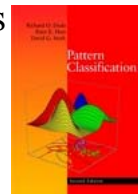


# Chapter 5: Linear Discriminant Functions

- ▼ Introduction
- ▼ Linear Discriminant Functions and Decisions Surfaces
- ▼ Generalized Linear Discriminant Functions



All materials used in this course were taken from the textbook "*Pattern Classification*" by Duda et al., John Wiley & Sons, 2001 with the permission of the authors and the publisher

## ▼ Introduction

- In chapter 3, the underlying probability densities were known (or given)
- The training sample was used to estimate the parameters of these probability densities (ML, MAP estimations)
- In this chapter, we only know the proper forms for the discriminant functions: similar to non-parametric techniques
- They may not be optimal, but they are very simple to use
- They provide us with linear classifiers

### Linear discriminant functions and decisions surfaces

#### – Definition

It is a function that is a linear combination of the components of  $x$

$$g(x) = w^t x + w_0 \quad (1)$$

where  $w$  is the weight vector and  $w_0$  the bias

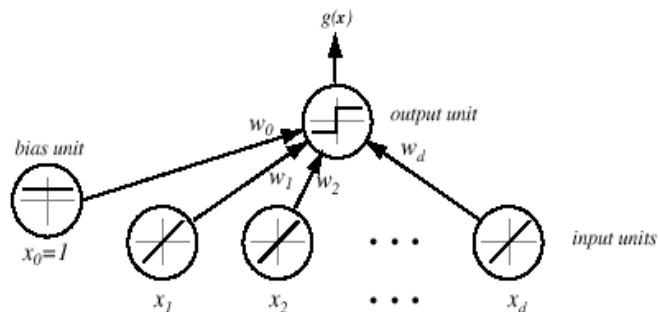
#### – A two-category classifier with a discriminant function of the form (1) uses the following rule:

Decide  $\omega_1$  if  $g(x) > 0$  and  $\omega_2$  if  $g(x) < 0$

$\Leftrightarrow$  Decide  $\omega_1$  if  $w^t x > -w_0$  and  $\omega_2$  otherwise

If  $g(x) = 0 \Rightarrow x$  is assigned to either class

2

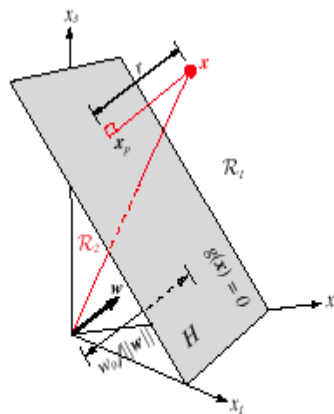


**FIGURE 5.1.** A simple linear classifier having  $d$  input units, each corresponding to the values of the components of an input vector. Each input feature value  $x_i$  is multiplied by its corresponding weight  $w_i$ ; the effective input at the output unit is the sum all these products,  $\sum w_i x_i$ . We show in each unit its effective input-output function. Thus each of the  $d$  input units is linear, emitting exactly the value of its corresponding feature value. The single bias unit unit always emits the constant value 1.0. The single output unit emits a +1 if  $w^t x + w_0 > 0$  or a -1 otherwise. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

2

- The equation  $g(\mathbf{x}) = 0$  defines the decision surface that separates points assigned to the category  $\omega_1$  from points assigned to the category  $\omega_2$
- When  $g(\mathbf{x})$  is linear, the decision surface is a hyperplane
- Algebraic measure of the distance from  $\mathbf{x}$  to the hyperplane (interesting result!)

