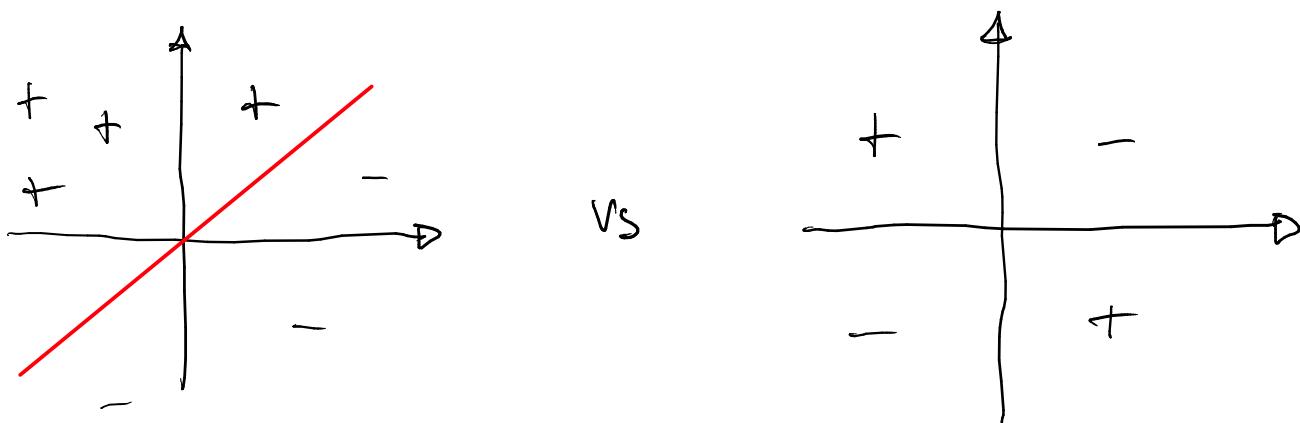


## 7. Linear model for classification

Learning a function  $f: X \rightarrow Y$  with  $X \subseteq \mathbb{R}^d$  and  $Y = \{c_1, \dots, c_K\}$ . The idea of these methods is to consider the ML problem as an optimization problem and find the solution to this one. We assume that data are linearly separable, this is a particular feature of the data.

INSTANCES are LINEARLY SEPARABLE if

$\exists$  a hyperplane that divide the instance space into two regions such that differently classified instances are separated

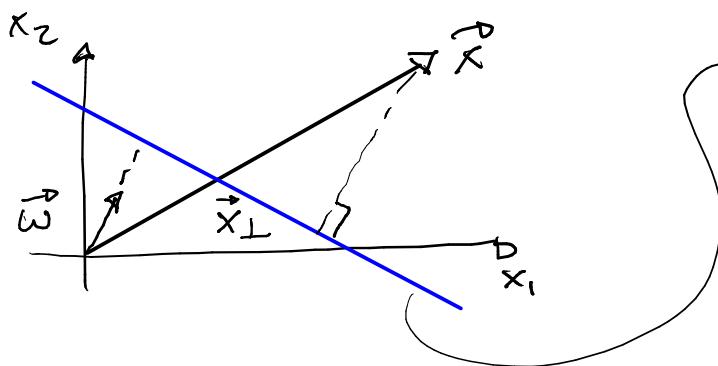


Graphical representation of the dataset  $x \subseteq \mathbb{R}^2$  and  $y = \{+, -\}$

If you have more than two classes the model is described by  $K$ -different functions, the multi-class function is:  $y_k(\vec{x}) = \vec{w}_k^\top \vec{x} + w_{k0}$   
In two class we need only one function:

$$y(\vec{x}) = \vec{w}^\top \vec{x} + w_0 = \vec{w}^\top \vec{x} + w_0 \text{ with } \vec{w} = \begin{pmatrix} w_0 \\ \vec{w} \end{pmatrix}, \vec{x} = \begin{pmatrix} 1 \\ \vec{x} \end{pmatrix}$$

Let's see the geometrical interpretation of those linear functions, starting by the binary classification problem:



line in which  $y(\vec{x}) = 0$

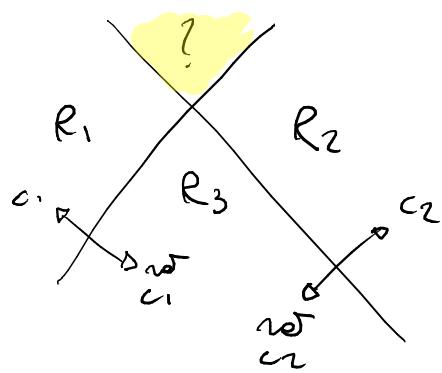
$\vec{w}$  is  $\perp$  to this line.

This line is partitioning the space and on the sides you have different points depending on  $y(\vec{x}) > 0$  or  $< 0$

We want to find the best line since it will be used to classify.

