametric classification criterion (linear or quadratic discriminant function) is derived from a nonnormal population, the resulting error rate estimates may be biased.

If your quantitative variables are not normally distributed, or if you want to classify observations on the basis of categorical variables, you should consider using the CATMOD or LOGISTIC procedure to fit a categorical linear model with the classification variable as the dependent variable. Press and Wilson (1978) compare logistic regression and parametric discriminant analysis and conclude that logistic regression is preferable to parametric discriminant analysis in cases for which the variables do not have multivariate normal distributions within classes. However, if you do have normal within-class distributions, logistic regression is less efficient than parametric discriminant analysis. Efron (1975) shows that with two normal populations having a common covariance matrix, logistic regression is between one half and two thirds as effective as the linear discriminant function in achieving asymptotically the same error rate.

Do not confuse discriminant analysis with cluster analysis. All varieties of discriminant analysis require prior knowledge of the classes, usually in the form of a sample from each class. In cluster analysis, the data do not include information on class membership; the purpose is to construct a classification. See Chapter 8, "Introduction to Clustering Procedures."

Canonical discriminant analysis is a dimension-reduction technique related to principal components and canonical correlation, and it can be performed by both the CANDISC and DISCRIM procedures. A discriminant criterion is always derived in PROC DISCRIM. If you want canonical discriminant analysis without the use of a discriminant criterion, you should use PROC CANDISC. Stepwise discriminant analysis is a variable-selection technique implemented by the STEPDISC procedure. After selecting a subset of variables with PROC STEPDISC, use any of the other dis-

criminant procedures to obtain more detailed analyses. PROC CANDISC and PROC STEPDISC perform hypothesis tests that require the within-class distributions to be approximately normal, but these procedures can be used descriptively with nonnormal data.

Another alternative to discriminant analysis is to perform a series of univariate oneway ANOVAs. All three discriminant procedures provide summaries of the univariate ANOVAs. The advantage of the multivariate approach is that two or more classes that overlap considerably when each variable is viewed separately may be more distinct when examined from a multivariate point of view.

## Example: Contrasting Univariate and Multivariate Analyses

Consider the two classes indicated by 'H' and 'O' in Figure 7.1. The results are shown in Figure 7.2.

```
data random;
   drop n;

   Group = 'H';
   do n = 1 to 20;
      X = 4.5 + 2 * normal(57391);
      Y = X + .5 + normal(57391);
      output;
   end;

   Group = 'O';
   do n = 1 to 20;
      X = 6.25 + 2 * normal(57391);
      Y = X - 1 + normal(57391);
      output;
   end;

run;

symbol1 v='H' c=blue;
symbol2 v='O' c=yellow;
proc gplot;
   plot Y*X=Group / cframe=ligr nolegend;
run;

proc candisc anova;
   class Group;
   var X Y;
run;
```
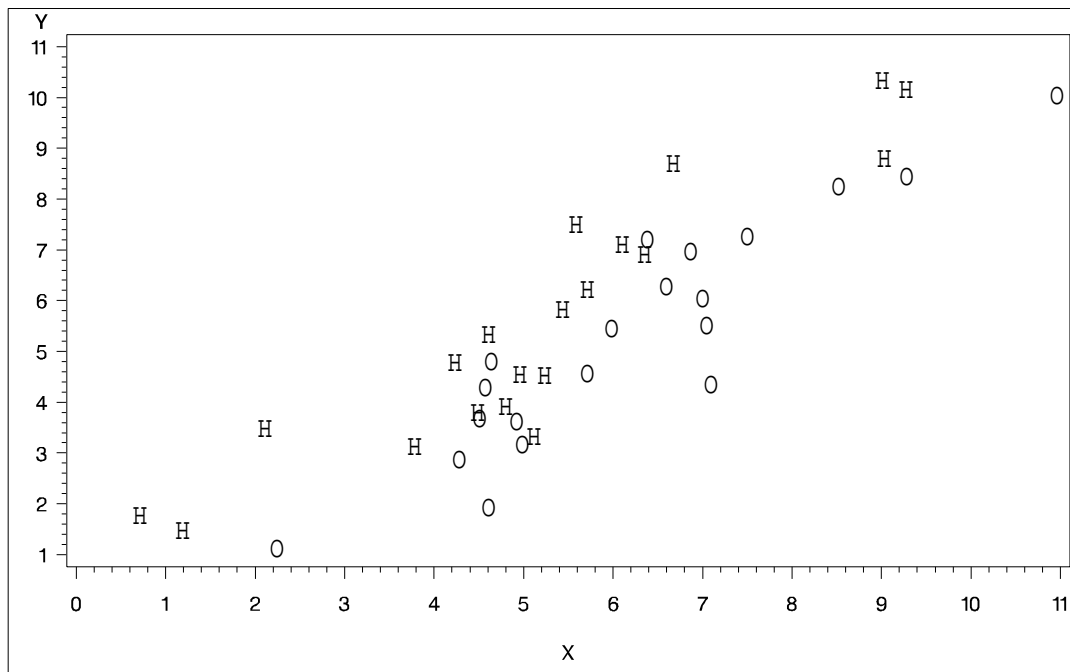
**Figure 7.1.**   Groups for Contrasting Univariate and Multivariate Analyses

```
                        The CANDISC Procedure

        Observations        40        DF Total              39
        Variables            2        DF Within Classes     38
        Classes              2        DF Between Classes      1


                     Class Level Information

                 Variable
        Group    Name        Frequency        Weight     Proportion

        H        H                  20       20.0000       0.500000
        O        O                  20       20.0000       0.500000
```