2



**FIGURE 5.2.** The linear decision boundary  $H$ , where  $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = 0$ , separates the feature space into two half-spaces  $\mathcal{R}_1$  (where  $g(\mathbf{x}) > 0$ ) and  $\mathcal{R}_2$  (where  $g(\mathbf{x}) < 0$ ). From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

3

$$\mathbf{x} = \mathbf{x}_p + \frac{\mathbf{r} \cdot \mathbf{w}}{\|\mathbf{w}\|} \quad (\text{since } \mathbf{w} \text{ is colinear with } \mathbf{x} - \mathbf{x}_p \text{ and } \frac{\mathbf{w}}{\|\mathbf{w}\|} = \mathbf{1})$$

$$\text{since } \mathbf{g}(\mathbf{x}) = 0 \text{ and } \mathbf{w}^t \cdot \mathbf{w} = \|\mathbf{w}\|^2$$

$$\text{therefore } \mathbf{r} = \frac{\mathbf{g}(\mathbf{x})}{\|\mathbf{w}\|}$$

$$\text{in particular } d(\mathbf{0}, \mathbf{H}) = \frac{\mathbf{w}_0}{\|\mathbf{w}\|}$$

- In conclusion, a linear discriminant function divides the feature space by a hyperplane decision surface
- The orientation of the surface is determined by the normal vector  $\mathbf{w}$  and the location of the surface is determined by the bias

2

### – The multi-category case

- We define  $c$  linear discriminant functions

$$\mathbf{g}_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + \mathbf{w}_{i0} \quad i = 1, \dots, c$$

and assign  $\mathbf{x}$  to  $\omega_i$  if  $\mathbf{g}_i(\mathbf{x}) > \mathbf{g}_j(\mathbf{x}) \forall j \neq i$ ; in case of ties, the classification is undefined

- In this case, the classifier is a “linear machine”
- A linear machine divides the feature space into  $c$  decision regions, with  $\mathbf{g}_i(\mathbf{x})$  being the largest discriminant if  $\mathbf{x}$  is in the region  $\mathbf{R}_i$
- For a two contiguous regions  $\mathbf{R}_i$  and  $\mathbf{R}_j$ ; the boundary that separates them is a portion of hyperplane  $\mathbf{H}_{ij}$  defined by:

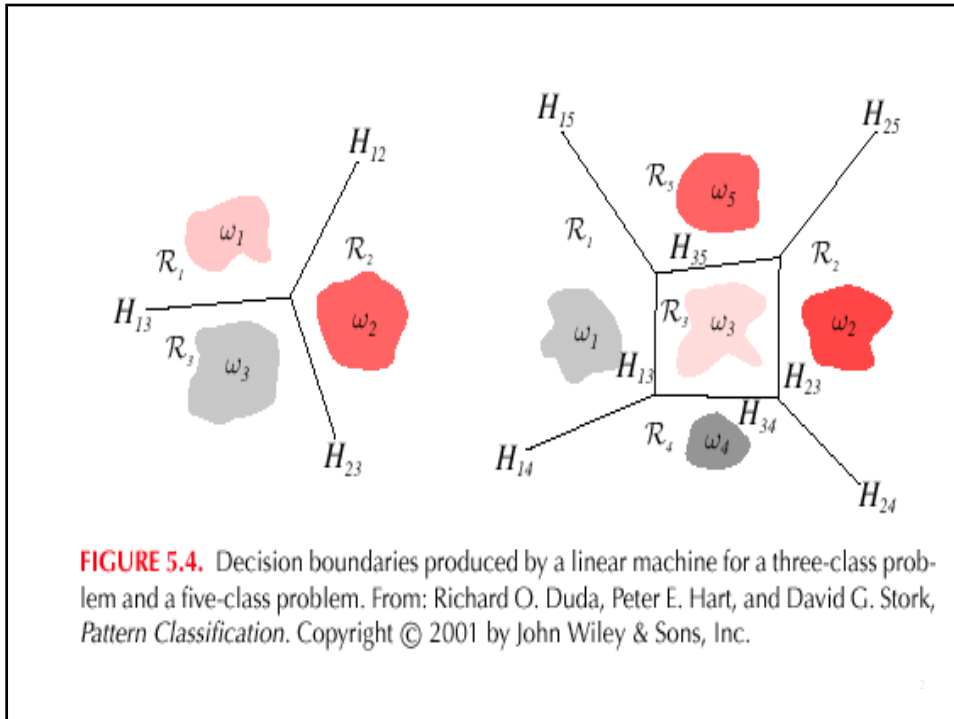
$$\mathbf{g}_i(\mathbf{x}) = \mathbf{g}_j(\mathbf{x})$$

$$\Leftrightarrow (\mathbf{w}_i - \mathbf{w}_j)^t \mathbf{x} + (\mathbf{w}_{i0} - \mathbf{w}_{j0}) = 0$$

