



**HAL**  
open science

# Two Classes Linear Discrimination, a Min-Max Approach

Philippe Castagliola, Bernard Dubuisson

► **To cite this version:**

Philippe Castagliola, Bernard Dubuisson. Two Classes Linear Discrimination, a Min-Max Approach. Pattern Recognition Letters, 1989, 10 (5), pp.281-287. 10.1016/0167-8655(89)90030-5 . hal-00340822

**HAL Id: hal-00340822**

**<https://hal.science/hal-00340822>**

Submitted on 19 May 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# *Two classes linear discrimination A min-max approach*

P. CASTAGLIOLA and B. DUBUISSON

*Université de Technologie de Compiègne, URA CNRS. 817, B.P. 649, 60206 Compiègne Cedex, France*

Received 11 April 1989

Revised 22 June 1989

*Abstract:* In this paper we propose a new dichotomic linear discrimination algorithm based on a criterion recursively defined from the min and max functions. This criterion allows us, using a gradient algorithm, to find two hyperplanes which split two classes in an optimal way. We give two examples of the possibilities of the method.

*Key words:* Linear discrimination, min and max function, gradient algorithm.

## 1. Introduction

The main goal of pattern recognition is to assign a new observed vector into one class out of  $M$  known classes. This decision can be made with the help of discriminant surfaces obtained from the knowledge of sample vectors in classes.

Let us suppose that we have to discriminate only between two classes (dichotomy) in  $\mathbb{R}^d$ . Different possibilities exist; the optimum one is the Bayes discriminator. Associate costs are  $\{0,1\}$  (cost of good classification is 0, cost of an error is 1), the optimum decision is associate  $\mathbf{x}$  to  $\omega_i$  if:

$$p(\omega_i | \mathbf{x}) = \max p(\omega_j | \mathbf{x}), \quad j = 1,2.$$

$p(\omega_i | \mathbf{x})$  is the a posteriori probability of class  $\omega_i$  given  $\mathbf{x}$ .

If these probabilities are unknown, they have to be estimated using samples of different classes. This set of samples is called the learning set.

Another possibility is to first decide on the form of the boundary between the two classes. One of the most studied discriminant functions is the linear one. Classes are said to be linearly separable if we can compute the equation of a hyperplane in  $\mathbb{R}^d$  separating the two classes.

The first linear discrimination algorithm was developed by Rosenblatt (Rosenblatt, 1957) in the mid 1950s, and was inspired from the neuronal model of Mc Culloch and Pitts (Mc Culloch & Pitts, 1943). Its usual name is the Perceptron. This algorithm is considered as a special case of the reward-and-punishment procedure. Unfortunately, if the classes are not linearly separable, this algorithm does not converge and no indication of this fact is given to the user.

The Ho–Kashyap procedure (Ho & Kashyap, 1965) gives a better solution to this problem. It always converges, indicating that the classes are linearly separable or not; if they are, this algorithm furnishes a hyperplane equation.

New solutions are currently proposed to solve non-linear separability by using the connexionist approach

(Widrow, Hoff, 1960) (neuron models). Although this non-probabilistic approach has interesting new developments, we are not interested in this approach.

Our goal is to find two parallel hyperplanes  $\mathbb{P}_1$  and  $\mathbb{P}_2$  separating the two classes in an optimal way. This problem is completely equivalent to find an optimal unit vector  $\mathbf{u}^*$ , orthogonal to  $\mathbb{P}_1$  and  $\mathbb{P}_2$ , and maximising a criterion  $d(\mathbf{u})$  which represents a signed distance between the classes.

We will define this signed distance from the difference between some min and max functions, and use it to find the optimal solution.

If the signed distance  $d(\mathbf{u}^*)$ , corresponding to the optimal solution is positive, we will conclude that the classes are linearly separable; if it is not, we will conclude that they are not linearly separable.

## 2. Principle of the method

Let  $\omega_1$  and  $\omega_2$  be two classes of  $\mathbb{R}^d$ . We want to discriminate between these two classes.

The learning set  $\mathcal{L}$  is composed of  $n$  vectors,  $n_1$  from  $\omega_1$  and  $n_2$  from  $\omega_2$ :

$$\omega_1 = \{\mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,n_1}\} = \{\mathbf{x}_{1,i}\}_{i=1}^{n_1} \quad \text{and} \quad \omega_2 = \{\mathbf{x}_{2,1}, \dots, \mathbf{x}_{2,n_2}\} = \{\mathbf{x}_{2,i}\}_{i=1}^{n_2}.$$

Let  $\mathbf{u}$  be a unitary vector orthogonal to  $\mathbb{P}_1$  and  $\mathbb{P}_2$  (see Figure 1). And let  $\hat{x}_{i,k} = \mathbf{u}^t \mathbf{x}_{i,k}$ .