When we have  $k$ -classes, with  $k \geq 2$  we cannot use simple tricks

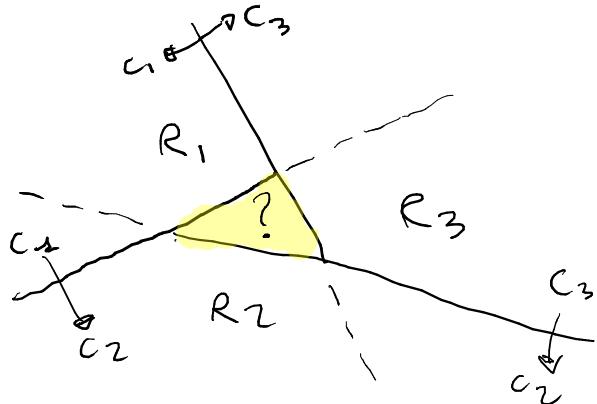
1 - VS - REST



You use binary classifier  $(k-1)$  in general. The problem is that there is a region in which we cannot determine the right class.

1 - VS - 1

Quadratic number of classifiers  $k(k-1)/2$  and you classify by considering each possible pair



We still have a region of uncertainty.

We can transform  $k$ -class classifier into a set of binary classifiers.

We have to change the model:

$$y_k(\vec{x}) = \vec{w}_k^\top \vec{x} + w_{k0} \quad \text{compute the } \vec{W}$$

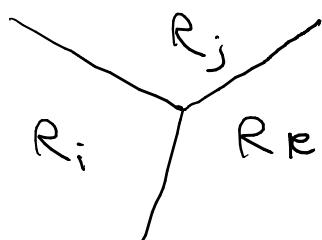
Then assign to a new instance  $\vec{x}$  class  $C_i$  if:

$$\boxed{y_k(\vec{x}) > y_j(\vec{x}) \quad \forall j \neq k}$$

Given this formulation we can derive the equation of the hyperplane separating the two classes  $C_k$  and  $C_j$ :

$$(\vec{w}_k - \vec{w}_j)^\top \vec{x} + (w_{0k} - w_{0j}) = 0$$

And the representation is like this:



This separates correctly the regions:  
The output of the model is a partition of the space. If you don't model the problem like this you have trouble!

Compute solution:

$$y_k(\vec{x}) = \vec{w}_k^\top \vec{x} + w_{0k} \quad \vec{x} \text{ not in dataset}$$

$$\hat{y}(\vec{x}) = \begin{pmatrix} y_1(\vec{x}) \\ \vdots \\ y_k(\vec{x}) \end{pmatrix} = \begin{pmatrix} \vec{w}_1^\top \vec{x} \\ \vdots \\ \vec{w}_k^\top \vec{x} \end{pmatrix} = \vec{W}^\top \vec{x}$$

$$\vec{w}_k = \begin{pmatrix} w_{0k} \\ \vec{w}_k \end{pmatrix} \quad \vec{x} = \begin{pmatrix} 1 \\ \vec{x} \end{pmatrix} \text{ dummy component to get } w_{0k}$$

Given a multi-class classification problem and a dataset  $D$  with linearly separable data, determine  $\vec{W}$  such that  $\vec{y}(\vec{x}) = \vec{W}^\top \vec{x}$  is the  $k$ -class decision

## LEAST SQUARES

We have a dataset  $D$  with inputs and outputs,  $D = \{(\vec{x}_n, \vec{t}_n)\}_{n=1}^N\}$  and we use a particular convention to represent  $\vec{t}$ , in particular we use the 1-of-K coding scheme (only one component is 1):

$$x \in C_K \rightarrow t_K = 1, t_j = 0 \quad \forall j \neq K, \vec{t}_n = (0 \dots 1 \dots 0)^T$$

So now we can represent our dataset in a more formal:

$$\tilde{X} = \begin{pmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_N^T \end{pmatrix} \quad T = \begin{pmatrix} \vec{t}_1^T \\ \vdots \\ \vec{t}_N^T \end{pmatrix}$$

$\tilde{X}$  design matrix

We define now an error function, where the error is the difference from the predictions and the value that is in the dataset.

We consider the LEAST SQUARES ERROR FUNCTION:

$$E(\tilde{W}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{X} \tilde{W} - T)^T (\tilde{X} \tilde{W} - T) \right\}$$

$\downarrow$  prediction on the input       $\downarrow$  real value of the input

$E(\tilde{W}) = 0$  when  $\tilde{W}$  predicts the same value of  $T$  for each value in the dataset (for each sample)

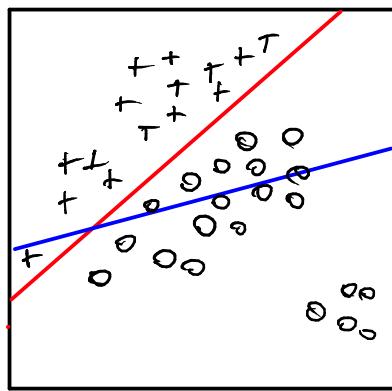
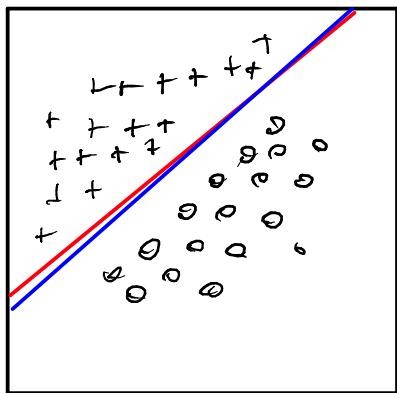
TASK : We want to find the value of  $\tilde{W}$  that minimizes the error. This is a quadratic problem, & can be solved with a close form solution:

$$\tilde{W} = \underbrace{(\tilde{X}^T \tilde{X})^{-1}}_{\tilde{X}^\# \text{ pseudoinverse of } X} \tilde{X}^T T$$

$$\vec{y}(\vec{x}) = \vec{\tau}^T (\tilde{\vec{x}}^*)^T \tilde{\vec{x}} \rightarrow \text{simple and easy.}$$

INCONVENIENTS: Works very well for some dataset but is NOT ROBUST to OUTLIERS.  
 (outliers = points not coming from the same distribution)

Least square considers a way of minimizing the distance from the computed line.