- $\mathbf{w}_i - \mathbf{w}_j$  is normal to  $\mathbf{H}_{ij}$  and

$$d(\mathbf{x}, \mathbf{H}_{ij}) = \frac{\mathbf{g}_i - \mathbf{g}_j}{\|\mathbf{w}_i - \mathbf{w}_j\|}$$

3



- It is easy to show that the decision regions for a linear machine are convex, this restriction limits the flexibility and accuracy of the classifier

## Class Exercises

- ✓ Ex. 13 p.159
- ✓ Ex. 3 p.201
- ✓ Write a C/C++/Java program that uses a k-nearest neighbor method to classify input patterns. Use the table on p.209 as your training sample.

Experiment the program with the following data:

- $k = 3$      $x_1 = (0.33, 0.58, -4.8)$   
                    $x_2 = (0.27, 1.0, -2.68)$   
                    $x_3 = (-0.44, 2.8, 6.20)$
- Do the same thing with  $k = 11$
- Compare the classification results between  $k = 3$  and  $k = 11$   
 (use the most dominant class voting scheme amongst the  $k$  classes)

2

## ✓ Generalized Linear Discriminant Functions

- Decision boundaries which separate between classes may not always be linear
- The complexity of the boundaries may sometimes request the use of highly non-linear surfaces
- A popular approach to generalize the concept of linear decision functions is to consider a generalized decision function as:

$$g(x) = w_1 f_1(x) + w_2 f_2(x) + \dots + w_N f_N(x) + w_{N+1} \quad (1)$$

where  $f_i(x)$ ,  $1 \leq i \leq N$  are scalar functions of the pattern  $x$ ,  
 $x \in \mathbb{R}^n$

3

- Introducing  $f_{n+1}(x) = 1$  we get:

$$\mathbf{g}(\mathbf{x}) = \sum_{i=1}^{N+1} \mathbf{w}_i f_i(\mathbf{x}) = \mathbf{w}^T \cdot \dot{\mathbf{x}}$$

where  $\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N, \mathbf{w}_{N+1})^T$  and  $\dot{\mathbf{x}} = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_N(\mathbf{x}), f_{N+1}(\mathbf{x}))^T$

- This latter representation of  $g(x)$  implies that any decision function defined by equation (1) can be treated as linear in the  $(N + 1)$  dimensional space ( $N + 1 > n$ )
- $g(x)$  maintains its non-linearity characteristics in  $\mathbb{R}^n$

3

- The most commonly used generalized decision function is  $g(x)$  for which  $f_i(x)$  ( $1 \leq i \leq N$ ) are polynomials

$$\mathbf{g}(\mathbf{x}) = (\dot{\mathbf{w}})^T \mathbf{x} \quad \text{T: is the vector transpose form}$$

Where  $\dot{\mathbf{w}}$  is a new weight vector, which can be calculated from the original  $\mathbf{w}$  and the original linear  $f_i(x)$ ,  $1 \leq i \leq N$

- Quadratic decision functions for a 2-dimensional feature space

$$\mathbf{g}(\mathbf{x}) = \mathbf{w}_1 \mathbf{x}_1^2 + \mathbf{w}_2 \mathbf{x}_1 \mathbf{x}_2 + \mathbf{w}_3 \mathbf{x}_2^2 + \mathbf{w}_4 \mathbf{x}_1 + \mathbf{w}_5 \mathbf{x}_2 + \mathbf{w}_6$$

here :  $\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_6)^T$  and  $\dot{\mathbf{x}} = (\mathbf{x}_1^2, \mathbf{x}_1 \mathbf{x}_2, \mathbf{x}_2^2, \mathbf{x}_1, \mathbf{x}_2, 1)^T$

3

- For patterns  $x \in \mathbb{R}^n$ , the most general quadratic decision function is given by:

$$g(x) = \sum_{i=1}^n w_{ii} x_i^2 + \sum_{i=1}^{n-1} \sum_{j=i+1}^n w_{ij} x_i x_j + \sum_{i=1}^n w_i x_i + w_{n+1} \quad (2)$$