Now we have to choose a criterion for our problem. This criterion represents a signed distance:

$$d(\mathbf{u}) = s(\mathbf{u}) - b(\mathbf{u}) \tag{1}$$

where

$$s(\mathbf{u}) = \min(\hat{x}_{2,1}, \dots, \hat{x}_{2,n_2}) = \min(\hat{x}_{2,i})_1^{n_2}, \quad b(\mathbf{u}) = \max(\hat{x}_{1,1}, \dots, \hat{x}_{1,n_1}) = \max(\hat{x}_{1,i})_1^{n_1}.$$

Let  $\mathbf{u}^*$  be the direction which maximises the criterion  $d(\mathbf{u})$ . Two opposite cases may happen:

*Case 1:*  $d(\mathbf{u}^*) > 0$ . Then  $\omega_1$  and  $\omega_2$  are linearly separable (Figure 2a). We can easily find two parallel hyperplanes, orthogonal to the direction  $\mathbf{u}^*$ , with the equations:

$$\mathbb{P}_1: \mathbf{x}^t \mathbf{u}^* - b(\mathbf{u}^*) \quad \text{and} \quad \mathbb{P}_2: \mathbf{x}^t \mathbf{u}^* - s(\mathbf{u}^*).$$

If the classes are interchanged, the optimal vector  $\mathbf{u}^*$  becomes  $-\mathbf{u}^*$  and the difference  $d(\mathbf{u}^*)$  is still positive. For example, if we have first  $b(\mathbf{u}^*) = 2$  and  $s(\mathbf{u}^*) = 3$  then  $d(\mathbf{u}^*) = 1 > 0$ . After interchanging the classes,  $\mathbf{u}^*$  becomes  $-\mathbf{u}^*$ , and  $b(\mathbf{u}^*) = -3$ ,  $s(\mathbf{u}^*) = -2$ , then  $d(\mathbf{u}^*) = -2 - (-3) = 1$  is still positive.

*Case 2:*  $d(\mathbf{u}^*) < 0$ . Then  $\omega_1$  and  $\omega_2$  are not linearly separable with a zero error probability (Figure 2b). The same considerations can be reformulated about the interchange of  $\omega_1$  and  $\omega_2$ .

Let us first consider the functions min and max used above in the definition of the functions  $s(\mathbf{u})$  and  $b(\mathbf{u})$ .

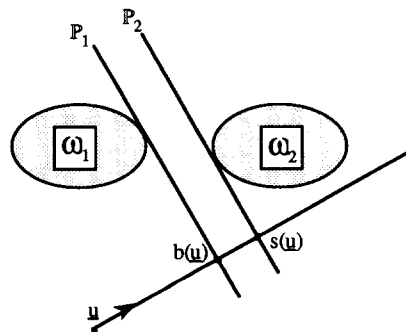


Figure 1.

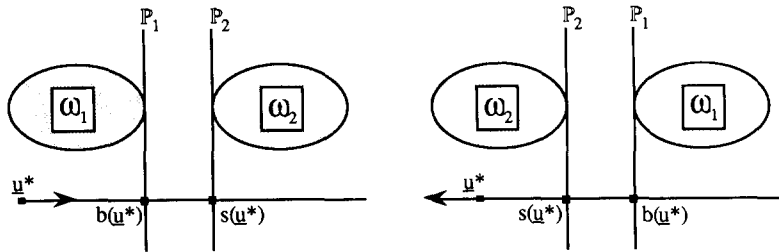


Figure 2a.

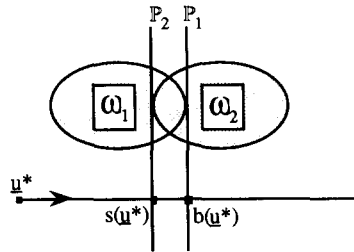


Figure 2b.

