

A short horizontal bar with a teal segment on the left and an orange segment on the right.

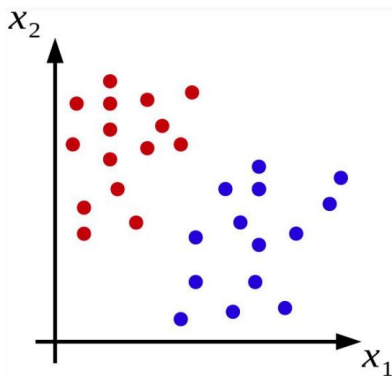
Machine Learning

Classifiers based on support vectors - Part 1

2023-2Q

Support Vector Machines (SVM)

It is a **supervised learning** algorithm used for both **classification** and regression problems .



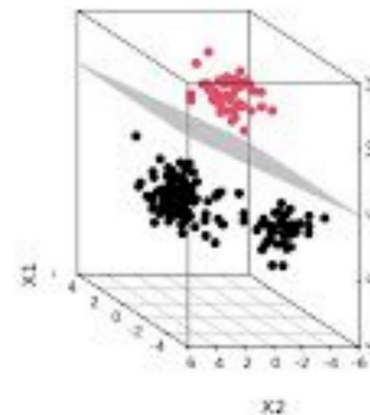
Its **main objective** is to find the **hyperplane optimal** that best **separates the different classes**.

Hyperplane concept

Equation of a hyperplane

In a **p-dimensional** space a hyperplane is defined by

$$b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p = 0$$



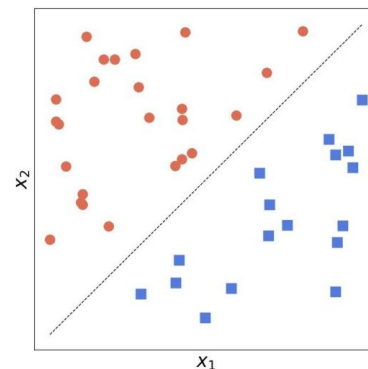
Equation of a hyperplane in R2

In **R2**, the line is defined by $b_0 + b_1 x_1 + b_2 x_2 = 0$

• (x_1, x_2) on the line • $x = (x_1, x_2)$

such that $b_0 + b_1 x_1 + b_2 x_2 > 0$

• (x_1, x_2) such that $b_0 + b_1 x_1 + b_2 x_2 < 0$



Linear separability

Suppose we have a set of **n examples** of **p attributes** x_i and **class** y_i ($1 \leq i \leq n$):

$$x_1 = (x_{1,1}, x_{1,2}, \dots, x_{1,p}),$$

$$y_1 \quad x_2 = (x_{2,1}, x_{2,2}, \dots, x_{2,p}), y_2$$

...

$$x_n = (x_{n,1}, x_{n,2}, \dots, x_{n,p}), y_n$$

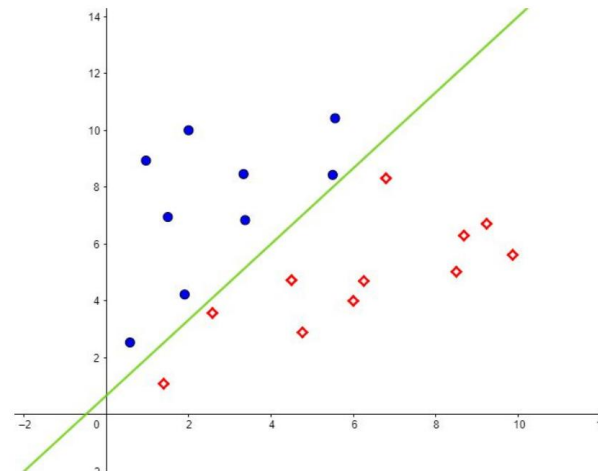
where each y_i belongs to $\{1, \dots, K\}$.

The goal is to find a classifier that correctly classifies each of these examples according to their class.

Linear separability

Separation Hyperplane

If we obtain a hyperplane such that all the examples whose class is -1 are on one side of the hyperplane and all examples whose class is $+1$ remain from the other we will have achieved the objective.



Linear separability



So ...

Given x_i , $i = 1, \dots, n$ if it happens that

$$b_0 + b_1 x_{i,1} + b_2 x_{i,2} + \dots + b_p x_{i,p} > 0 \text{ when } y_i = 1$$

and

$$b_0 + b_1 x_{i,1} + b_2 x_{i,2} + \dots + b_p x_{i,p} < 0 \text{ when } y_i = -1$$

then we can guarantee that $y_i (b_0 + b_1 x_{i,1} + b_2 x_{i,2} + \dots + b_p x_{i,p}) > 0$

Linear separability

$$y_i (b_0 + b_1 x_{i,1} + b_2 x_{i,2})$$

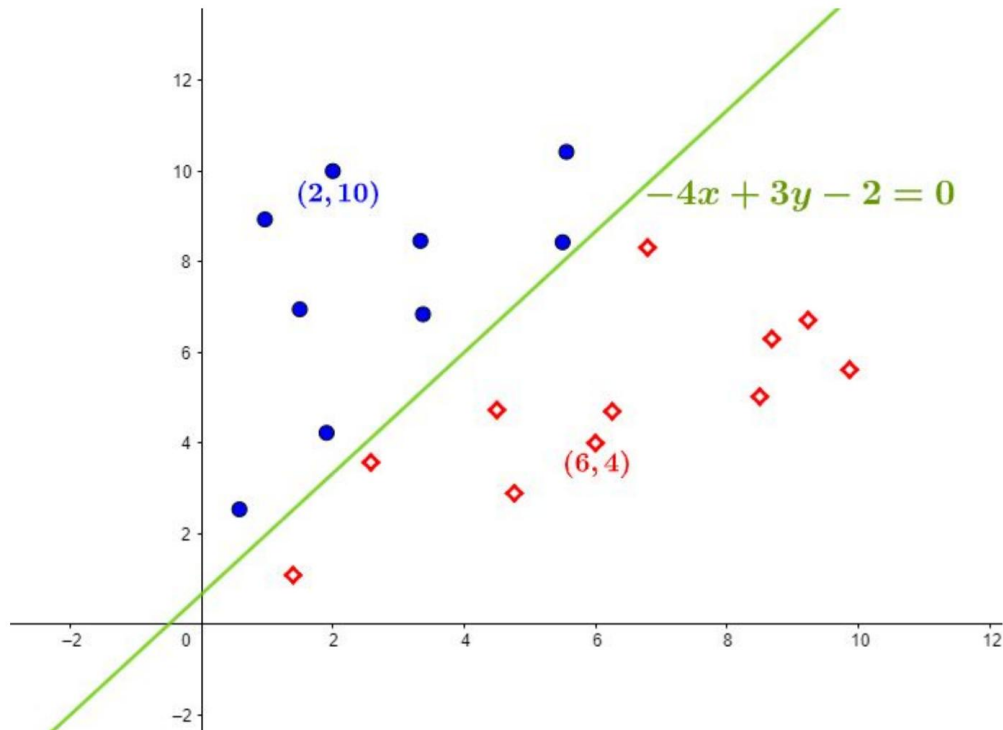
$$> 0 \quad y_i (-2 - 4x_{i,1} + 3x_{i,2}) > 0$$

$$1 * (-2 - 4*2 + 3*10) > 0$$

$$* (20) > 0 \quad 1$$

$$-1 * (-2 - 4*6 + 3*4) > 0$$

$$-1 * (-14) > 0$$



Linear separability



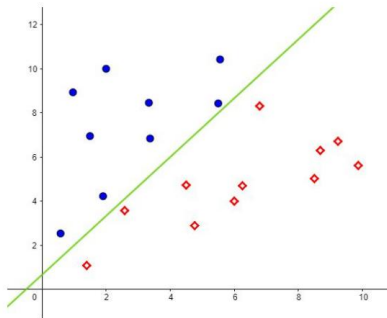
Given a test observation \mathbf{x}_{y} with p attributes we can classify it according to:

