

EX 1

1.1

confusion matrix reports how many times a class  $C_i$  is classified  
in class  $C_j$

		$C_1$	$C_2$	$\leftarrow$ predicted class
true classes	$\leftarrow$	$C_1$		
		$C_2$		

in the main diagonal we have the

accuracy = # element of diagonal / # num of all elements

the number of elements in a row must be equal numbers of samples of  
for class  $C_j$

1.2

1.3

$\nabla^P$	A	B	C
T	140	50	10
A	70%	25%	5%
B	10%	80%	10%
C	5%	15%	80%

accuracy =  $480/600 = 0.8 = 80\%$ .

Ex. 2

2.1

①      ②      ③      ④

$$D = \{(FFF, F), (FTT, T), (TTF, T), (TFT, T)\}$$

$$h_1 = (x_1 \wedge \neg x_2 \wedge x_3) \vee x_2$$

$$h_2 = (\neg x_1 \wedge x_2 \wedge x_3) \vee x_1$$

an hypothesis is consistent if and only if given a dataset  $D$  an hypothesis  $h$  and a target concept  $c$   $h(x) = c(x)$  for each training sample  $(x, c(x))$

$$\text{consistent}(h, D) = (\forall x \in D) h(x) = c(x)$$

• check with  $h_1$  ~~and  $h_2$~~

$x_1$	$x_2$	$x_3$	$h_1$
F	F	F	F
F	T	T	T
T	T	F	T
T	F	T	T

• check with  $h_2$

$x_1$	$x_2$	$x_3$	$h_2$
F	F	F	F
F	T	T	T
T	T	F	T
T	F	T	T

both  $h_1$  and  $h_2$  are consistent

(5)

2.2

$$P(D|h_1) = 0.6 \quad P(D|h_2) = 0.8 \quad P(h_1) = 0.2 \quad P(h_2) = 0.1$$

posterior hypothesis  $P(h_1|D) = P(D|h_1)P(h_1) = 0.6 \cdot 0.2 = 0.12$   
 $P(h_2|D) = P(D|h_2)P(h_2) = 0.8 \cdot 0.1 = 0.08$

$$\text{argmax } \{ (0.12), (0.08) \} \rightarrow h_1 = \text{hmap}$$

EX. 3

3.1

using a non linear functions of input variables

$$y(x, w) = \sum_{n=1}^M w_n \phi_n(x) \quad \text{with } w = [w_1, \dots, w_M]^T$$

and  $\phi = [\phi_0(x), \dots, \phi_M(x)]^T$  non linear in  $w$

number of trainable parameters =  $M+1$

↳ Bias

3.2

target value is effected by additive noise  $t = y(x, w) + \epsilon$

assuming gaussian noise  $P(\epsilon|\beta) = N(\epsilon|0, \beta^{-1})$  we have

$$P(t|w, x, \beta) = N(t|y(x, w), \beta^{-1})$$

assuming distribution iid

$$P(t_1, \dots, t_N | x_1, \dots, x_N, w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1}) = \\ = -\frac{\beta}{2} \sum_{n=1}^N [t_n - w^T \phi(x_n)]^2 - \frac{N}{2} \ln(2\pi\beta^{-1})$$

maximum likelihood correspond

$$\text{argmax}_w P(t_1, \dots, t_N | x_1, \dots, x_N, w, \beta)$$

correspond least squares error

$$\text{argmin}_w E_w(w) = \text{argmin}_w \frac{1}{2} \sum_{n=1}^N [t_n - w^T \phi(x_n)]^2$$

a way to solve it is using a numerical algorithm

$$\hat{w} \leftarrow \hat{w} - \eta \nabla E_w \Rightarrow \hat{w} \leftarrow \hat{w} - \eta [t_n - w^T \phi(x_n)] \phi(x_n)$$

(6)