### 3. Definition and computation of MIN and MAX

#### 3.1. The MIN function

Let  $E = \{a_1, \dots, a_n\}$ , be a set of  $n$  real values. The min function is defined by the following:

$$\min\{a_1, \dots, a_n\} = a_i \quad \text{if } a_i < a_j \quad \text{for } j = 1, \dots, n, \quad j \neq i.$$

If the cardinality of  $E$  is one, it is obvious that  $\min(a_1) = a_1$ . Let us choose the cardinality of  $E$  equal to 2. It is also obvious that:

$$\min(a_1, a_2) = \frac{1}{2}(a_1 + a_2 - |a_1 - a_2|).$$

So recursively we can write when  $\text{card } E = n$ :

$$\min(a_1, \dots, a_n) = \frac{1}{2}(\min(a_1, \dots, a_{n-1}) + a_n - |\min(a_1, \dots, a_{n-1}) - a_n|).$$

Using the notation  $\min(a_1, \dots, a_n) = \min(a_i)_1^n$ :

$$\min(a_i)_1^1 = a_1, \quad \min(a_i)_1^n = \frac{1}{2}(\min(a_i)_1^{n-1} + a_n - |(\min(a_i)_1^{n-1} - a_n)|). \tag{2}$$

#### 3.2. The MAX function

The same considerations may be used for the max function.

$$\max(a_1) = a_1, \quad \max(a_1, a_2) = \frac{1}{2}(a_1 + a_2 + |a_1 - a_2|),$$

$$\max(a_i)_1^1 = a_1, \quad \max(a_i)_1^n = \frac{1}{2}(\max(a_i)_1^{n-1} + a_n + |(\max(a_i)_1^{n-1} - a_n)|). \tag{3}$$

#### 4. Computation of the optimal vector $u^*$

##### 4.1. A gradient algorithm

We search for a unit vector  $u^*$  that maximizes  $d(u)$ . This optimisation can be solved using a fixed step gradient:

$$u_{n+1} = u_n + p \cdot \frac{\partial}{\partial u} d(u), \quad p \text{ being the step to be fixed.}$$

If we use this method directly, the vector  $u_n$  tends towards the right direction, but unfortunately its norm tends towards  $+\infty$  (unconstrained solution).

We used the unitary norm constraint:  $\|u^*\| = 1$ . So, we must normalize the vector  $u$  after each gradient algorithm iteration. Thus, the modified gradient algorithm is:

$$v_n = u_n + p \cdot \frac{\partial}{\partial u} d(u), \quad u_{n+1} = \frac{v_n}{\|v_n\|}.$$

##### 4.2. Gradient of $s(u)$ , $b(u)$ and $d(u)$

The functions  $s(u)$  and  $b(u)$  are recursively defined, so are their gradient functions.

**Remark.** Let  $f(x)$  be a function from  $\mathbb{R}^n$  to  $\mathbb{R}$ , derivable from  $\mathbb{R}^n$ . Theoretically,  $(\partial/\partial x)|f(x)|$  is not defined on  $f(x) = 0$ . So, we will fix by definition:

$$\frac{\partial}{\partial x}|f(x)| = 0 \quad \text{when } f(x) = 0,$$

which is completely equivalent to:

$$\frac{\partial}{\partial x}|f(x)| = \text{sign}(f(x)) \frac{\partial}{\partial x} f(x)$$

with

$$\text{sign}(f(x)) = -1 \quad \text{if } f(x) < 0, \quad \text{sign}(f(x)) = 0 \quad \text{if } f(x) = 0, \quad \text{sign}(f(x)) = 1 \quad \text{if } f(x) > 0.$$

We have

$$\frac{\partial}{\partial u} s(u) = \frac{\partial}{\partial u} \min(\hat{x}_{2,i}^{n_2})$$

and from (2)

$$= \frac{1}{2} \left( \frac{\partial}{\partial u} \min(\hat{x}_{2,i}^{n_2-1} + x_{2,n_2} - \text{sign}(\min(\hat{x}_{2,i}^{n_2-1} - x_{2,n_2})) \left( \frac{\partial}{\partial u} \min(\hat{x}_{2,i}^{n_2-1} - x_{2,n_2}) \right) \right).$$

So

$$= \frac{1}{2} \left( \frac{\partial}{\partial u} \min(\hat{x}_{2,i}^{n_2-1} (1 - \text{sign}(\min(\hat{x}_{2,i}^{n_2-1} - x_{2,n_2}))) + x_{2,n_2} (1 + \text{sign}(\min(\hat{x}_{2,i}^{n_2-1} - x_{2,n_2}))) \right) \quad \text{with } \frac{\partial}{\partial u} \min(\hat{x}_{2,i}^1) = x_{2,1}. \tag{4}$$