- \mathbf{x}_{y} will be of **class 1** if $b_0 + b_1 x_{\text{y}1} + b_2 x_{\text{y}2} + \dots$

$$b_0 + b_1 x_{\text{y}1} + b_2 x_{\text{y}2} + \dots > 0$$

- \mathbf{x}_{y} will be of **class -1** if $b_0 + b_1 x_{\text{y}1} + b_2 x_{\text{y}2} + \dots$

$$b_0 + b_1 x_{\text{y}1} + b_2 x_{\text{y}2} + \dots < 0$$



Linear separability

Let $f(\mathbf{x}) = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p$. Furthermore, we can say that:

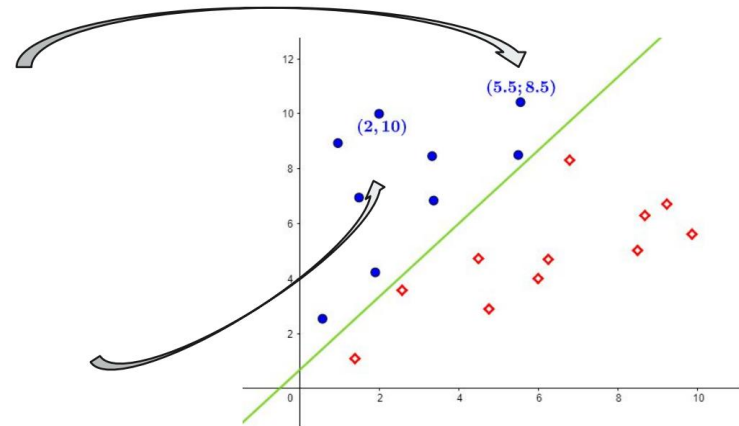
- If $f(\mathbf{x})$ is a value close to 0, \mathbf{x} will be close to the hyperplane
- If $f(\mathbf{x})$ is a value far from 0, \mathbf{x} will be far from the hyperplane.

$$-4.5 + 3.8 - 2 = -1.5$$

-1.5 near 0 \Rightarrow (5.5; 8.5) near the hyperplane

$$-4.2 + 3.10 - 2 = 20$$

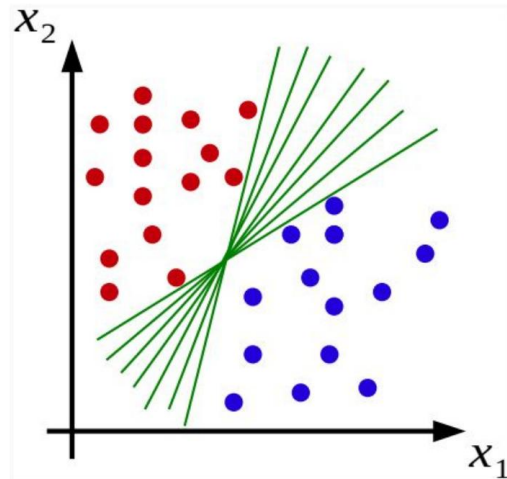
20 far from 0 \Rightarrow (2.10) far from the hyperplane



The maximum margin classifier

How do we separate the classes?

There may be more than one hyperplane separating them, in
In general, infinite planes that separate both classes.

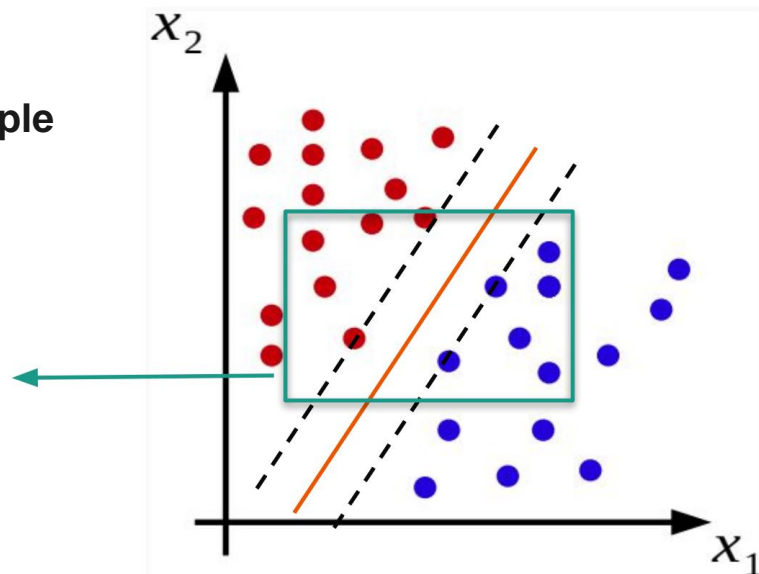
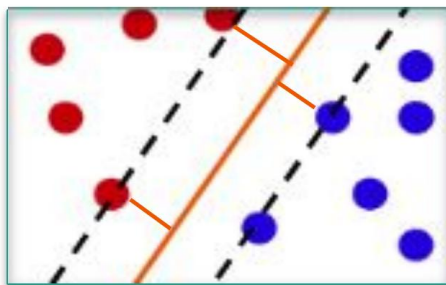


However, there is a hyperplane that has a particular property.

