



# Interactive Stepwise Discriminant Analysis \*

Dimitar Vandev

## ABSTRACT

ldagui.m is an interactive tool for linear and quadratic discriminant analysis. The reason for developing such a tool is the inconsistency of conventional statistical programs in following aspects:

- treating missing data;
- interaction with the user;
- testing the quality of obtained model.

\*Supported by contracts:PRO-ENBIS: GTC1-2001-43031 and WINE DB: G6RD-CT-2001-00646

Sozopol 22-28.06.03

D. Vandev

Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



# Interactive Stepwise Discriminant Analysis \*

Dimitar Vandev

## ABSTRACT

ldagui.m is an interactive tool for linear and quadratic discriminant analysis. The reason for developing such a tool is the inconsistency of conventional statistical programs in following aspects:

- treating missing data;
- interaction with the user;
- testing the quality of obtained model.

\*Supported by contracts:PRO-ENBIS: GTC1-2001-43031 and WINE DB: G6RD-CT-2001-00646

Sozopol 22-28.06.03

D. Vandev

Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

# 1. Overview

Discriminant analysis (DA) is a very popular tool in applied statistics.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1. Overview

Discriminant analysis (DA) is a very popular tool in applied statistics.

The program `Ldagui.m` is developed in the frame of MATLAB. It is used with the help of menus, shortcuts, listboxes and a slider.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1. Overview

Discriminant analysis (DA) is a very popular tool in applied statistics.

The program `Ldagui.m` is developed in the frame of MATLAB. It is used with the help of menus, shortcuts, listboxes and a slider.

- First we will shortly outline the mathematics behind DA.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1. Overview

Discriminant analysis (DA) is a very popular tool in applied statistics.

The program `Ldagui.m` is developed in the frame of MATLAB. It is used with the help of menus, shortcuts, listboxes and a slider.

- First we will shortly outline the mathematics behind DA.
- Then we will describe menus and shortcuts of the program.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1. Overview

Discriminant analysis (DA) is a very popular tool in applied statistics.

The program `Ldagui.m` is developed in the frame of MATLAB. It is used with the help of menus, shortcuts, listboxes and a slider.

- First we will shortly outline the mathematics behind DA.
- Then we will describe menus and shortcuts of the program.
- Finally a small demonstration will be done to illustrate other features of the program.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

Let suppose we have observed two random variables:

1. **continuous**  $\xi \in R^p$ ;





Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

Let suppose we have observed two random variables:

1. **continuous**  $\xi \in R^p$ ;
2. **discrete** (or categorical)  $\eta$  with values in  $1, 2, \dots, G$ .



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

Let suppose we have observed two random variables:

1. **continuous**  $\xi \in R^p$ ;
2. **discrete** (or categorical)  $\eta$  with values in  $1, 2, \dots, G$ .
3. they have joined distribution (DA model):
  - $\mathbf{P}(\eta = g) = p(g)$
  - Conditional distribution of  $\xi \in R^p$  given  $\eta = g$  is described by the density  $f(x, m(g), C(g))$ .



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

Let suppose we have observed two random variables:

1. **continuous**  $\xi \in R^p$ ;
2. **discrete** (or categorical)  $\eta$  with values in  $1, 2, \dots, G$ .
3. they have joined distribution (DA model):
  - $\mathbf{P}(\eta = g) = p(g)$
  - Conditional distribution of  $\xi \in R^p$  given  $\eta = g$  is described by the density  $f(x, m(g), C(g))$ .

Here  $f$  is the density of Gauss distribution in  $R^p$  described by two parameters:

- mean -  $m(g)$ ;
- covariance -  $C(g)$ ,



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.1. Bayes formula

Suppose we know the parameters of this model:

1. The **prior** probabilities -  $\{p(g)\}$ ;



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.1. Bayes formula

Suppose we know the parameters of this model:

1. The **prior** probabilities -  $\{p(g)\}$ ;
2. **Group means** -  $\{m(g)\}$ ;



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.1. Bayes formula

Suppose we know the parameters of this model:

1. The **prior** probabilities -  $\{p(g)\}$ ;
2. **Group means** -  $\{m(g)\}$ ;
3. Within group **covariance matrices** -  $C(g)$ ;