The number of terms at the right-hand side is:

$$l = N + 1 = n + \frac{n(n-1)}{2} + n + 1 = \frac{(n+1)(n+2)}{2}$$

This is the total number of weights which are the free parameters of the problem

- If for example  $n = 3$ , the vector  $\mathbf{x}$  is 10-dimensional
- If for example  $n = 10$ , the vector  $\mathbf{x}$  is 65-dimensional

3

- In the case of polynomial decision functions of order  $m$ , a typical  $f_i(x)$  is given by:

$$f_i(x) = x_{i_1}^{e_1} x_{i_2}^{e_2} \dots x_{i_m}^{e_m}$$

where  $1 \leq i_1, i_2, \dots, i_m \leq n$  and  $e_i, 1 \leq i \leq m$  is 0 or 1.

- It is a polynomial with a degree between 0 and  $m$ . To avoid repetitions, we request  $i_1 \leq i_2 \leq \dots \leq i_m$

$$g^m(x) = \sum_{i_1=1}^n \sum_{i_2=i_1}^n \dots \sum_{i_m=i_{m-1}}^n w_{i_1 i_2 \dots i_m} x_{i_1} x_{i_2} \dots x_{i_m} + g^{m-1}(x)$$

(where  $g^0(x) = w_{n+1}$ ) is the most general polynomial decision function of order  $m$

3



Example 1: Let  $n = 3$  and  $m = 2$  then:

$$\begin{aligned} g^2(\mathbf{x}) &= \sum_{i_1=1}^3 \sum_{i_2=i_1}^3 w_{i_1 i_2} x_{i_1} x_{i_2} + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 \\ &= w_{11} x_1^2 + w_{12} x_1 x_2 + w_{13} x_1 x_3 + w_{22} x_2^2 + w_{23} x_2 x_3 + w_{33} x_3^2 \\ &\quad + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 \end{aligned}$$

Example 2: Let  $n = 2$  and  $m = 3$  then:

$$\begin{aligned} g^3(\mathbf{x}) &= \sum_{i_1=1}^2 \sum_{i_2=i_1}^2 \sum_{i_3=i_2}^2 w_{i_1 i_2 i_3} x_{i_1} x_{i_2} x_{i_3} + g^2(\mathbf{x}) \\ &= w_{111} x_1^3 + w_{112} x_1^2 x_2 + w_{122} x_1 x_2^2 + w_{222} x_2^3 + g^2(\mathbf{x}) \\ \text{where } g^2(\mathbf{x}) &= \sum_{i_1=1}^2 \sum_{i_2=i_1}^2 w_{i_1 i_2} x_{i_1} x_{i_2} + g^1(\mathbf{x}) \\ &= w_{11} x_1^2 + w_{12} x_1 x_2 + w_{22} x_2^2 + w_1 x_1 + w_2 x_2 + w_3 \end{aligned}$$

- The commonly used quadratic decision function can be represented as the general  $n$ - dimensional quadratic surface:

$$g(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{b} + c$$

where the matrix  $\mathbf{A} = (a_{ij})$ , the vector

$\mathbf{b} = (b_1, b_2, \dots, b_n)^T$  and  $c$ , depends on the weights  $w_{ii}$ ,  $w_{ij}$ ,  $w_i$  of equation (2)

- If  $\mathbf{A}$  is positive definite then the decision function is a hyperellipsoid with axes in the directions of the eigenvectors of  $\mathbf{A}$ 
  - In particular: if  $\mathbf{A} = \mathbf{I}_n$  (Identity), the decision function is simply the  $n$ -dimensional hypersphere

- If  $A$  is negative definite, the decision function describes a hyperhyperboloid
- In conclusion: it is only the matrix  $A$  which determines the shape and characteristics of the decision function