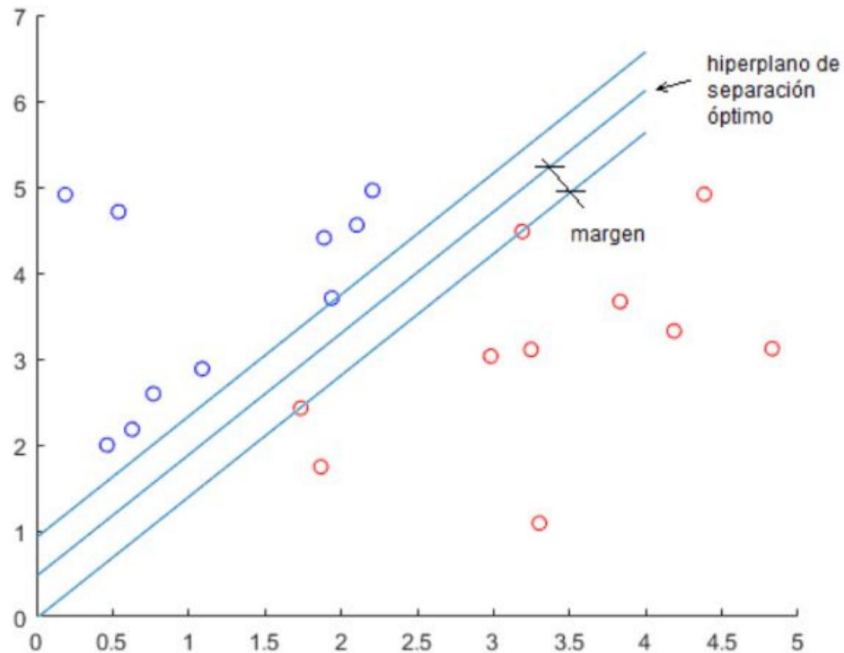
The maximum margin classifier

Let H be a hyperplane that separates both classes, let us consider the distances of each one of the examples to H .

We will define **margin** as the **distance of the example closest to H** .

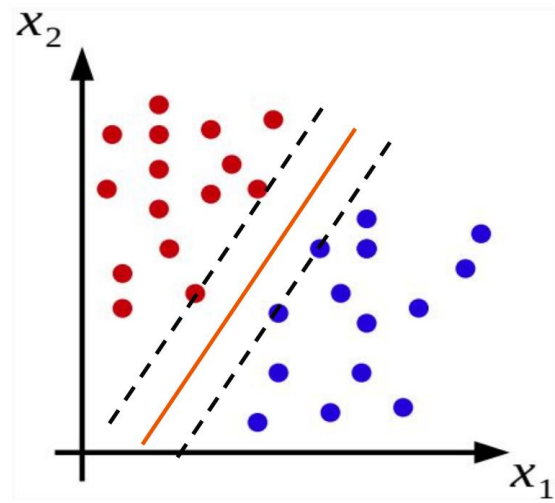
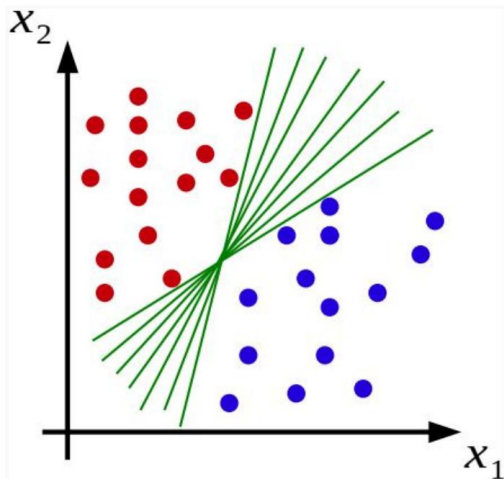
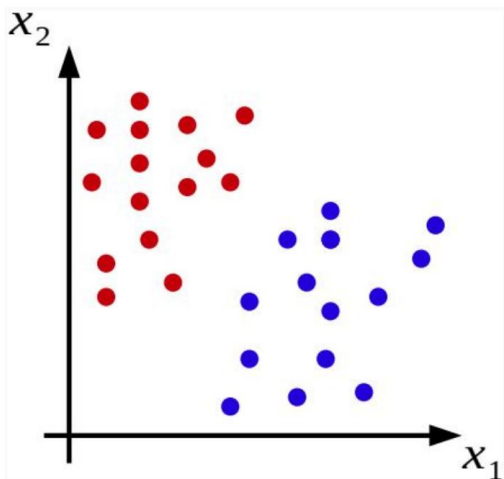


The maximum margin classifier



The maximum margin classifier

We are interested in the hyperplane that has a maximum margin, that is, **Hyperplane with maximum margin**
o **Optimal separation hyperplane.**

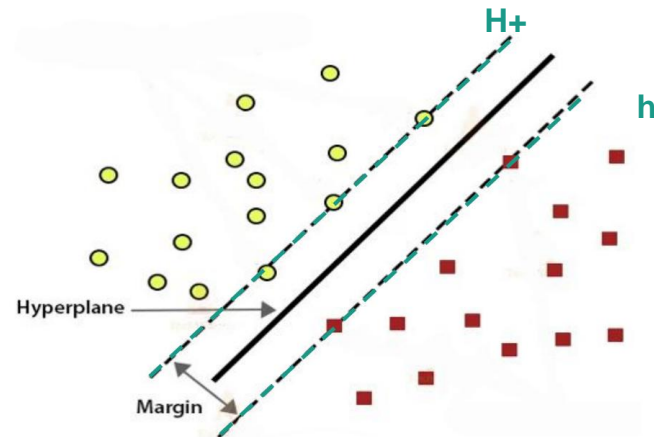


The maximum margin classifier

If we use the **maximal margin hyperplane** to separate both classes, the classifier is called **Maximal Margin Classifier**.

In addition, two hyperplanes are defined more, equidistant from H , H_+ and H_- .

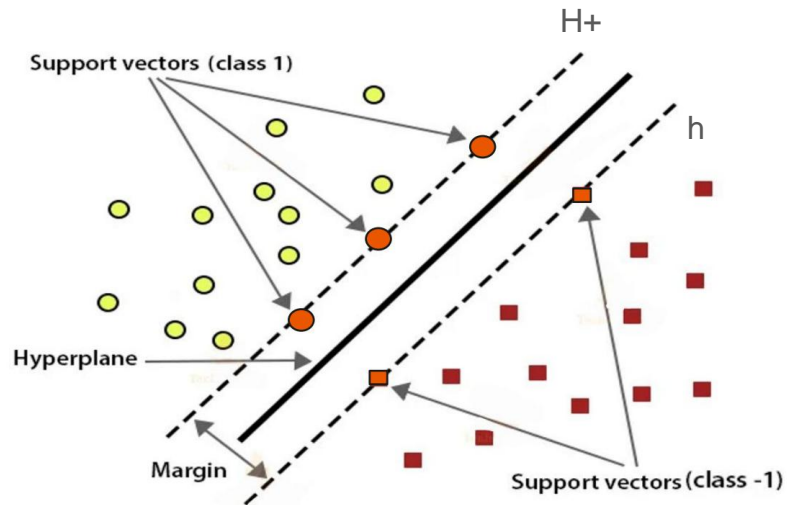
The distance between H_+ and H , and between H_- and H , is the same: **the margin**.



The maximum margin classifier

On H^+ and H^- you can see examples of both classes that are above them are called **support vectors**.

The **optimal separation hyperplane depends only on the support vectors** and not on the rest of the class examples.