That is, the set of numbers:  $\{p_g, m_g, C(g), g = 1, 2, \dots, G\}$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.1. Bayes formula

Suppose we know the parameters of this model:

1. The **prior** probabilities -  $\{p(g)\}$ ;
2. **Group means** -  $\{m(g)\}$ ;
3. Within group **covariance matrices** -  $C(g)$ ;

That is, the set of numbers:  $\{p_g, m_g, C(g), g = 1, 2, \dots, G\}$

Then according of the famous formula of Bayes we may write down the conditional probability of  $\eta = g$  given  $x$ :

$$\mathbf{P}(\eta = g | \xi = g) = q(g|x) = c(x).p(g).f(x, m(g), C(g)),$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.1. Bayes formula

Suppose we know the parameters of this model:

1. The **prior** probabilities -  $\{p(g)\}$ ;
2. **Group means** -  $\{m(g)\}$ ;
3. Within group **covariance matrices** -  $C(g)$ ;

That is, the set of numbers:  $\{p_g, m_g, C(g), g = 1, 2, \dots, G\}$

Then according of the famous formula of Bayes we may write down the conditional probability of  $\eta = g$  given  $x$ :

$$\mathbf{P}(\eta = g | \xi = g) = q(g|x) = c(x) \cdot p(g) \cdot f(x, m(g), C(g)), \quad (1)$$

where  $c$  is a normalizing constant, such that  $\sum q(g|x) = 1$ .





Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.1. Bayes formula

Suppose we know the parameters of this model:

1. The **prior** probabilities -  $\{p(g)\}$ ;
2. **Group means** -  $\{m(g)\}$ ;
3. Within group **covariance matrices** -  $C(g)$ ;

That is, the set of numbers:  $\{p_g, m_g, C(g), g = 1, 2, \dots, G\}$

Then according of the famous formula of Bayes we may write down the conditional probability of  $\eta = g$  given  $x$ :

$$\mathbf{P}(\eta = g | \xi = g) = q(g|x) = c(x) \cdot p(g) \cdot f(x, m(g), C(g)), \quad (1)$$

where  $c$  is a normalizing constant, such that  $\sum q(g|x) = 1$ .

We call this probability **posterior** and say that the observation  $x$  belongs to the group  $g$  with probability  $q(g|x)$ .



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.1. Bayes formula

Suppose we know the parameters of this model:

1. The **prior** probabilities -  $\{p(g)\}$ ;
2. **Group means** -  $\{m(g)\}$ ;
3. Within group **covariance matrices** -  $C(g)$ ;

That is, the set of numbers:  $\{p_g, m_g, C(g), g = 1, 2, \dots, G\}$

Then according of the famous formula of Bayes we may write down the conditional probability of  $\eta = g$  given  $x$ :

$$\mathbf{P}(\eta = g | \xi = g) = q(g|x) = c(x) \cdot p(g) \cdot f(x, m(g), C(g)), \quad (1)$$

where  $c$  is a normalizing constant, such that  $\sum q(g|x) = 1$ .

We call this probability **posterior** and say that the observation  $x$  belongs to the group  $g$  with probability  $q(g|x)$ .

According the maximum likelihood principle the classification rule should then be:

$$\hat{g} = \underset{h}{\operatorname{argmax}} : q(h). \quad (2)$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

According the maximum likelihood principle the classification rule should then be:

$$\hat{g} = \underset{h}{\operatorname{argmax}} : q(h). \quad (2)$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.2. Linear Discriminant Analysis

Suppose that within group covariance  $C(g)$  are equal:

$$C(g) = C, \quad (g = 1, 2, \dots, G)$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.2. Linear Discriminant Analysis

Suppose that within group covariance  $C(g)$  are equal:

$$C(g) = C, \quad (g = 1, 2, \dots, G) \quad (3)$$

Then the maximum likelihood rule (2) becomes a set of inequalities:

$$p(\hat{g}) \cdot f(x, m(\hat{g}), C) \geq p(h) \cdot f(x, m(h), C), \quad (h = 1, 2, \dots, G).$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page