EX. 4

4.1

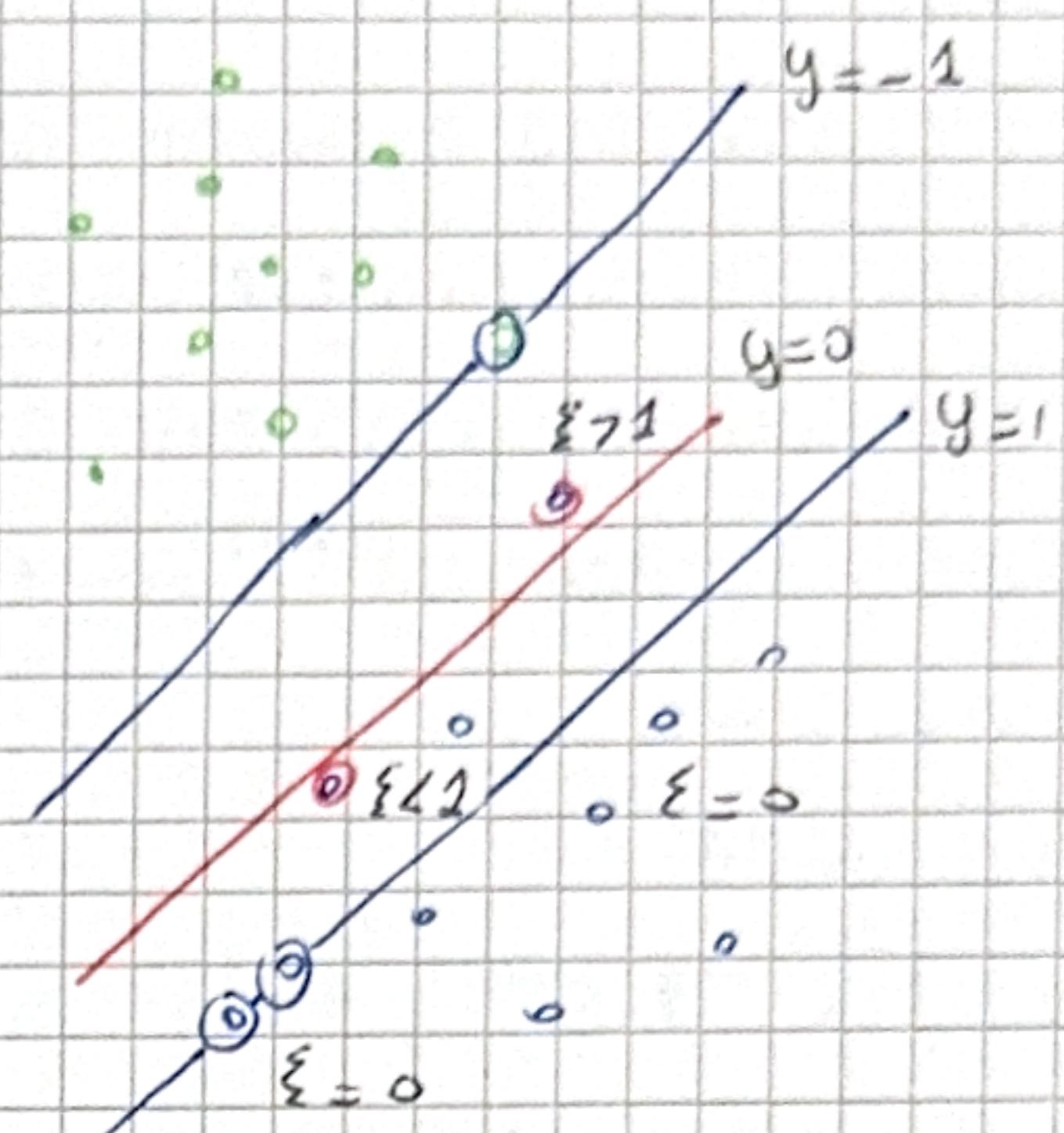
for have sum with soft margin constraints we introduce