The maximum margin classifier

Construction of the Maximum Margin Classifier

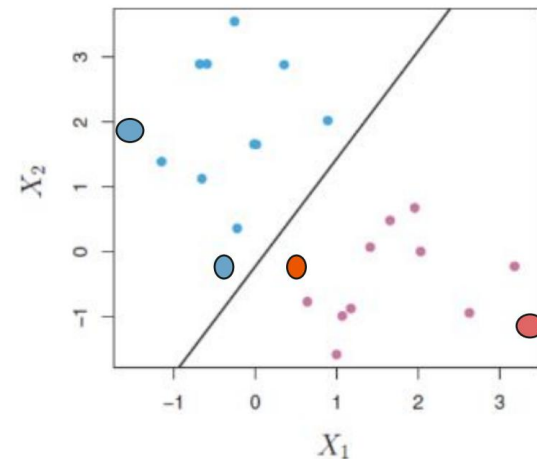
Let M be the margin (which we want to maximize) and since b defines the optimal separation hyperplane then what we want to do is find the values of b such that they maximize M subject to

- $y_i * (b_0 + b_1 x_{i,1} + b_2 x_{i,2} + \dots + b_p x_{i,p}) \geq M, \forall i, 1 \leq i \leq n,$
- $\sum_{j=1}^p b_j^2 = 1.$

Support Vector Classifier

The distance of a test observation from the optimal separation hyperplane gives us an idea of the **confidence** we can have in the classification.

- If the **distance is great** we will have **more confidence** (the observation is quite inside the class).
- If the **distance is small**, close to 0, we will have **less confidence** (the observation is close to the class boundary and therefore close to the other class).



Is the optimal separation hyperplane always optimal?

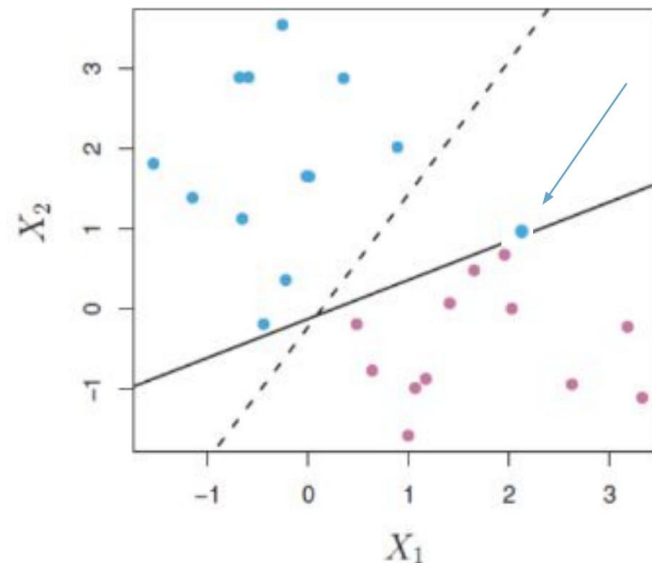


Suppose two classes that are well separated, that is

That is, they have a considerable margin.

We added an example that makes the margin
reduce considerably.

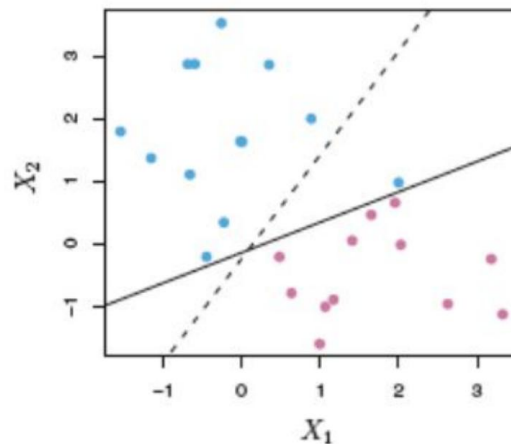
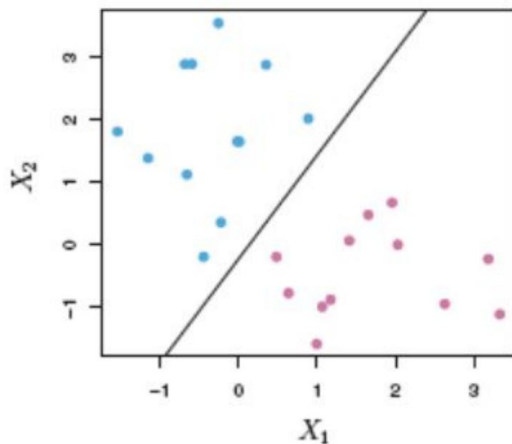
Examples that with the initial hyperplane would have been
ranked with greater confidence will now be
classified with less confidence.



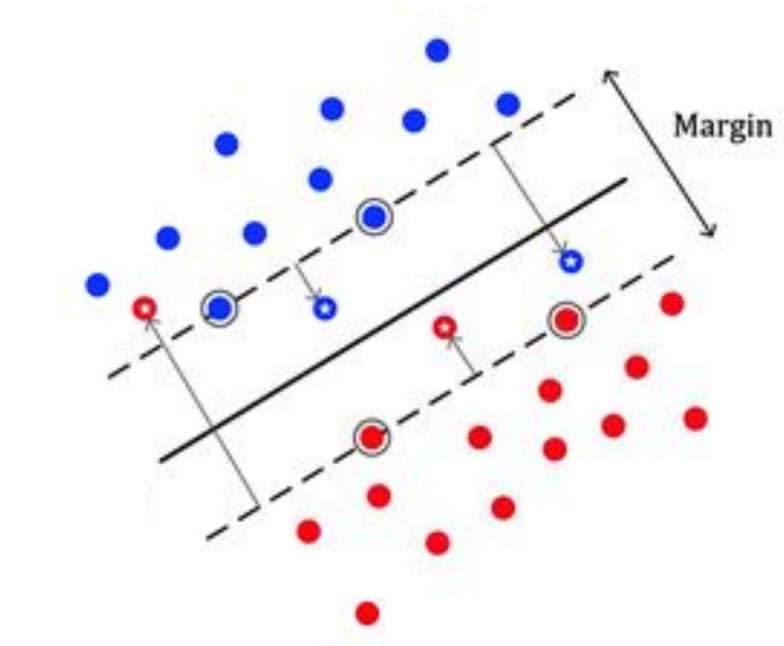
Is the optimal separation hyperplane always optimal?

Wouldn't it be worth using the initial hyperplane (the one that didn't take into account the example that makes the margin very small) instead of using the new hyperplane that divides but

Does it lead to a less reliable test?



Classifier with tolerant margin



Classifier with tolerant margin



Goals:

- That the individual observations are classified with robustness (that the distance to the hyperplane is not critical)
- Better classification of most training examples (assuming non-linearly separable classes).

Classifier with tolerant margin

Find the values of b and $\tilde{y}_1, \dots, \tilde{y}_n$ such that they maximize M subject to

- $y_i * (b_0 + b_1 x_{i,1} + b_2 x_{i,2} + \dots + b_p x_{i,p}) \geq M * (1 - \epsilon_i), \forall i, 1 \leq i \leq n,$
- $\sum_{j=1}^p b_j^2 = 1.$
- $\forall i, \epsilon_i \geq 0 \wedge \sum_{i=1}^n \epsilon_i \leq C.$

where C is a method tuning parameter.

Each \tilde{y}_i allows you to classify example x_i in the wrong place if necessary.