◀▶

◀▶

Go Back

Full Screen

Close

Quit

## 1.2. Linear Discriminant Analysis

Suppose that within group covariance  $C(g)$  are equal:

$$C(g) = C, \quad (g = 1, 2, \dots, G) \quad (3)$$

Then the maximum likelihood rule (2) becomes a set of inequalities:

$$p(\hat{g}) \cdot f(x, m(\hat{g}), C) \geq p(h) \cdot f(x, m(h), C), \quad (h = 1, 2, \dots, G). \quad (4)$$

or (what is the same) to:

$$b(\hat{g})'x + a(\hat{g}) \geq b(h)'x + a(h), \quad (h = 1, 2, \dots, G), \quad (5)$$

We decide that the observation  $x$  belongs to the group  $g$ , if for each  $h$  the inequality (5) holds:

$$L_g(x) \geq L_h(x), \quad (h = 1, 2, \dots, G), \quad (6)$$

The functions  $L$  are called **discriminant functions**.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

The vector  $b(g)$  and the number  $a(g)$  in this case are calculated explicitly:

$$b(h) = m(h)'C^{-1}, \quad a(h) = \log p(h) - m(h)'C^{-1}m(h). \quad (7)$$

This is why DA takes the name **Linear** - the **discriminant functions** are **linear** functions.





Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

The vector  $b(g)$  and the number  $a(g)$  in this case are calculated explicitly:

$$b(h) = m(h)'C^{-1}, \quad a(h) = \log p(h) - m(h)'C^{-1}m(h). \quad (7)$$

This is why DA takes the name **Linear** - the **discriminant functions** are **linear** functions.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 1.3. Quadratic Discriminant Analysis

When the assumption (3):  $C(g) = C$  is not appropriate, the corresponding functions become quadratic.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 1.3. Quadratic Discriminant Analysis

When the assumption (3):  $C(g) = C$  is not appropriate, the corresponding functions become quadratic.

If one has equal prior probabilities  $p(h) = 1/G$ , the solution of the classification problem (2) is equivalent to the minimization of so called **Mahalanobis distances** of the observation to the group means:

$$h(x, g) = (x - m(g))'C_g^{-1}(x - m(g))$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 1.3. Quadratic Discriminant Analysis

When the assumption (3):  $C(g) = C$  is not appropriate, the corresponding functions become quadratic.

If one has equal prior probabilities  $p(h) = 1/G$ , the solution of the classification problem (2) is equivalent to the minimization of so called **Mahalanobis distances** of the observation to the group means:

$$h(x, g) = (x - m(g))' C_g^{-1} (x - m(g)) \quad (8)$$

One uses Mahalanobis distances (8) to classify the observation to the closest group:

$$\hat{g} = \underset{h}{\operatorname{argmin}} h(x, h).$$

In general however, the Bayes rule (1) is better.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.4. Estimation

Let the training sample consists of vectors  $(g_i, x_i), i = 1.2. \dots, n.$

Denote by  $I(g)$  the set  $\{i : g_i = g\}$  and let  $n(g) = |I(g)|.$

First the standard calculations - averages:

$$m(g) = \frac{1}{n(g)} \sum_{i \in I(g)} x_i.$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.4. Estimation

Let the training sample consists of vectors  $(g_i, x_i), i = 1, 2, \dots, n$ .

Denote by  $I(g)$  the set  $\{i : g_i = g\}$  and let  $n(g) = |I(g)|$ .

First the standard calculations - averages:

$$m(g) = \frac{1}{n(g)} \sum_{i \in I(g)} x_i. \quad (9)$$

Cross products:

$$SS(g) = \sum_{i \in I(g)} (x_i - m(g))(x_i - m(g))', \quad SS_{in} = \sum_{g \in G} SS(g)$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page

◀▶

◀▶

Go Back

Full Screen

Close

Quit

## 1.4. Estimation

Let the training sample consists of vectors  $(g_i, x_i), i = 1.2. \dots, n$ .

Denote by  $I(g)$  the set  $\{i : g_i = g\}$  and let  $n(g) = |I(g)|$ .