Likewise, we can write:

$$\frac{\partial}{\partial \mathbf{u}} b(\mathbf{u}) = \frac{\partial}{\partial \mathbf{u}} \max(\hat{x}_{1,i})_1^{n_1^2}$$

and from (3)

$$= \frac{1}{2} \left( \frac{\partial}{\partial \mathbf{u}} \max(\hat{x}_{1,i})_1^{n_1-1} + \mathbf{x}_{1,n_1} + \text{sign}(\max(\hat{x}_{1,i})_1^{n_1-1} - \mathbf{x}_{1,n_1}) \left( \frac{\partial}{\partial \mathbf{u}} \max(\hat{x}_{1,i})_1^{n_1-1} - \mathbf{x}_{1,n_1} \right) \right).$$

So

$$= \frac{1}{2} \left( \frac{\partial}{\partial \mathbf{u}} \max(\hat{x}_{1,i})_1^{n_1-1} (1 + \text{sign}(\max(\hat{x}_{1,i})_1^{n_1-1} - \mathbf{x}_{1,n_1})) \right. \\ \left. + \mathbf{x}_{1,n_1} (1 - \text{sign}(\max(\hat{x}_{1,i})_1^{n_1-1} - \mathbf{x}_{1,n_1})) \right) \quad \text{with } \frac{\partial}{\partial \mathbf{u}} \max(\hat{x}_{1,i})_1^1 = \mathbf{x}_{1,1}. \quad (5)$$

Finally, using (1):

$$\frac{\partial}{\partial \mathbf{u}} d(\mathbf{u}) = \frac{\partial}{\partial \mathbf{u}} s(\mathbf{u}) - \frac{\partial}{\partial \mathbf{u}} b(\mathbf{u}). \quad (6)$$

## 5. Algorithm implementation

We give a simplified skeleton of the min-max algorithm, using (4), (5), (6). We must fix the gradient algorithm step, the parameter  $\varepsilon$  which controls the calculation accuracy, and the initial unit vector.

**Do**

**min1** =  $\hat{x}_{2,1}$ ; **grads** =  $\mathbf{x}_{2,1}$

**For**  $i = 2$  **to**  $n_2$

**min2** =  $\hat{x}_{2,i}$

**grads** =  $\frac{1}{2}((1 - \text{sign}(\text{min1} - \text{min2})) \mathbf{grads} + (1 + \text{sign}(\text{min1} - \text{min2})) \mathbf{x}_{2,i})$

**min1** =  $\min(\text{min1}, \text{min2})$

**End\_For**

**max1** =  $\hat{x}_{1,1}$ ; **gradb** =  $\mathbf{x}_{1,1}$

**For**  $i = 2$  **to**  $n_1$

**max2** =  $\hat{x}_{1,i}$

**gradb** =  $\frac{1}{2}((1 + \text{sign}(\text{max1} - \text{max2})) \mathbf{gradb} + (1 - \text{sign}(\text{max1} - \text{max2})) \mathbf{x}_{1,i})$

**max1** =  $\max(\text{max1}, \text{max2})$

**End\_For**

**gradd** = **grads** - **gradb**

**u** = **u** + **step** · **gradd**; **u** = **u** /  $\|\mathbf{u}\|$

**diff1** = **min1** - **max1**

**While** (**diff1** - **diff2**) >  $\varepsilon$

**If** (**diff1** < 0) **then write** " $\omega_1$  and  $\omega_2$  are not linearly separable" **End**

**Else**

**Write** " $\omega_1$  and  $\omega_2$  are linearly separable"

**Write** "Hyperplane equations:"  $\mathbf{u}'\mathbf{x} - \text{min1}$ ,  $\mathbf{u}'\mathbf{x} - \text{max1}$

**End\_Else**

### 6. Experimental results

The following are two examples of the use of this algorithm. In both configurations, each class contains 100 vectors coming from two normal populations of  $\mathbb{R}^2$ , where the variance-covariance matrices are set to the identity matrix.

**Configuration 1** (Figure 3): linearly separable classes (there exists a hyperplane separating the two classes with a zero estimated error probability).

$$\omega_1: \text{mean} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \omega_2: \text{mean} = \begin{pmatrix} 5 \\ 5 \end{pmatrix}.$$

parameters: step =  $10^{-1}$ ,  $\epsilon = 10^{-8}$ , total iteration number = 52.

This algorithm converges in an efficient way, concluding that the classes are actually linearly separable (Difference = + 1.674 > 0). This maximum difference gives a good indication of the absolute gap between the two classes.