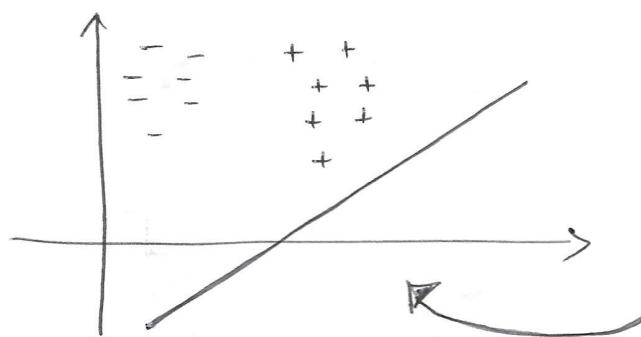
Outliers affect the line.

## Fischer's linear discriminant

(5)

For this problem we consider only 2 classes and a bidimensional space. We want to determine  $y = \vec{w}^T \vec{x}$  and classify  $\vec{x} \in C_1$  if  $y \geq -w_0$ ,  $\vec{x} \in C_2$  if  $y < -w_0$ .

The geometry of  $\vec{w}^T \vec{x}$  is just a projection of  $\vec{x}$  on the line determined by  $\vec{w}$ . One way to compute the separation line is the following:



We can consider all projections of the points on a line. According to different lines the projection is diff.  
IF THE LINE IS NOT "WELL-ORIENTED" THE PROJECTIONS OVERLAP!

We look for a line that minimize the overlapping or better that MAXIMIZE THE SEPARATION of the points! You can consider - the center of mass for each class (centroids / means);

$$\vec{m}_1 = \frac{1}{N_1} \sum_{m \in C_1} \vec{x}_m \quad \vec{m}_2 = \frac{1}{N_2} \sum_{m \in C_2} \vec{x}_m$$

Choose  $\vec{w}$  that maximizes  $J(\vec{w}) = \vec{w}^T (\vec{m}_2 - \vec{m}_1)$  s.t.  $\|\vec{w}\| = 1$ . You project only the mean and you do it for all possible function and looking for the one that maximizes separation.

GOOD! BUT STILL NOT PERFECT because you consider only the means, you DO NOT CONSIDER THE COVARIANCE. If we consider the covariance you can rotate the line in such a way ~~so that~~ it aligns with covariances. COVARIANCE ALLOW THE MODEL TO ROTATE!

$$J(\vec{w}) = \frac{\vec{w}^T S_B \vec{w}}{\vec{w}^T S_W \vec{w}}$$

Fischer criterion.

$$S_B = \left( \vec{m}_2 - \vec{m}_1 \right) \left( \vec{m}_2 - \vec{m}_1 \right)^T$$

Between classes

$$S_W = \sum_{m \in C_1} (\vec{x}_m - \vec{m}_1)^T (\vec{x}_m - \vec{m}_1)$$

...  $m \in C_2$

$$S_W = \sum_{m \in C_2} (\vec{x}_m - \vec{m}_2)^T (\vec{x}_m - \vec{m}_2)$$

...  $m \in C_1$

$$\boxed{\vec{w}^* \propto S_w^{-1} (\vec{m}_2 - \vec{m}_1)}$$

⑥

$$J(\vec{w}) = \{(W^T S_w W)^{-1} (W^T S_B W)\} \text{ with } k \geq 2 \text{ classes.}$$

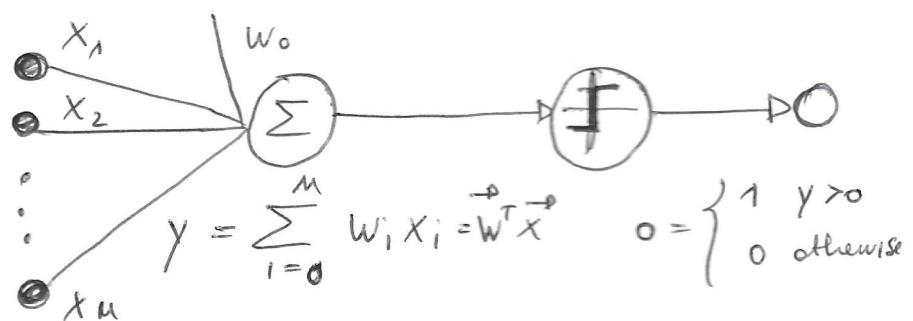
BETTER THAN ~~LEAST~~ LEAST SQUARES, more robust to outliers and better with more than two classes.

## PERCEPTRON

BASIC UNIT FOR A NEURAL NETWORK.

No close form, it will use an iterative method. We still define an optimization problem, but we solve it with an iterative process.

$$x_0 = 1$$



The model is defined as a linear combination of the input dimensions with some weights. We have another portion in which we apply the sign function, that is discontinuous.