First the standard calculations - averages:

$$m(g) = \frac{1}{n(g)} \sum_{i \in I(g)} x_i. \quad (9)$$

Cross products:

$$SS(g) = \sum_{i \in I(g)} (x_i - m(g))(x_i - m(g))', \quad SS_{in} = \sum_{g \in G} SS(g) \quad (10)$$

Standard maximum likelihood estimates are:

$$C(g) = \frac{1}{n(g) - 1} SS(g), \quad C = \frac{1}{n - G} SS_{in}. \quad (11)$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

We propose to **correct** the within group covariance estimate considering instead the mixture:

$$C_g := (1 - \alpha) * C + \alpha * C_g. \quad (12)$$

The parameter  $0 \leq \alpha \leq 1$  is to be chosen in interactive way via slider. Such corrections are not new (see for example (Lauter, 159-168) in (Fedorov, 1992)).





Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

We propose to **correct** the within group covariance estimate considering instead the mixture:

$$C_g := (1 - \alpha) * C + \alpha * C_g. \quad (12)$$

The parameter  $0 \leq \alpha \leq 1$  is to be chosen in interactive way via slider. Such corrections are not new (see for example (Lauter, 159-168) in (Fedorov, 1992)).



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 1.5. Selecting variables

The standard Fisher approach was to maximize the between group variance or (what is the same) to minimize common within group variance:

$$SS = \sum_{g \in G} \sum_{i \in I(g)} (x_i - m(g))(x_i - m(g))'$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page

◀▶

◀▶

Go Back

Full Screen

Close

Quit

## 1.5. Selecting variables

The standard Fisher approach was to maximize the between group variance or (what is the same) to minimize common within group variance:

$$SS = \sum_{g \in G} \sum_{i \in I(g)} (x_i - m(g))(x_i - m(g))'$$

One may use trace or determinant to find corresponding variables. Now in all programs the so called Wilks lambda is used:

$$\Lambda = \frac{\det(SS_{in})}{\det(SS_{total})}$$



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page

◀▶

◀▶

Go Back

Full Screen

Close

Quit

## 1.5. Selecting variables

The standard Fisher approach was to maximize the between group variance or (what is the same) to minimize common within group variance:

$$SS = \sum_{g \in G} \sum_{i \in I(g)} (x_i - m(g))(x_i - m(g))'$$

One may use trace or determinant to find corresponding variables. Now in all programs the so called Wilks lambda is used:

$$\Lambda = \frac{\det(SS_{in})}{\det(SS_{total})}$$

It is easy to calculate (see (Jennrich, 1977)) and convenient to update when new variable is to be chosen.

## 2. File

Now we will go through the menus of the program.



Overview

**File**

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 2. File

Now we will go through the menus of the program.

The File drop down menu may be used separately in order to fill the missing data with within group means.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 2.1. Open Data

The program loads a data (.csv) file.

The program will ask you to supply one variable to be used for classification. You should decide. Otherwise the use of `Ldagui.m` is impossible.



*Overview*

*File*

*Model*

*Diagnostics*

*Use*

*Algorithms*

*References*

*Home Page*

*Title Page*



*Go Back*

*Full Screen*

*Close*

*Quit*



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 2.1. Open Data

The program loads a data (.csv) file.

The program will ask you to supply one variable to be used for classification. You should decide. Otherwise the use of `Ldagui.m` is impossible.

### Comma separated values

These (.csv) files are common for many applications. They are easily exported and imported by Excel and Statistica programs.





Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 2.1. Open Data

The program loads a data (.csv) file.

The program will ask you to supply one variable to be used for classification. You should decide. Otherwise the use of `Ldagui.m` is impossible.

### Comma separated values

These (.csv) files are common for many applications. They are easily exported and imported by Excel and Statistica programs.

`Ldagui.m` assumes that:

- the first row contains strings for variable names;
- the first column contains strings for case names.

All other fields should contain numbers (or be empty for missing values).

## Categorical variables

The categorical variables should have consecutive positive integer values. When exporting from Statistica you should say **integers** instead of **text** values.



*Overview*

*File*

*Model*

*Diagnostics*

*Use*

*Algorithms*

*References*

*Home Page*

*Title Page*



*Go Back*

*Full Screen*

*Close*

*Quit*



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## Categorical variables

The categorical variables should have consecutive positive integer values. When exporting from Statistica you should say **integers** instead of **text** values.