It also gives two hyperplane equations (here they are straight lines) that give two linear boundaries for the classes  $\omega_1$  and  $\omega_2$ :

$$P_1: 0.711 x_1 + 0.703 x_2 - 4.373 = 0, \quad P_2: 0.711 x_1 + 0.703 x_2 - 2.699 = 0.$$

**Configuration 2** (Figure 4): nonlinearly separable classes.

$$\omega_1: \text{mean} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \omega_2: \text{mean} = \begin{pmatrix} 3 \\ 3 \end{pmatrix}.$$

parameters: step =  $10^{-3}$ ,  $\epsilon = 10^{-8}$ , total iteration number = 1062.

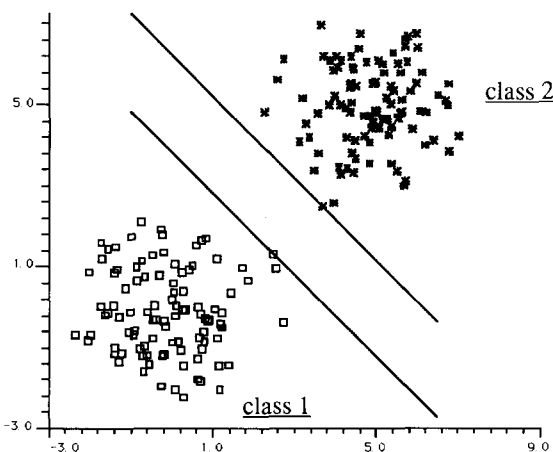


Figure 3. Direction = 7.1118283E-01, 7.0300710E-01. Direction = 7.1119996E-01, 7.0298977E-01. Result after 52 iterations: Min = 4.3736852E + 00, Max = 2.6990202E + 00, Difference = 1.6746650E + 00, Direction = 7.1119996E-01, 7.0298977E-01.

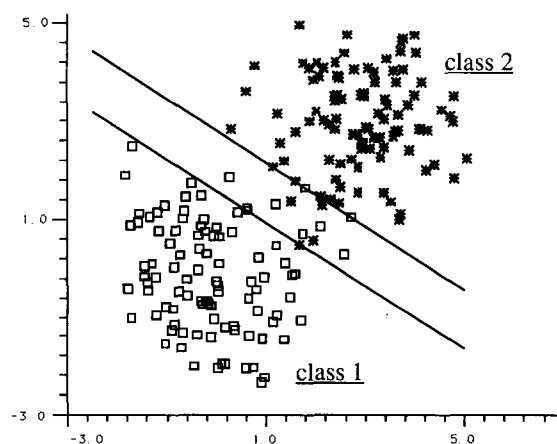


Figure 4. Direction = 5.4219291E-01, 8.4025404E-01. Direction = 5.4171130E-01, 8.4056461E-01. Result after 1062 iterations: Min = 1.3271238E + 00, Max = 2.3541315E + 00, Difference = - 1.0270076E + 00, Direction = 5.4171130E-01, 8.4056461E-01.

We point out that it was necessary to use a small step value ( $10^{-3}$ ). From that beginning, the algorithm converges, and then it oscillates around the solution. This approximate convergence allows us to conclude that the two classes are not linearly separable (Difference =  $-1.027 < 0$ ). This maximum difference gives a good indication of the classes overlapping.

## 7. Conclusions

This algorithm based on the new criterion  $d(\mathbf{u})$ , allows us to know, if the two classes are linearly separable or not. In the two cases, it gives a notion about an 'absolute distance' between the classes (gap or overlapping) and if they are linearly separable. It also gives two hyperplane equations which can be used in pattern recognition algorithms.

We noticed that the algorithm behaviour depends on the classes overlapping size: it would converge rapidly if the classes do not overlap.

In practice, we have to set the step value: if it is too big, the algorithm will converge first, and then will oscillate around the solution (with insufficient accuracy); if too small it will converge slowly. This heuristic step setting is the main default of this algorithm.

## References

- Rosenblatt, F. (1957). The Perceptron: a perceiving and recognizing automaton. Project PARA, Cornell Aeronaut. Lab. Rep. VG-1196-G-4.
- Ho, Y.C. and Kashyap, R.L. (1965). An algorithm for linear inequalities and its applications. *IEEE Trans. Electronic Computers* 14 (5), 683–688.
- McCulloch, W.S. and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *Bull. Math. Biophysics* 5, 115–133.
- Widrow, B. and Hoff, M.E. (1960). Adaptive switching circuits. 1960 *IRE WESCON Convention record*. New York, 96–104.