OUR GOAL IS TO FIND  $w_i$  BY USING AN ITERATIVE PROCESS, At some in the execution of the algorithm you receive new information, in particular new samples from dataset, and you want to update  $w_i$ :

## TRAINING RULE

$$w_i \leftarrow w_i + \Delta w_i, \quad \Delta w_i = \eta (t - o) x_i$$

↑ value in D      ↑ prediction

The bigger is the error  
the bigger is the update.

Unthreshold unit!

$$E(\vec{w}) = \frac{1}{2} \sum_{m=1}^N (t_m - o_m)^2 = \frac{1}{2} \sum_{m=1}^N \left( t_m - \vec{w}^T \vec{x}_m \right)^2$$

learn  $w_i$  to minimize the squared error.

To minimize we can use a gradient method, computing the derivative with respect to each component of  $\vec{w}$ ! ④

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_{m=1}^N (t_m - \vec{w}^T \vec{x}_m)^2 = \sum_{m=1}^N (t_m - \vec{w}^T \vec{x}_m) (-x_{i,m})$$

$$\Delta w_i = \Theta \eta \frac{\partial E}{\partial w_i} = -\eta \sum_{m=1}^N (t_m - \vec{w}^T \vec{x}_m) (-x_{i,m})$$

$\eta$  is a small constant called "LEARNING RATE". We have a minus because we want to go in the negative direction of the gradient and we want to minimize the error. We can apply the algorithm in:

### BATCH MODE

We compute the gradient over all samples in  $D$

or

### MINI-BATCH

use a subset  $S \subset D$

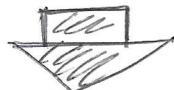


Correctly compute the gradient.

VS

### INCREMENTAL MODE

We consider only one sample at the time.

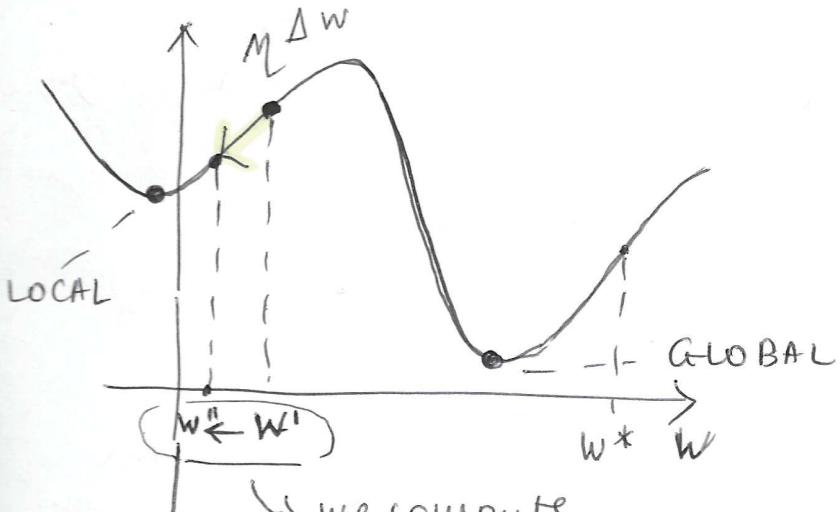


Better in ML because it converges faster!  
less sensitive to local minima.

### TERMINATION CONDITIONS:

- Predefined number of iterations
- Threshold on changes in the loss function  $E(\vec{w})$ .

Let's consider the following plot with  $w$  with only one dimension. (8)



THIS PROCESS WILL BRING TO THE LOCAL MINIMUM, NOT TO A GLOBAL

The solution of the iteration process depends on where you start. IF you start in  $w = w'$  you will get to a local minimum, if you start in  $w = w^*$  you will go to a global one.

### STEPS of PERCEPTRON :

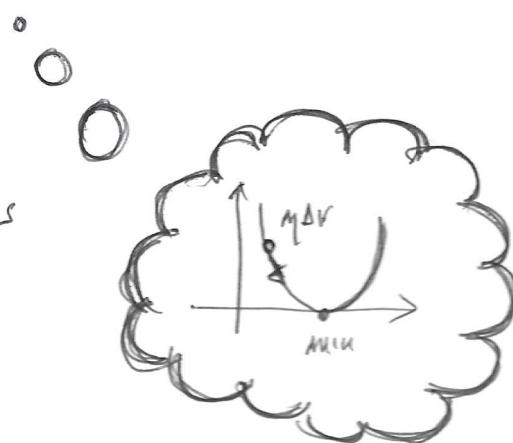
1. Initialize  $\hat{w}$  (small random values)
2. Repeat until termination:  
$$\hat{w}_i \leftarrow \hat{w}_i + \Delta w_i;$$
3. Output  $\hat{w}$ .