Any information about text values they may have in Statistica is lost and moreover their values in `Ldagui.m` are changed to first natural numbers: 1,2,3...

## 2.2. Filling Missing Data

You will be asked to supply selection variable. This is not obligatory. Then you should supply obligatory classification variable.



*Overview*

*File*

*Model*

*Diagnostics*

*Use*

*Algorithms*

*References*

*Home Page*

*Title Page*



*Go Back*

*Full Screen*

*Close*

*Quit*

## 2.2. Filling Missing Data

You will be asked to supply selection variable. This is not obligatory. Then you should supply obligatory classification variable.

Both classification and selection variables will be used in the algorithm for filling missing data.



[Overview](#)

[File](#)

[Model](#)

[Diagnostics](#)

[Use](#)

[Algorithms](#)

[References](#)

[Home Page](#)

[Title Page](#)



[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 2.2. Filling Missing Data

You will be asked to supply selection variable. This is not obligatory. Then you should supply obligatory classification variable.

Both classification and selection variables will be used in the algorithm for filling missing data.

**Filling Missing Data** is done automatically by LDAgui upon reading of .csv data. They are replaced by within group means. These means are formed by each combination of values of classification and selection variables.

## 2.3. Save Data

Saves the data file in a form of comma separated file for later import in Excel or Statistica.



*Overview*

*File*

*Model*

*Diagnostics*

*Use*

*Algorithms*

*References*

*Home Page*

*Title Page*



*Go Back*

*Full Screen*

*Close*

*Quit*

## 2.3. Save Data

Saves the data file in a form of comma separated file for later import in Excel or Statistica.

## 2.4. Exit to MATLAB

Saves the MATLAB workspace in tempmodel.mat for later examination and use.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



### 2.3. Save Data

Saves the data file in a form of comma separated file for later import in Excel or Statistica.

### 2.4. Exit to MATLAB

Saves the MATLAB workspace in tempmodel.mat for later examination and use.

### 2.5. Quit MATLAB



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3. Model

Under model we understand:

1. the training sample (with no missing values);



*Overview*

*File*

*Model*

*Diagnostics*

*Use*

*Algorithms*

*References*

*Home Page*

*Title Page*



*Go Back*

*Full Screen*

*Close*

*Quit*



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3. Model

Under model we understand:

1. the training sample (with no missing values);
2. a subset of cases having fixed value of the selection variable (if any);



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3. Model

Under model we understand:

1. the training sample (with no missing values);
2. a subset of cases having fixed value of the selection variable (if any);
3. a subset of variables chosen for predictors (may be empty);



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3. Model

Under model we understand:

1. the training sample (with no missing values);
2. a subset of cases having fixed value of the selection variable (if any);
3. a subset of variables chosen for predictors (may be empty);
4. fixed value of the parameter  $\alpha$  (12) of nonlinearity.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3. Model

Under model we understand:

1. the training sample (with no missing values);
2. a subset of cases having fixed value of the selection variable (if any);
3. a subset of variables chosen for predictors (may be empty);
4. fixed value of the parameter  $\alpha$  (12) of nonlinearity.
5. the estimated parameters of DA model.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 3. Model

Under model we understand:

1. the training sample (with no missing values);
2. a subset of cases having fixed value of the selection variable (if any);
3. a subset of variables chosen for predictors (may be empty);
4. fixed value of the parameter  $\alpha$  (12) of nonlinearity.
5. the estimated parameters of DA model.

### 3.1. Build Model

Performs all preliminary calculations for an empty model with no selection variable taken into account. To activate this option click on the Selection Listbox.

## 3.2. Load Model

Loads previously saved model (or workspace).



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit





Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 3.2. Load Model

Loads previously saved model (or workspace).

## 3.3. Save Model

Save the current model with data, names, selected groups, predictors, etc. for later use. In fact the current workspace of MATLAB is saved.

### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.
- Variables in model with their  $f$ - and  $p$ - values ordered;



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.
- Variables in model with their  $f$ - and  $p$ - values ordered;
- Value of parameter  $\alpha$  responsible for nonlinearity;



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.
- Variables in model with their  $f$ - and  $p$ - values ordered;
- Value of parameter  $\alpha$  responsible for nonlinearity;
- Value and p-value of Wilks  $\Lambda$ ;

### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.
- Variables in model with their  $f$ - and  $p$ - values ordered;
- Value of parameter  $\alpha$  responsible for nonlinearity;
- Value and p-value of Wilks  $\Lambda$ ;
- Results of the classification of the training sample - number of errors;



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.
- Variables in model with their  $f$ - and  $p$ - values ordered;
- Value of parameter  $\alpha$  responsible for nonlinearity;
- Value and p-value of Wilks  $\Lambda$ ;
- Results of the classification of the training sample - number of errors;
- Cases classified with probability below .8;

### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.
- Variables in model with their  $f$ - and  $p$ - values ordered;
- Value of parameter  $\alpha$  responsible for nonlinearity;
- Value and p-value of Wilks  $\Lambda$ ;
- Results of the classification of the training sample - number of errors;
- Cases classified with probability below .8;
- Estimated power of the model by groups.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



### 3.4. Print Results

The following results are printed in the MATLAB command window:

- File name - the name of file (data or model) you have loaded recently;
- Model name - the name of corresponding value of selection variable if any;
- Number of cases in training sample.
- Variables in model with their  $f$ - and  $p$ - values ordered;
- Value of parameter  $\alpha$  responsible for nonlinearity;
- Value and p-value of Wilks  $\Lambda$ ;
- Results of the classification of the training sample - number of errors;
- Cases classified with probability below .8;
- Estimated power of the model by groups.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



[Overview](#)

[File](#)

[Model](#)

[Diagnostics](#)

[Use](#)

[Algorithms](#)

[References](#)

[Home Page](#)

[Title Page](#)



[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Finally, a huge sample with 6000 observations per group is produced according estimated within group means and covariance matrices. The sample is classified and results reported on the MATLAB command window. This may be considered as an estimate of the theoretical power of the model.

### 3.5. Clear Model

Clears any information for the model. You should start with Build model step.



[Overview](#)

[File](#)

[Model](#)

[Diagnostics](#)

[Use](#)

[Algorithms](#)

[References](#)

[Home Page](#)

[Title Page](#)



[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Finally, a huge sample with 6000 observations per group is produced according estimated within group means and covariance matrices. The sample is classified and results reported on the MATLAB command window. This may be considered as an estimate of the theoretical power of the model.

### 3.5. Clear Model

Clears any information for the model. You should start with Build model step.



*Overview*

*File*

*Model*

*Diagnostics*

*Use*

*Algorithms*

*References*

*Home Page*

*Title Page*



*Go Back*

*Full Screen*

*Close*

*Quit*

## 4. Diagnostics

The tools proposed for making adequate decision are:

- Test - (Ctrl-t) - produces a test random sample;



*Overview*

*File*

*Model*

*Diagnostics*

*Use*

*Algorithms*

*References*

*Home Page*

*Title Page*



*Go Back*

*Full Screen*

*Close*

*Quit*

## 4. Diagnostics

The tools proposed for making adequate decision are:

- Test - (Ctrl-t) - produces a test random sample;
- Leave-One-Out - checks the model against deleting each of observations;

## 4. Diagnostics

The tools proposed for making adequate decision are:

- Test - (Ctrl-t) - produces a test random sample;
- Leave-One-Out - checks the model against deleting each of observations;
- Plot - (Ctrl-g) - makes two plots over canonical variables.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 4. Diagnostics

The tools proposed for making adequate decision are:

- Test - (Ctrl-t) - produces a test random sample;
- Leave-One-Out - checks the model against deleting each of observations;
- Plot - (Ctrl-g) - makes two plots over canonical variables.

### 4.1. Test

A small sample with 100 observations per group is produced according estimated within group means and covariance matrices. The sample is classified and results reported on the MATLAB command window. This may be considered as an estimate of the power of the model. One may repeat this step in order to be sure or use print menu with larger sample **3.4.**



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 4. Diagnostics

The tools proposed for making adequate decision are:

- Test - (Ctrl-t) - produces a test random sample;
- Leave-One-Out - checks the model against deleting each of observations;
- Plot - (Ctrl-g) - makes two plots over canonical variables.