It could go through:

- be within the class margin
- being on the wrong side of the hyperplane.

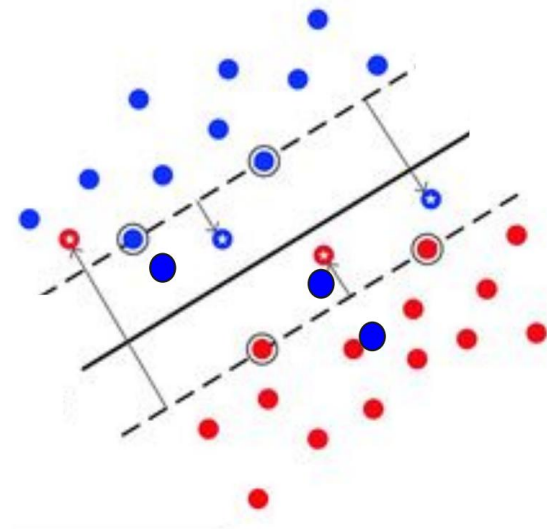
Classifier with tolerant margin

Given an observation x_i if

$$y_i * (b_0 + b_1 x_{i,1} + b_2 x_{i,2} + \dots + b_p x_{i,p}) \geq M * (1 - \epsilon_i), \forall i, 1 \leq i \leq n,$$

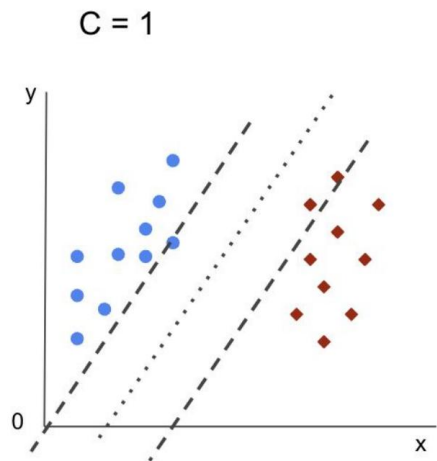
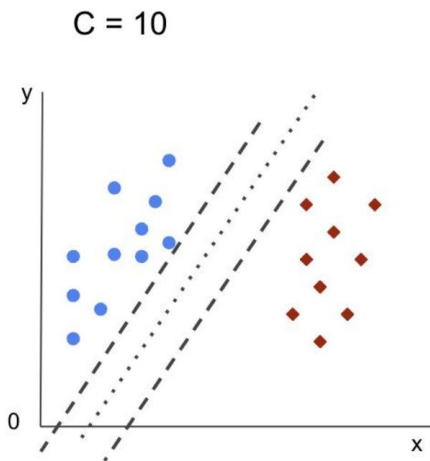
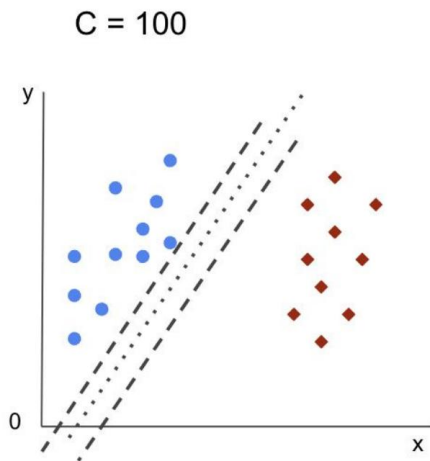
is satisfied for

- $\tilde{y}_i = 0$ ÿ the observation is on the **correct side** of the **margin**
- $\tilde{y}_i > 0$ ÿ the observation is on the **wrong side** of the **margin**
- $\tilde{y}_i > 1$ ÿ the observation is on the **wrong side** of the **hyperplane**



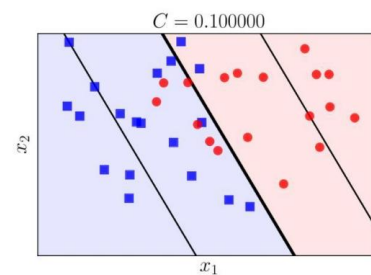
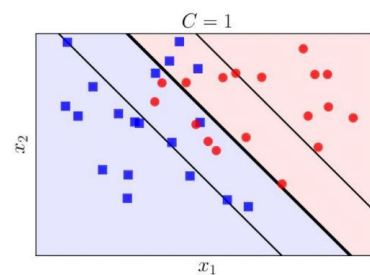
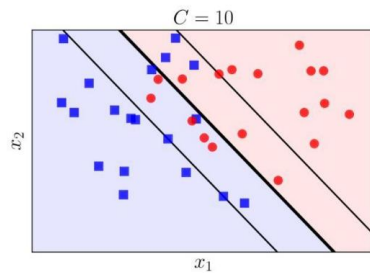
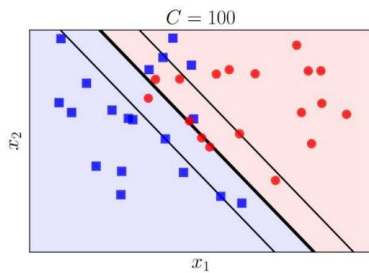
Classifier with tolerant margin

The value of **C** is a parameter that says how much the observations (as a whole) will be allowed to **violate the margin or hyperplane**.



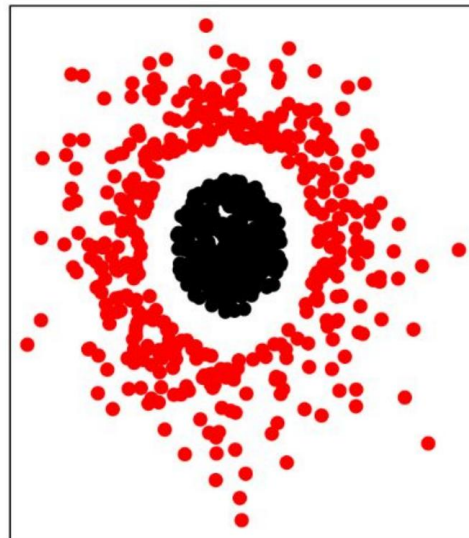
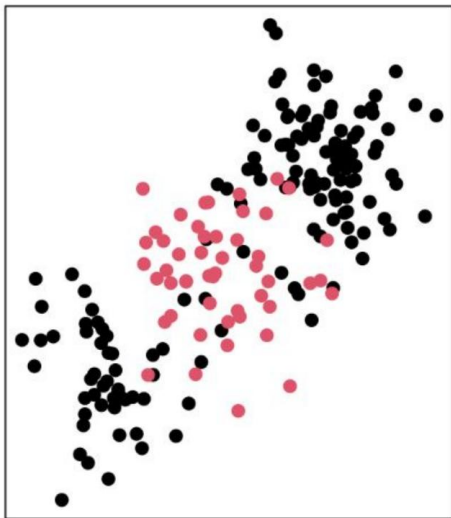
Classifier with tolerant margin

- If $C = 0$, the Tolerant Margin Classifier becomes a Margin Classifier maximal.
- One way to find C is with cross validation.