**Figure 7.2.**   Contrasting Univariate and Multivariate Analyses

```
                       The CANDISC Procedure

                     Univariate Test Statistics

                 F Statistics,    Num DF=1,   Den DF=38

            Total     Pooled    Between
          Standard   Standard   Standard            R-Square
Variable  Deviation  Deviation  Deviation  R-Square  / (1-RSq)  F Value  Pr > F

X          2.1776     2.1498     0.6820    0.0503     0.0530     2.01   0.1641
Y          2.4215     2.4486     0.2047    0.0037     0.0037     0.14   0.7105


                        Average R-Square

                Unweighted              0.0269868
                Weighted by Variance    0.0245201


           Multivariate Statistics and Exact F Statistics

                     S=1    M=0    N=17.5

Statistic                     Value    F Value    Num DF    Den DF    Pr > F

Wilks' Lambda              0.64203704   10.31       2         37      0.0003
Pillai's Trace             0.35796296   10.31       2         37      0.0003
Hotelling-Lawley Trace     0.55754252   10.31       2         37      0.0003
Roy's Greatest Root        0.55754252   10.31       2         37      0.0003
```

```
                       The CANDISC Procedure

                     Adjusted     Approximate        Squared
            Canonical   Canonical     Standard      Canonical
            Correlation Correlation     Error       Correlation

     1       0.598300    0.589467     0.102808       0.357963

                    Eigenvalues of Inv(E)*H
                      = CanRsq/(1-CanRsq)

          Eigenvalue   Difference   Proportion   Cumulative

     1       0.5575                   1.0000       1.0000

        Test of H0: The canonical correlations in the
          current row and all that follow are zero

          Likelihood    Approximate
              Ratio       F Value    Num DF   Den DF    Pr > F

     1     0.64203704      10.31       2        37      0.0003

              NOTE: The F statistic is exact.
```

```
                    The CANDISC Procedure

              Total Canonical Structure

              Variable                Can1

              X                  -0.374883
              Y                   0.101206


              Between Canonical Structure

              Variable                Can1

              X                  -1.000000
              Y                   1.000000


          Pooled Within Canonical Structure

              Variable                Can1

              X                  -0.308237
              Y                   0.081243
```

```
                    The CANDISC Procedure

        Total-Sample Standardized Canonical Coefficients

              Variable                Can1

              X               -2.625596855
              Y                2.446680169


     Pooled Within-Class Standardized Canonical Coefficients

              Variable                Can1

              X               -2.592150014
              Y                2.474116072


              Raw Canonical Coefficients

              Variable                Can1

              X               -1.205756217
              Y                1.010412967


          Class Means on Canonical Variables

               Group                Can1

               H           0.7277811475
               O          -.7277811475
```

The univariate $R^2$s are very small, 0.0503 for X and 0.0037 for Y, and neither variable shows a significant difference between the classes at the 0.10 level.

The multivariate test for differences between the classes is significant at the 0.0003 level. Thus, the multivariate analysis has found a highly significant difference, whereas the univariate analyses failed to achieve even the 0.10 level. The Raw Canonical Coefficients for the first canonical variable, Can1, show that the classes differ most widely on the linear combination -1.205756217 X + 1.010412967 Y or approximately Y - 1.2 X. The $R^2$ between Can1 and the class variable is 0.357963 as given by the Squared Canonical Correlation, which is much higher than either univariate $R^2$.

In this example, the variables are highly correlated within classes. If the within-class correlation were smaller, there would be greater agreement between the univariate and multivariate analyses.

# References

Cooley, W.W. and Lohnes, P.R. (1971), *Multivariate Data Analysis*, New York: John Wiley & Sons, Inc.

Dillion, W. and Goldstein, M. (1984), *Multivariate Analysis: Methods and Applications*, New York: John Wiley & Sons, Inc.

Efron, B. (1975), "The Efficiency of Logistic Regression Compared to Normal Discriminant Analysis," *Journal of the American Statistical Association*, 70, 892–898.

Fisher, R.A. (1936), "The Use of Multiple Measurements in Taxonomic Problems," *Annals of Eugenics*, 7, 179–188.

Gnanadesikan, R. (1977), *Methods for Statistical Data Analysis of Multivariate Observations*, New York: John Wiley & Sons, Inc.

Hand, D.J. (1981), *Discrimination and Classification*, New York: John Wiley & Sons, Inc.

Hand, D.J. (1982), *Kernel Discriminant Analysis*, New York: Research Studies Press.

Hora, S.C. and Wilcox, J.B. (1982), "Estimation of Error Rates in Several-Population Discriminant Analysis," *Journal of Marketing Research*, XIX, 57–61.

Klecka, W.R. (1980), *Discriminant Analysis*, Sage University Paper Series on Quantitative Applications in the Social Sciences, 07-019. Beverly Hills, CA: Sage Publications.

Kshirsagar, A.M. (1972), *Multivariate Analysis*, New York: Marcel Dekker, Inc.

Lachenbruch, P.A. (1975), *Discriminant Analysis*, New York: Hafner.

Lachenbruch, P.A. (1979), "Discriminant Analysis," *Biometrics*, 35, 69–85.

Press, S.J. and Wilson, S. (1978), "Choosing Between Logistic Regression and Discriminant Analysis," *Journal of the American Statistical Association*, 73, 699–705.

Silverman, B.W. (1986), *Density Estimation for Statistics and Data Analysis*, New York: Chapman and Hall.

Tatsuoka, M.M. (1971), *Multivariate Analysis*, New York: John Wiley & Sons, Inc.