### 4.1. Test

A small sample with 100 observations per group is produced according estimated within group means and covariance matrices. The sample is classified and results reported on the MATLAB command window. This may be considered as an estimate of the power of the model. One may repeat this step in order to be sure or use print menu with larger sample **3.4.**





Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 4.2. Leave-One-Out

A standard procedure is performed:

1. For each observation in the training sample a model with the same variables is build but without this particular observation.
2. The training sample is classified with this new model and classification errors counted.
3. The errors for all observations are summarized and reported.

## 4.3. Plot

Second (upper plot) and third canonical variables are plotted against the first (on horizontal axes). The training sample is plotted with different colors for the groups



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 4.2. Leave-One-Out

A standard procedure is performed:

1. For each observation in the training sample a model with the same variables is build but without this particular observation.
2. The training sample is classified with this new model and classification errors counted.
3. The errors for all observations are summarized and reported.

## 4.3. Plot

Second (upper plot) and third canonical variables are plotted against the first (on horizontal axes). The training sample is plotted with different colors for the groups

## 5. Use

### 5.1. Load sample

A standard data (.csv) file is loaded which should not contain missing values in the columns used for recognition. Columns to use should have the same variable names.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 5. Use

### 5.1. Load sample

A standard data (.csv) file is loaded which should not contain missing values in the columns used for recognition. Columns to use should have the same variable names.

### 5.2. Print results

Results of classification are printed.



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit

## 5. Use

### 5.1. Load sample

A standard data (.csv) file is loaded which should not contain missing values in the columns used for recognition. Columns to use should have the same variable names.

### 5.2. Print results

Results of classification are printed.

### 5.3. Save sample

The sample is saved in a data (.csv) file with resulting classification in the first column.

## 6. Algorithms

The calculations are based on the paper of (Jennrich, 1977) in the classical collection of (Einslein, Ralston et al., 1977) being in the foundations of the package BMDP(see (Dixon, 1981)).



[Overview](#)

[File](#)

[Model](#)

[Diagnostics](#)

[Use](#)

**[Algorithms](#)**

[References](#)

[Home Page](#)

[Title Page](#)



[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

## 6. Algorithms

The calculations are based on the paper of (Jennrich, 1977) in the classical collection of (Einslein, Ralston et al., 1977) being in the foundations of the package BMDP(see (Dixon, 1981)).



[Overview](#)

[File](#)

[Model](#)

[Diagnostics](#)

[Use](#)

**[Algorithms](#)**

[References](#)

[Home Page](#)

[Title Page](#)



[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

## References

DIXON, J. (ed.) (1981). *BMDP Statistical Software - 81*. University of California, Los Angeles.

EINSLEIN, K., RALSTON, A., ET AL. (eds.) (1977). *Statistical Methods for Digital Computers*. John Wiley & Sons, New York.

FEDOROV, V. (ed.) (1992). *Model oriented data analysis: a survey of recent methods*. Physica-Verlag, Heidelberg.

JENNRICH, R. I. (1977). Stepwise discriminant analysis. In: (Einslein, Ralston et al., 1977), pp. 76–95.

LAUTER, H. (159-168). Bootstrap and estimation of nonlinear parameters. In: (Fedorov, 1992).



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit



## References

DIXON, J. (ed.) (1981). *BMDP Statistical Software - 81*. University of California, Los Angeles.

EINSLEIN, K., RALSTON, A., ET AL. (eds.) (1977). *Statistical Methods for Digital Computers*. John Wiley & Sons, New York.

FEDOROV, V. (ed.) (1992). *Model oriented data analysis: a survey of recent methods*. Physica-Verlag, Heidelberg.

JENNRICH, R. I. (1977). Stepwise discriminant analysis. In: (Einslein, Ralston et al., 1977), pp. 76–95.

LAUTER, H. (159-168). Bootstrap and estimation of nonlinear parameters. In: (Fedorov, 1992).



Overview

File

Model

Diagnostics

Use

Algorithms

References

Home Page

Title Page



Go Back

Full Screen

Close

Quit