ANOTHER PROBLEM : The learning rate

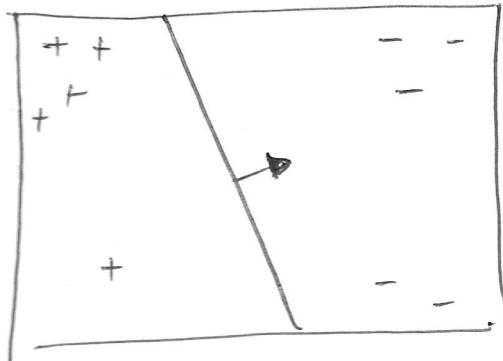
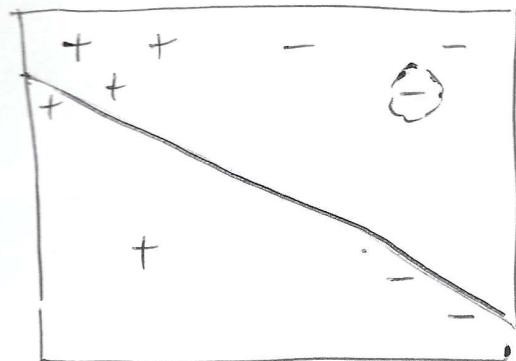
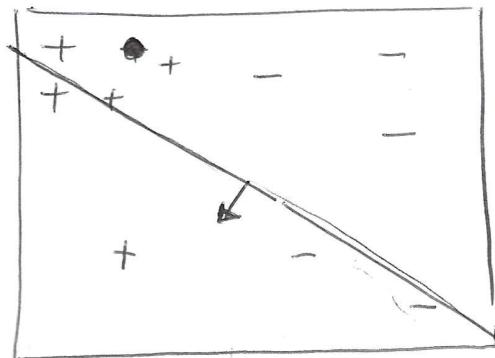
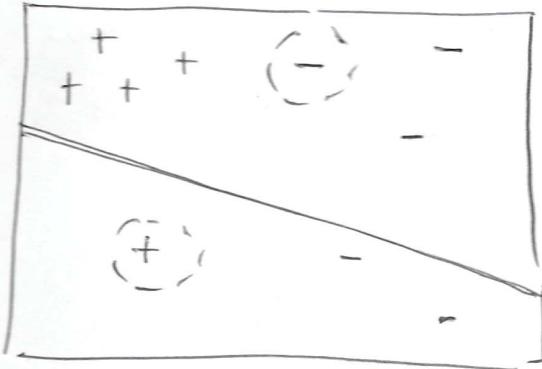
- if  $\eta$  is too small you will have too many steps
- if  $\eta$  is too big you may have oscillations and you can diverge

Convergence condition :

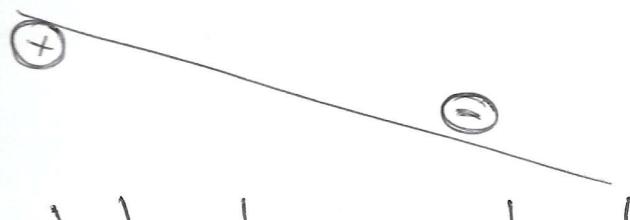
- data linearly separable
- $\eta$  small enough (usually means slow convergence)



Plot the evolution over time of perception. Remember (9) that every time we find one point for which the error is different than zero the parameter will change and the line is <sup>in</sup> some way attracted towards the data point. SINCE WE HAVE THE  $\eta$  TERM THE MOVEMENT IS LIMITED



The method stops when you find a line that correctly separates the dataset (the samples). This is true if  $\eta$  is small enough. IF  $\eta$  IS TOO SMALL YOU ALSO HAVE SOLUTION SENSITIVE TO NOISE, you will find a solution that has this shape



VERY RISKY | A line is close to the circles points!  
Usually this is not the line that separates the data.

PERCEPTRON IS THE ONLY METHOD THAT IS NOT AFFECTED BY OUTLIERS (since it no considers means, variances, ...)



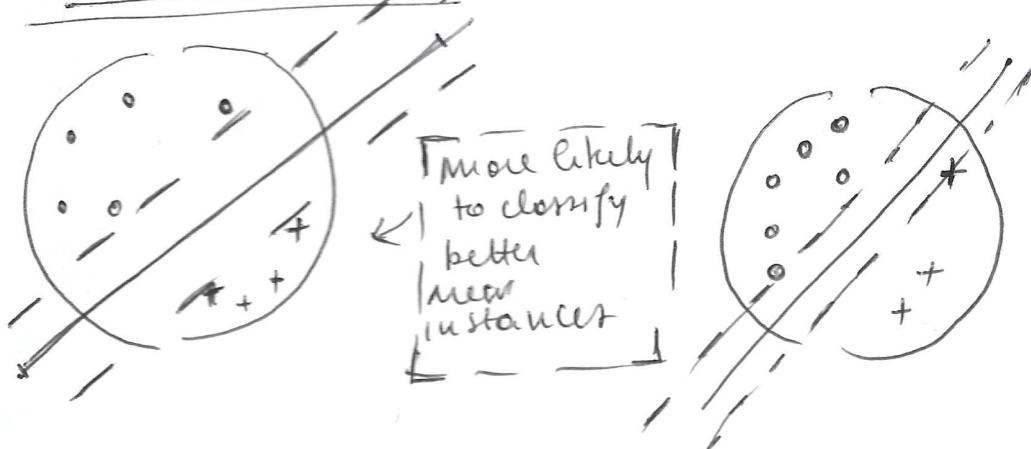
# SUPPORT VECTOR MACHINES

10

MAIN : Find the direction that guarantees the MAXIMUM  
GOAL : MARGIN between two categories.

Assumption: data linearly separable.