Classification with nonlinear decision boundaries

What we have seen so far is effective when the separation between classes is linear, but it doesn't work well in non-linear cases.



Classification with nonlinear decision boundaries

Let's consider a training set where

its examples x_i are of dimension $p = 2$ and its class is

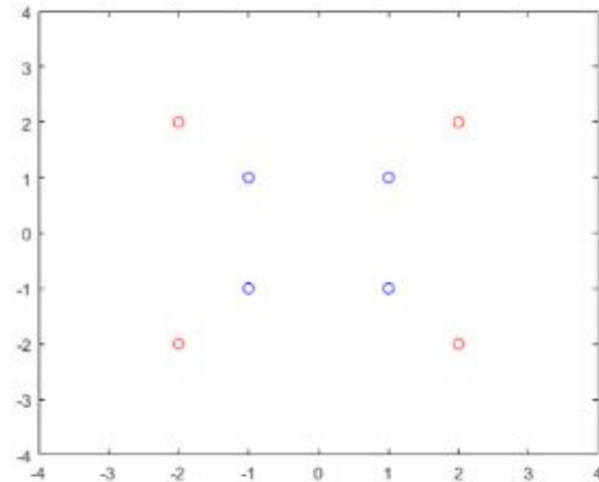
$y_i \in \{-1, 1\}$ (-1 in red and 1 in blue):

$X = \{(-2, 2), (-2, -2), (2, 2), (2, -2), (-1, 1), (-1, -1), (1, 1), (1, -1)\}$

$Y = \{-1, -1, 1, 1, 1, 1, 1, 1\}$

where a hyperplane cannot be established

linear separation between one class and the other.



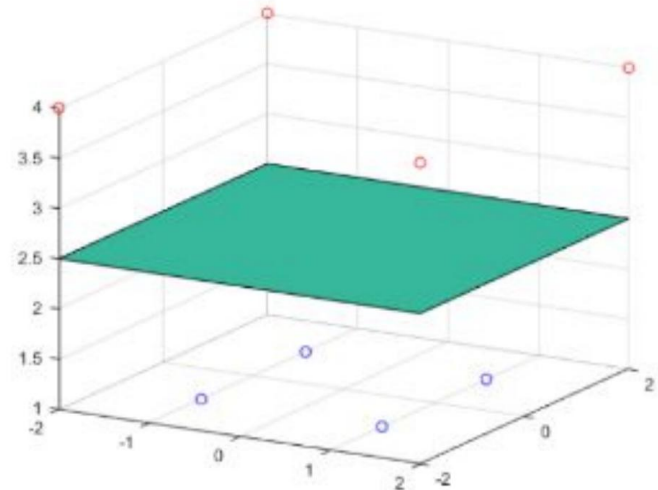
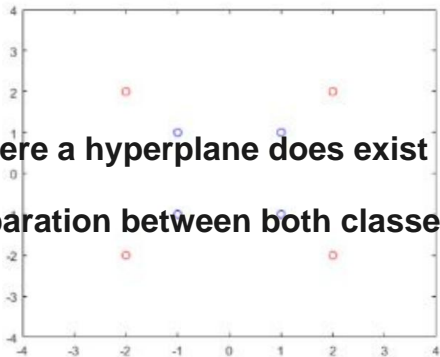
Classification with nonlinear decision boundaries

But if the examples $x_i = (x_{i1}, x_{i2})$ of _ _ _ _

$$X = \{(\ddot{y}2, \ddot{y}2, 4), (-2, 2, 4), (2, \ddot{y}2, 4), (2, 2, 4), (-1, \ddot{y}1, 1), (\ddot{y}1, 1, 1), (1, \ddot{y}1, 1), (1, 1, 1)\}$$

$$Y = \{\ddot{y}1, \ddot{y}1, \ddot{y}1, \ddot{y}1, 1, 1, 1, 1\}$$

Where a hyperplane does exist
separation between both classes.



Classification with nonlinear decision boundaries

For example

If we have examples in a p dimension $x_{i1}, x_{i2}, \dots, x_{ip}$ we could represent them in a $2p$ dimension according to

$x_{i1}, x_{i1}^2, x_{i2}, x_{i2}^2, \dots, x_{ip}, x_{ip}^2$ where there could be a hyperplane of dimension $2p-1$ that will separate them.

Classification with nonlinear decision boundaries

In this case, the problem to be solved is to find the values of $b = b_0, b_{11}, b_{12}, \dots, b_{p1}, b_{p2}$ and y_1, \dots, y_n such that they maximize M subject to

- $y_i * (b_0 + \sum_{j=1}^p b_{j1} * x_{ij} + \sum_{j=1}^p b_{j2} * x_{ij}^2) \geq (M - \epsilon_i), \forall i, 1 \leq i \leq n$
- $\sum_{j=1}^p \sum_{k=1}^2 b_{jk}^2 = 1.$
- $\epsilon_i \geq 0, \forall i, 1 \leq i \leq n \wedge \sum_{i=1}^n \epsilon_i \leq C.$

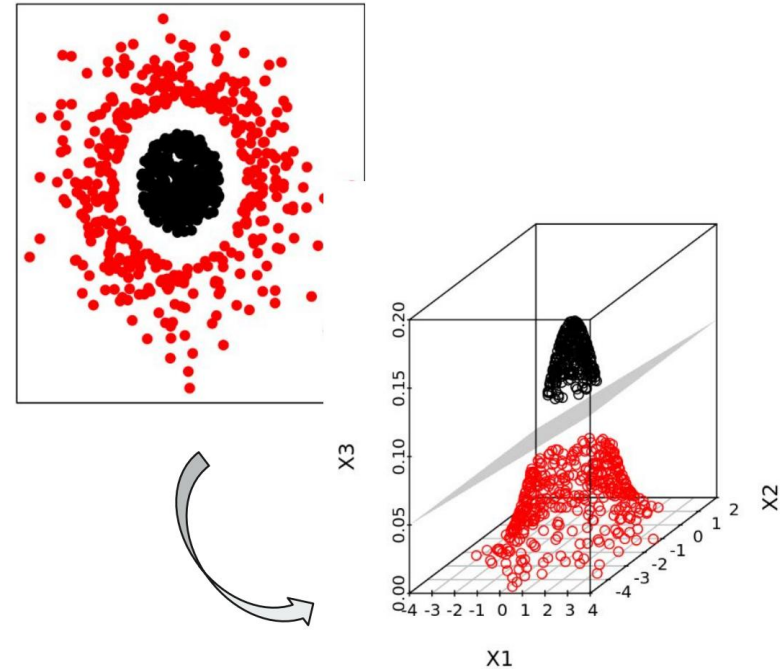
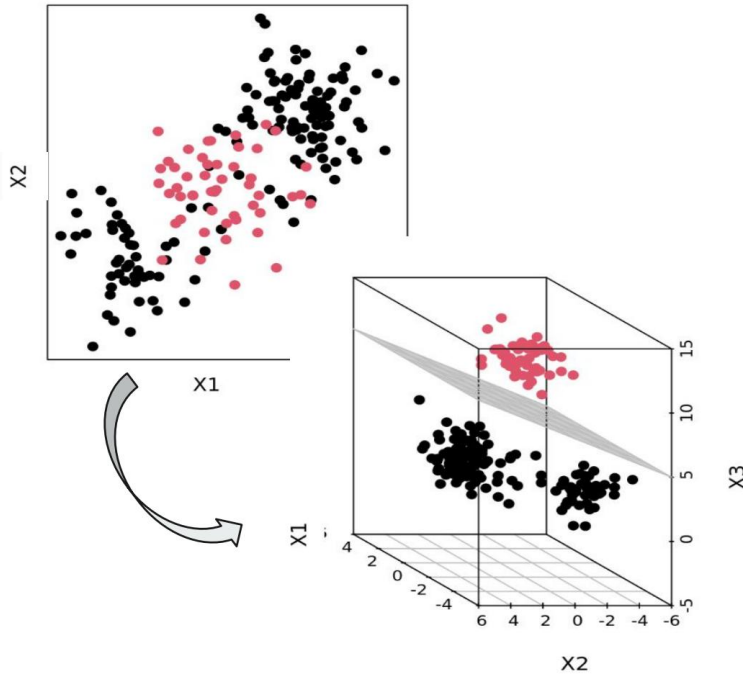
where C is a method tuning parameter.