3

**Problem:** Consider a 3 dimensional space and cubic polynomial decision functions

1. How many terms are needed to represent a decision function if only cubic and linear functions are assumed
2. Present the general 4<sup>th</sup> order polynomial decision function for a 2 dimensional pattern space
3. Let  $\mathbb{R}^3$  be the original pattern space and let the decision function associated with the pattern classes  $\omega_1$  and  $\omega_2$  be:

$$\mathbf{g}(\mathbf{x}) = 2\mathbf{x}_1^2 + \mathbf{x}_3^2 + \mathbf{x}_2\mathbf{x}_3 + 4\mathbf{x}_1 - 2\mathbf{x}_2 + 1$$

for which  $g(\mathbf{x}) > 0$  if  $\mathbf{x} \in \omega_1$  and  $g(\mathbf{x}) < 0$  if  $\mathbf{x} \in \omega_2$

- a) Rewrite  $g(\mathbf{x})$  as  $g(\mathbf{x}) = \mathbf{x}^T\mathbf{A}\mathbf{x} + \mathbf{x}^T\mathbf{b} + c$
- b) Determine the class of each of the following pattern vectors:  
(1,1,1), (1,10,0), (0,1/2,0)

3

### ▼ Positive Definite Matrices

1. A square matrix  $A$  is *positive definite* if  $x^T A x > 0$  for all nonzero column vectors  $x$ .
2. It is *negative definite* if  $x^T A x < 0$  for all nonzero  $x$ .
3. It is *positive semi-definite* if  $x^T A x \geq 0$ .
4. And *negative semi-definite* if  $x^T A x \leq 0$  for all  $x$ .

These definitions are hard to check directly and you might as well forget them for all practical purposes.

3

More useful in practice are the following properties, which hold when the matrix  $A$  is symmetric and which are easier to check.

The *ith principal minor* of  $A$  is the matrix  $A_i$  formed by the first  $i$  rows and columns of  $A$ . So, the first principal minor of  $A$  is the matrix  $A_1 = (a_{11})$ , the second principal minor is the matrix:

$$A_2 = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}, \text{ and so on.}$$

3

- The matrix  $A$  is positive definite if all its principal minors  $A_1, A_2, \dots, A_n$  have strictly positive determinants
- If these determinants are non-zero and alternate in signs, starting with  $\det(A_1) < 0$ , then the matrix  $A$  is negative definite
- If the determinants are all non-negative, then the matrix is positive semi-definite
- If the determinant alternate in signs, starting with  $\det(A_1) \leq 0$ , then the matrix is negative semi-definite

3

To fix ideas, consider a 2x2 symmetric matrix:

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}.$$

- It is positive definite if:
  - a)  $\det(A_1) = a_{11} > 0$
  - b)  $\det(A_2) = a_{11}a_{22} - a_{12}a_{12} > 0$
- It is negative definite if:
  - a)  $\det(A_1) = a_{11} < 0$
  - b)  $\det(A_2) = a_{11}a_{22} - a_{12}a_{12} > 0$
- It is positive semi-definite if:
  - a)  $\det(A_1) = a_{11} \geq 0$
  - b)  $\det(A_2) = a_{11}a_{22} - a_{12}a_{12} \geq 0$
- And it is negative semi-definite if:
  - a)  $\det(A_1) = a_{11} \leq 0$
  - b)  $\det(A_2) = a_{11}a_{22} - a_{12}a_{12} \geq 0$ .

3

**Exercise 1:** Check whether the following matrices are positive definite, negative definite, positive semi-definite, negative semi-definite or none of the above.

$$(a) \mathbf{A} = \begin{pmatrix} 2 & 1 \\ 1 & 4 \end{pmatrix}$$

$$(b) \mathbf{A} = \begin{pmatrix} -2 & 4 \\ 4 & -8 \end{pmatrix}$$

$$(c) \mathbf{A} = \begin{pmatrix} -2 & 2 \\ 2 & -4 \end{pmatrix}$$

$$(d) \mathbf{A} = \begin{pmatrix} 2 & 4 \\ 4 & 3 \end{pmatrix}$$

3

**Solutions of Exercise 1:**

$$\begin{aligned} \checkmark \quad A_1 &= 2 > 0 \\ A_2 &= 8 - 1 = 7 > 0 \quad \Rightarrow \mathbf{A} \text{ is positive definite} \end{aligned}$$

$$\begin{aligned} \checkmark \quad A_1 &= -2 \\ A_2 &= (-2 \times -8) - 16 = 0 \quad \Rightarrow \mathbf{A} \text{ is negative semi-positive} \end{aligned}$$

$$\begin{aligned} \checkmark \quad A_1 &= -2 \\ A_2 &= 8 - 4 = 4 > 0 \quad \Rightarrow \mathbf{A} \text{ is negative definite} \end{aligned}$$

$$\begin{aligned} \checkmark \quad A_1 &= 2 > 0 \\ A_2 &= 6 - 16 = -10 < 0 \quad \Rightarrow \mathbf{A} \text{ is none of the above} \end{aligned}$$