The separation surface should be as far as possible from the other points (Probability of classify new instances is good, is what we want).



MAXIMUM MARGIN HYPERPLANE IS WHAT SVM IS SUPPOSED TO FIND.

let's consider binary classification  $y: X \rightarrow \{-1, +1\}$ , with  $D = \{(\vec{x}_m, t_m)\}_{m=1}^N$ ,  $t_m \in \{-1, +1\}$  and a linear model:

$$y(\vec{x}) = \vec{w}^\top \vec{x} + w_0$$

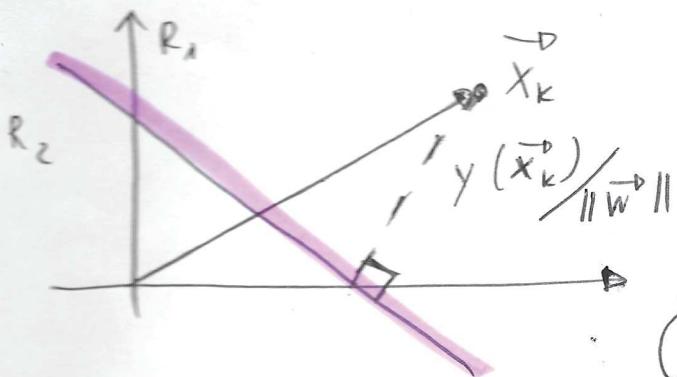
Assume  $D$  is linearly separable:

$$\exists \vec{w}, w_0 \text{ s.t. } \begin{cases} y(\vec{x}_m) > 0, & \text{if } t_m = +1 \\ y(\vec{x}_m) < 0, & \text{if } t_m = -1 \end{cases}$$

$$\Rightarrow t_m y(x_m) > 0 \quad \forall m = 1 \dots N$$

THE Margin is computed by considering the distance between the hyperplane and the closest point in the dataset.

The margin can be computed as  $\frac{\|y(\vec{x}_k)\|}{\|\vec{w}\|}$ . (11)



The closest point is the one that minimizes the distance between the point itself and the hyperplane.

$$\min_{m=1, \dots, N} \frac{|y(\vec{x}_m)|}{\|\vec{w}\|} = \text{Margin}$$

$$\min_{m=1, \dots, N} \frac{|y(\vec{x}_m)|}{\|\vec{w}\|} = \dots = \frac{1}{\|\vec{w}\|} \min_{m=1, \dots, N} [t_m (\vec{w}^T \vec{x}_m + w_0)]$$

$\vec{h} = \vec{w}^T \vec{x} + w_0$  is the hyperplane. So, what I want?

I WANT TO SELECT AMONG ALL POSSIBLE HYPERPLANES, THE ONE THAT MAXIMIZES THE MARGIN, aka the minimum distance between the hyperplane and any point in the dataset!

$h^* = \vec{w}^{*T} \vec{x} + w_0^*$  the hyperplane with maximum margin is computed as:

$$\vec{w}^{*T}, w_0^* = \underset{\vec{w}, w_0}{\operatorname{argmax}} \frac{1}{\|\vec{w}\|} \min_{m=1, \dots, N} [\vec{t}_m (\vec{w}^T \vec{x}_m + w_0)]$$

We can rescale all the points in such a way we have!  
(WE CAN SCALE BY A CONSTANT FACTOR; does not affect the solution)

$$t_k (\vec{w}^T \vec{x}_k + w_0) = 1$$

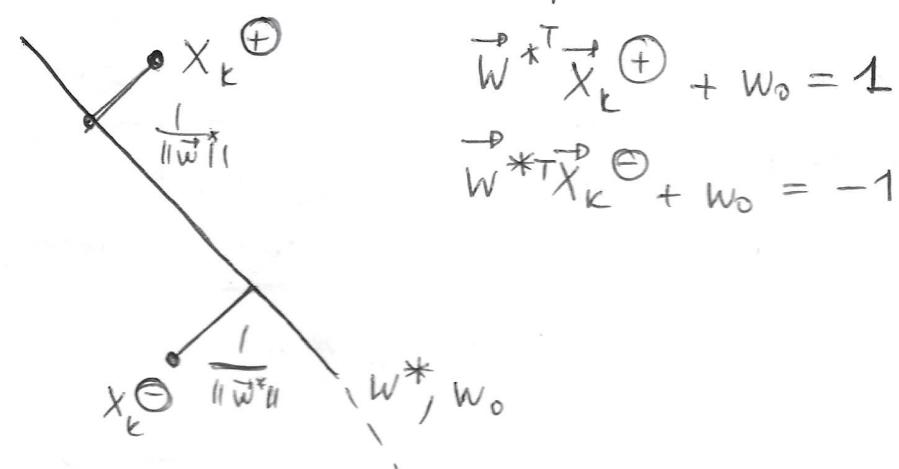
The margin is 1,  
 $\vec{x}_k$  is the closest point.

For all the other points the quantity is in general:

$$t_m (\vec{w}^T \vec{x}_m + w_0) \geq 1 \quad \forall m = 1, \dots, N$$

When we compute the optimal margin hyperplane we have at least 2 closest point (one for each side)!