Classification with nonlinear decision boundaries

In the example, we use x_{eh}^2 but we could use another degree or another function.

The idea is to obtain a space where there is linear separability, which implies nonlinear separability in the original space of the examples of dimension p .

Classification with nonlinear decision boundaries



Support Vector Machine



This proposal is a generalization of the Classification with limits of non-linear decisions and the way to generalize this idea is by introducing the concept of **Kernel**.

Support Vector Machine



The maximum margin classifier only depends on the vectors, y_k such that:

support, then, y, y_1, \dots

Given an observation x , if we want to know which class it belongs to, we calculate $f(x)$ as:

$$f(x) = b_0 + \sum_{k=1}^n y_i y_k x_i$$

where the x_i are the support vectors.

Support Vector Machine



We can write $f(x)$ as:

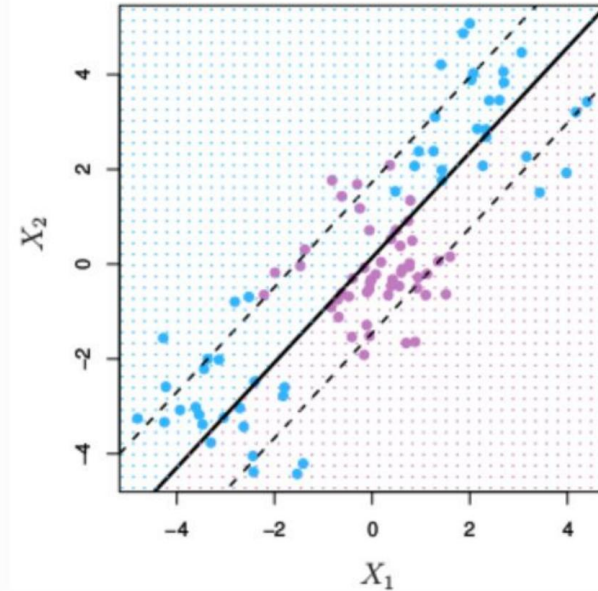
$$f(x) = b_0 + \sum_{i=1}^k \alpha_i K(x, x_i)$$

where $K(x, x_i) = \langle x, x_i \rangle$ and **K** is called **the Kernel**.

Support Vector Machine

Linear core

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{x} \cdot \mathbf{x}'$$



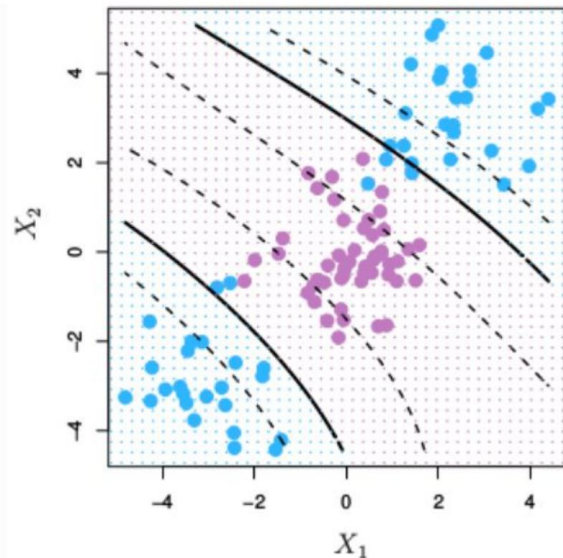
Support Vector Machine

Polynomial kernel

$$K(x', x_i) = \left(1 + \sum_{j=1}^p x_{ij}x'_j\right)^d$$

where d is the degree of the polynomial.

As d increases there will be more flexibility to find a linear separation in the new space of the examples (the expanded space).



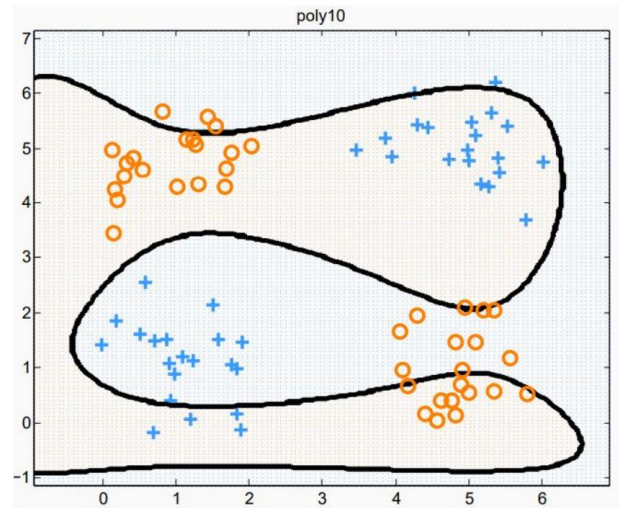
Support Vector Machine

Polynomial kernel

$$K(x', x_i) = \left(1 + \sum_{j=1}^p x_{ij}x'_j\right)^d$$

where d is the degree of the polynomial.

As d increases there will be more flexibility to find a linear separation in the new space of the examples (the expanded space).