3

## Exercise 2:

$$\text{Let } \mathbf{A} = \begin{pmatrix} 2 & 1 \\ 1 & 4 \end{pmatrix}$$

1. Compute the decision boundary assigned to the matrix  $\mathbf{A}$  ( $g(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{b} + c$ ) in the case where  $\mathbf{b}^T = (1, 2)$  and  $c = -3$
2. Solve  $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$  and find the shape and the characteristics of the decision boundary separating two classes  $\omega_1$  and  $\omega_2$
3. Classify the following points:
  - $\mathbf{x}^T = (0, -1)$
  - $\mathbf{x}^T = (1, 1)$

3

## Solution of Exercise 2:

$$\begin{aligned} 1. \quad g(\mathbf{x}) &= (\mathbf{x}_1, \mathbf{x}_2) \begin{pmatrix} 2 & 1 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} + (\mathbf{x}_1, \mathbf{x}_2) \begin{pmatrix} 1 \\ 2 \end{pmatrix} - 3 \\ &= (2\mathbf{x}_1 + \mathbf{x}_2, \mathbf{x}_1 + 4\mathbf{x}_2) \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} + \mathbf{x}_1 + 2\mathbf{x}_2 - 3 \\ &= 2\mathbf{x}_1^2 + \mathbf{x}_1\mathbf{x}_2 + \mathbf{x}_1\mathbf{x}_2 + 4\mathbf{x}_2^2 + \mathbf{x}_1 + 2\mathbf{x}_2 - 3 \\ &= 2\mathbf{x}_1^2 + 4\mathbf{x}_2^2 + 2\mathbf{x}_1\mathbf{x}_2 + \mathbf{x}_1 + 2\mathbf{x}_2 - 3 \end{aligned}$$

$$2. \quad \text{For } \lambda_1 = 3 + \sqrt{2} \text{ using } \begin{pmatrix} 2 - \lambda & 1 \\ 1 & 4 - \lambda \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \mathbf{0}, \text{ we obtain :}$$

$$\begin{cases} (-1 - \sqrt{2})\mathbf{x}_1 + \mathbf{x}_2 = 0 \\ \mathbf{x}_1 + (1 - \sqrt{2})\mathbf{x}_2 = 0 \end{cases} \Leftrightarrow (-1 - \sqrt{2})\mathbf{x}_1 + \mathbf{x}_2 = 0$$

This latter equation is a straight line colinear to the vector:

$$\vec{\mathbf{V}}_1 = (1, 1 + \sqrt{2})^T$$

3

For  $\lambda_2 = 3 + \sqrt{2}$  using  $\begin{pmatrix} 2-\lambda & 1 \\ 1 & 4-\lambda \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \mathbf{0}$ , we obtain :

$$\begin{cases} (\sqrt{2}-1)x_1 + x_2 = 0 \\ x_1 + (1+\sqrt{2})x_2 = 0 \end{cases} \Leftrightarrow (\sqrt{2}-1)x_1 + x_2 = 0$$

This latter equation is a straight line colinear to the vector:

$$\vec{V}_2 = (1, 1 - \sqrt{2})^T$$

The ellipsis decision boundary has two axes, which are respectively colinear to the vectors  $V_1$  and  $V_2$

$$3. \mathbf{X} = (0, -1)^T \Rightarrow g(0, -1) = -1 < 0 \Rightarrow \mathbf{x} \in \omega_2$$

$$\mathbf{X} = (1, 1)^T \Rightarrow g(1, 1) = 8 > 0 \Rightarrow \mathbf{x} \in \omega_1$$