(12)



The canonical representation of the maximum margin hyperp. can be found by solving the optimization problem:

$$\max \left( \frac{1}{\|w\|} \right) = \min \frac{1}{2} \|w\|^2$$

Subject to :

$$t_m (\vec{w}^T \vec{x}_m + w_0) \geq 1 \quad \forall m$$

useful since  
the problem  
is in a quadratic  
form that it is  
easy to solve

QP problem solved with Lagrangian method. The solution:

$$w^* = \sum_{m=1}^N a_m t_m \vec{x}_m$$

In this case the number of parameters is equal to N, that is the size of the dataset ( $a_i$  are the Lagrangian multipliers).

KARUSH KUHN TUCKER CONDITION:

$$\text{for each } \vec{x}_m \in X_D, \text{ either } a_m = 0 \text{ or } t_m y(\vec{x}_m) = 1$$

This means that many  $a_m = 0$  and few of them are different from zero.  $\vec{x}_m$  for which  $a_m \neq 0$  do not contribute to the solution!

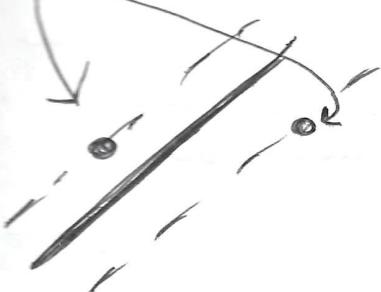
The solution is given only by those points  $x_i$  for which  $\alpha_i \neq 0$ , and these points are very few, sometimes they are just two. (13)

## SUPPORT VECTORS :

Vectors for which the corresponding lagrangian multipliers  $\alpha_i \neq 0$  or in other words they are the point at distance 1 from the margin hyperplane.

maximum

$$SV = \{ \vec{x}_k \in X_D \mid t_k y(\vec{x}_k) = 1 \}$$



SO THE SOLUTION IS !

$$y(\vec{x}) = \underbrace{\sum_{\vec{x}_j \in SV} \alpha_j t_j \vec{x}_j^T \vec{x}}_{\text{hyperplane expressed with } SV} + w_0 = 0$$

We can compute  $w_0$  because we can exploit the property that  $t_k y(\vec{x}_k) = 1$ . Instead of using one particular SV a more stable solution is obtained by averaging over all the support vectors :

$$w_0 = \frac{1}{|SV|} \sum_{\vec{x}_k \in SV} \left( t_k - \sum_{\vec{x}_j \in SV} \alpha_j t_j \vec{x}_j^T \vec{x}_j \right)$$

Note :

the original eq. for  $w_0$  is obtained by multiplying  $y$  for  $t_k$ :

$$t_k \left( \sum_{\vec{x}_j \in SV} \alpha_j \vec{x}_j^T \vec{x}_j + w_0 \right) = 1 \quad \text{and remembering that } t_k^2 = 1$$

REMEMBER THAT YOU HAVE AT LEAST TWO SV. The optimization problem for determining  $\vec{w}$  (dimension  $|x|$ ) has been transformed in an optimization problem for determining  $\vec{\alpha}$  (dimension  $|D|$ , most of them are zero)

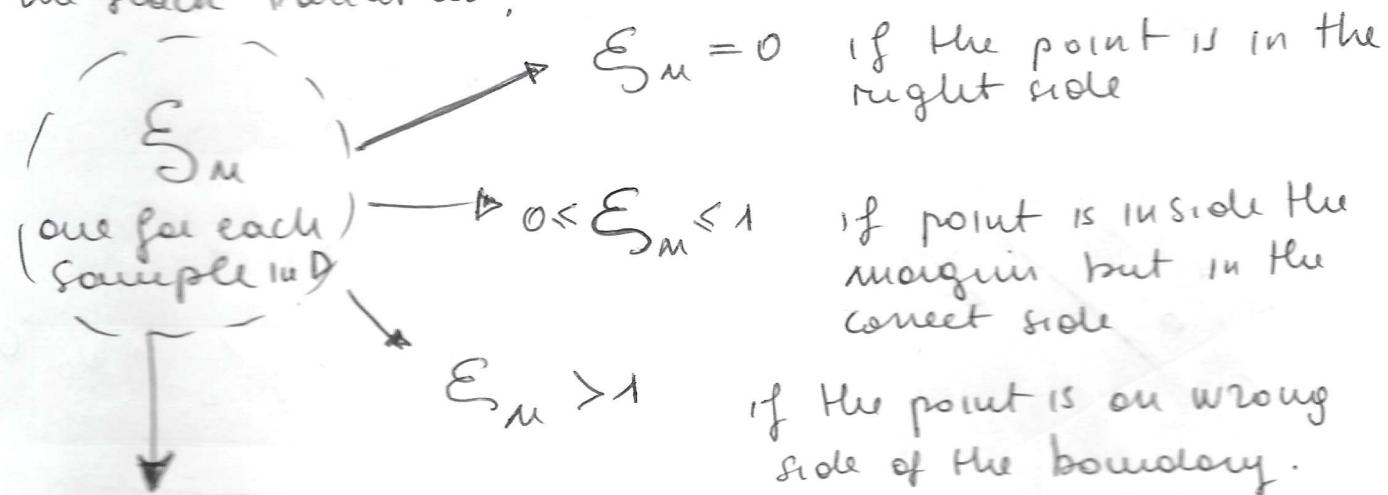
EFFICIENT when  $|x| < |D|$ , very useful when  $|x|$  is large or infinite.