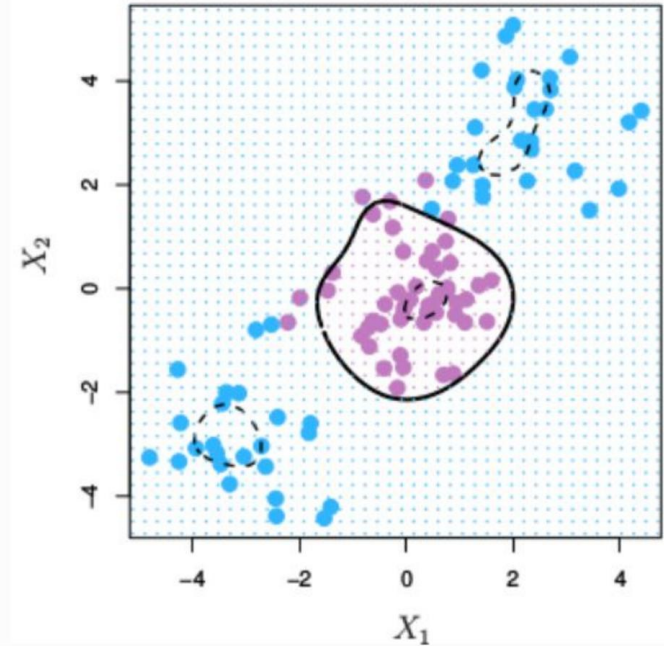


Support Vector Machine

Radial core

$$K(x', x_i) = e^{-\gamma \sum_{j=1}^p (x_{ij} - x'_j)^2}$$

where γ is a positive constant.

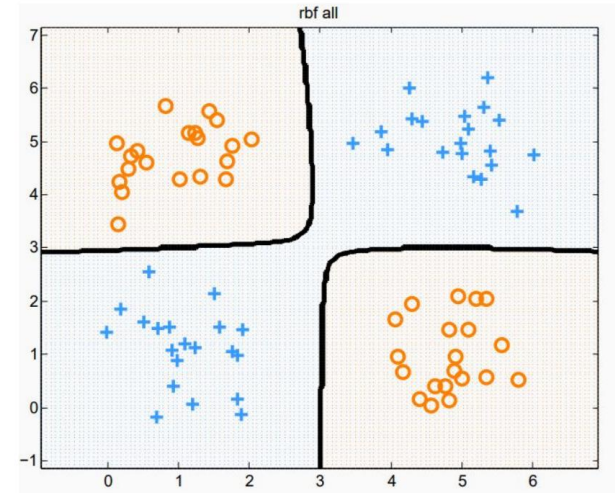


Support Vector Machine

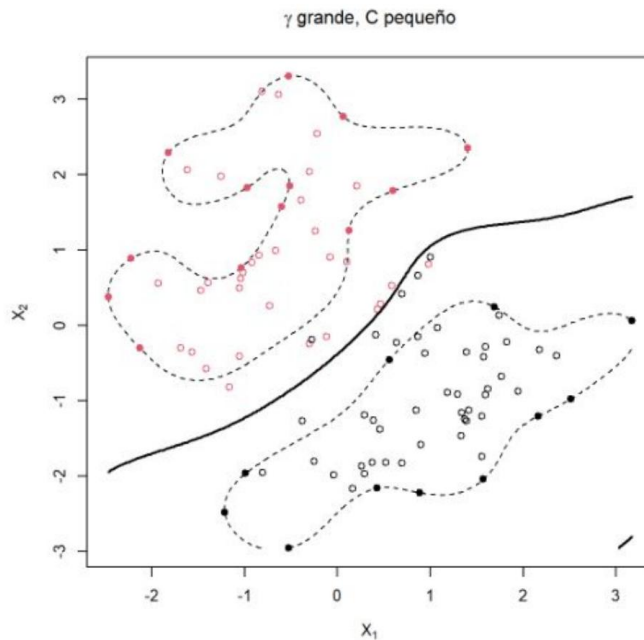
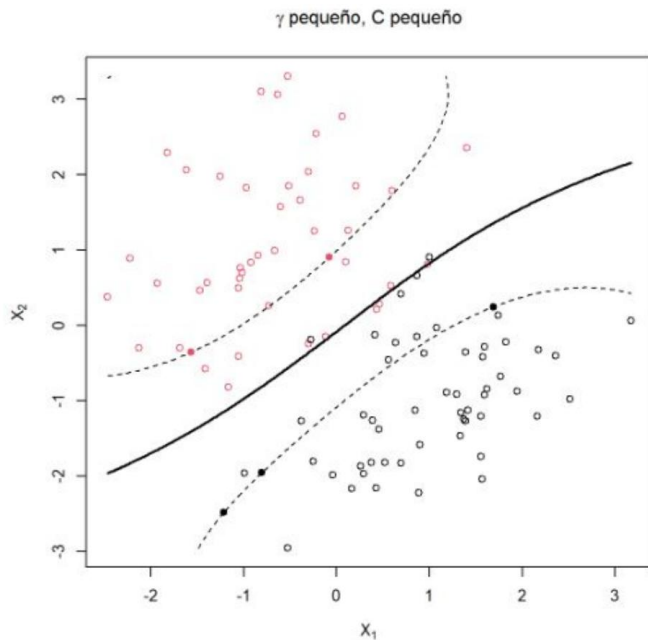
Radial core

$$K(x', x_i) = e^{-\gamma \sum_{j=1}^p (x_{ij} - x'_j)^2}$$

where γ is a positive constant.



Support Vector Machine



Support Vector Machine



Even if we change the Core, the formulation for the problem does not change.

The calculation of $f(x)$ will continue to be:

$$f(x) = b_0 + \sum_{i=1}^n \alpha_i K(x, x_i)$$

Support vector machine multiclass



What happens when our training set has more than 2 classes?

There are two approaches:

- One against one
- One against the rest

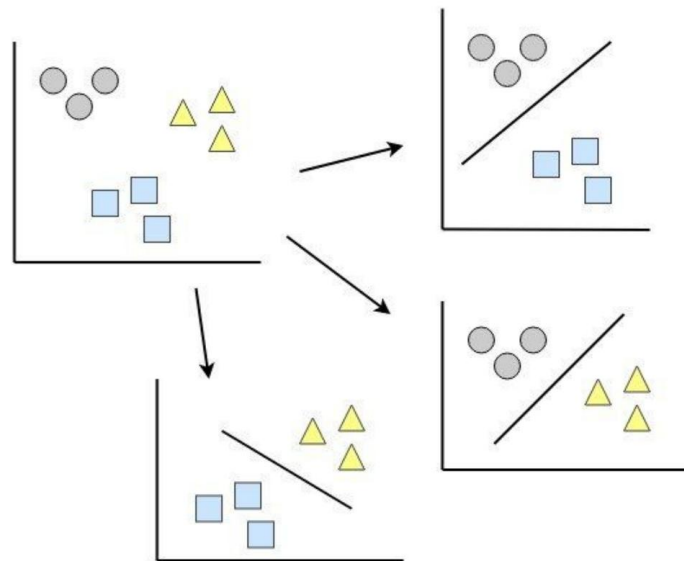
Support vector machine multiclass

One against one

An SVM is constructed for each distinct pair of classes (one class will be assigned +1 and the other -1).

The rest of the examples from the remaining classes will be ignored.

When a new observation is presented for each constructed SVM, a response is obtained and the class with the most votes is chosen.



SVM - Support multiclass vector machine

one against all

An SVM is constructed for each class (the class examples will be assigned $y_a + 1$ and the rest of the training set examples will be assigned $y_a - 1$).

When a new observation of each SVM constructed is presented, a

answer and the resulting class will be the one that corresponds to the SVM whose $f(x)$ is greater.