(4)

What if data are "almost" linear separable?

WE CAN  
TRANSFORM THE OPT.  
PROBLEM for SVM  
BY introducing  
the slack variables!

linear separable except for  
some points. Some methods  
do not work, SVM can be  
extended.



"RISK" that you take  
when the sample moves  
away from the safe region

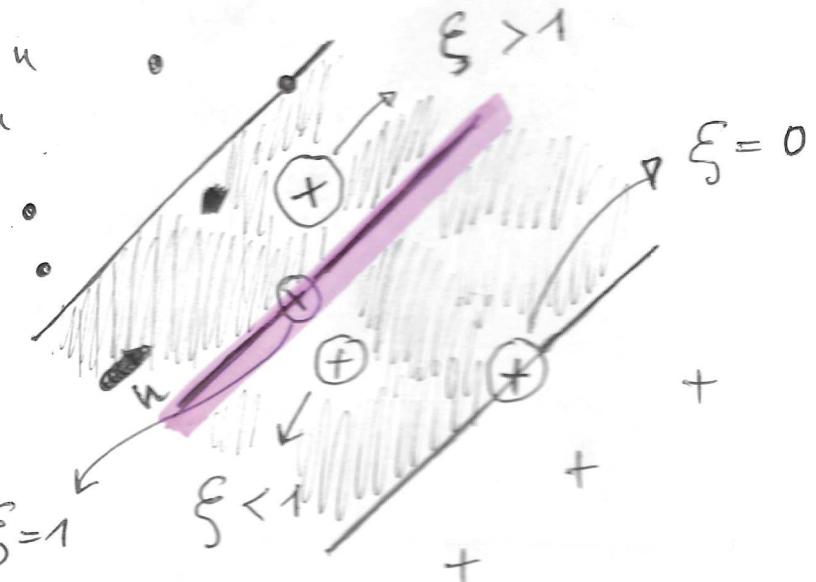
You want to reduce  
this variables (to min=  
size). But...

I can accept that  
some few points  
move away from  
the region by paying  
a cost (expressed by  
MARGIN CONSTRAINT)  
WRONG REGION;

WITH SOFT MARGIN CONSTRAINT I CAN ACCEPT POINTS IN THE

$$\boxed{t_m y(\vec{x}_u) \geq 1 - \xi_u} \quad \forall u=1\dots N$$

IS A RELAXATION.



The new optimization problem is given by the terms!

(15)

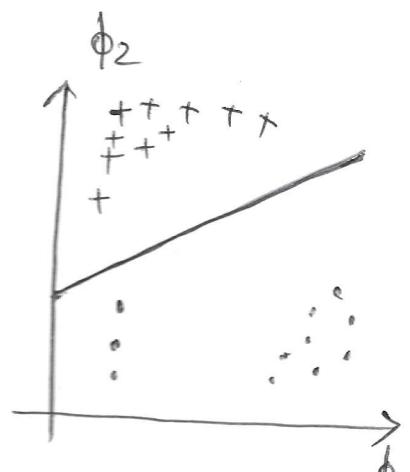
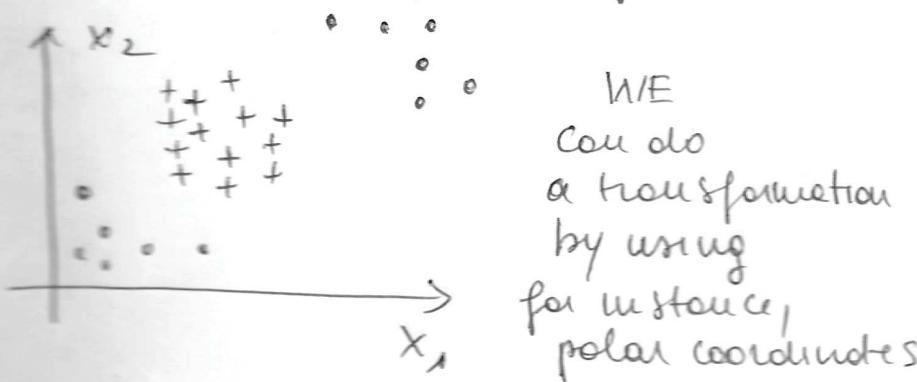
$$\boxed{\min \frac{1}{2} \|\vec{w}\|^2 + C \sum_{m=1}^N \xi_m}$$

s.t.  $t_m \gamma(\vec{x}_m) \geq 1 - \xi_m \quad \forall m = 1 \dots N$   
 $\xi_m \geq 0 \quad \forall m = 1 \dots N$

The solution is similar to the previous one, where most of  $a_m$  are 0; ~~all these methods are robust to outliers~~, all these methods are robust to outliers, that are all points for which  $a_m = 0$ .

### Basis function

Let's consider the following situation!



So far we considered models working directly on  $\vec{x}$ . All the results hold if we consider a non linear transformation of the input space,  $\phi(\vec{x})$  (BASIS FUNCTION).

Decision boundaries will be linear in the feature space and non-linear in the original space  $\vec{x}$ .

CLASSES THAT ARE UNARILY SEPARABLE IN THE FEATURE SPACE  $\phi$  MAY NOT BE SEPARABLE IN THE INPUT SPACE.

We can try different basis functions (kernel) to see which is the best.