a slack variable  $\xi_n \geq 0$  with  $n=1, \dots, N$

- $\xi_n = 0$  if point is ON or INSIDE the correct margin boundary

- $0 < \xi_n \leq 1$  if point is inside the margin but comet side

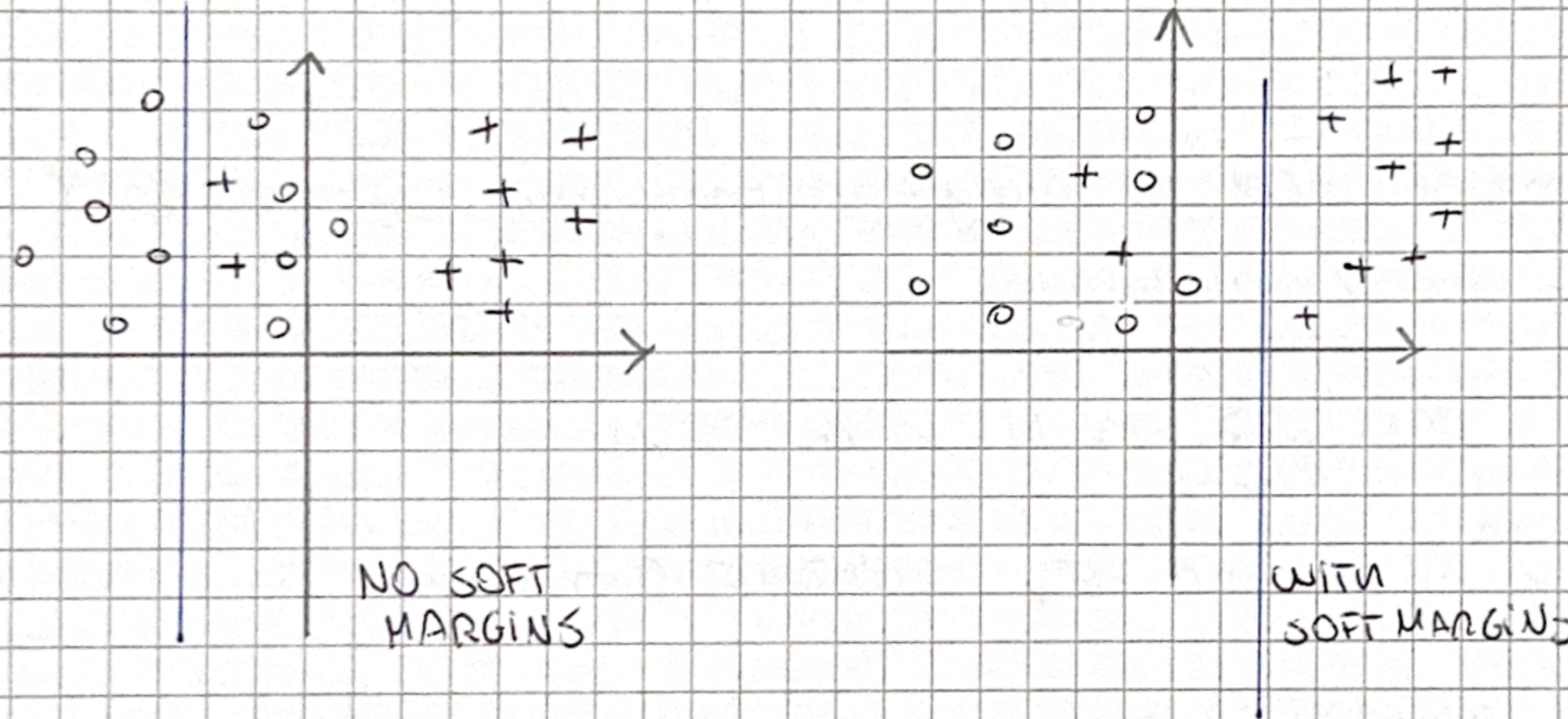
- $\xi_n > 1$  point in wrong side of boundary



optimizations with soft constraints

$$\omega^*, \omega^0 = \text{argmin} \frac{1}{2} \|\omega\|^2 + C \sum_{n=1}^N \xi_n$$

4.2



EX. 5

5.1

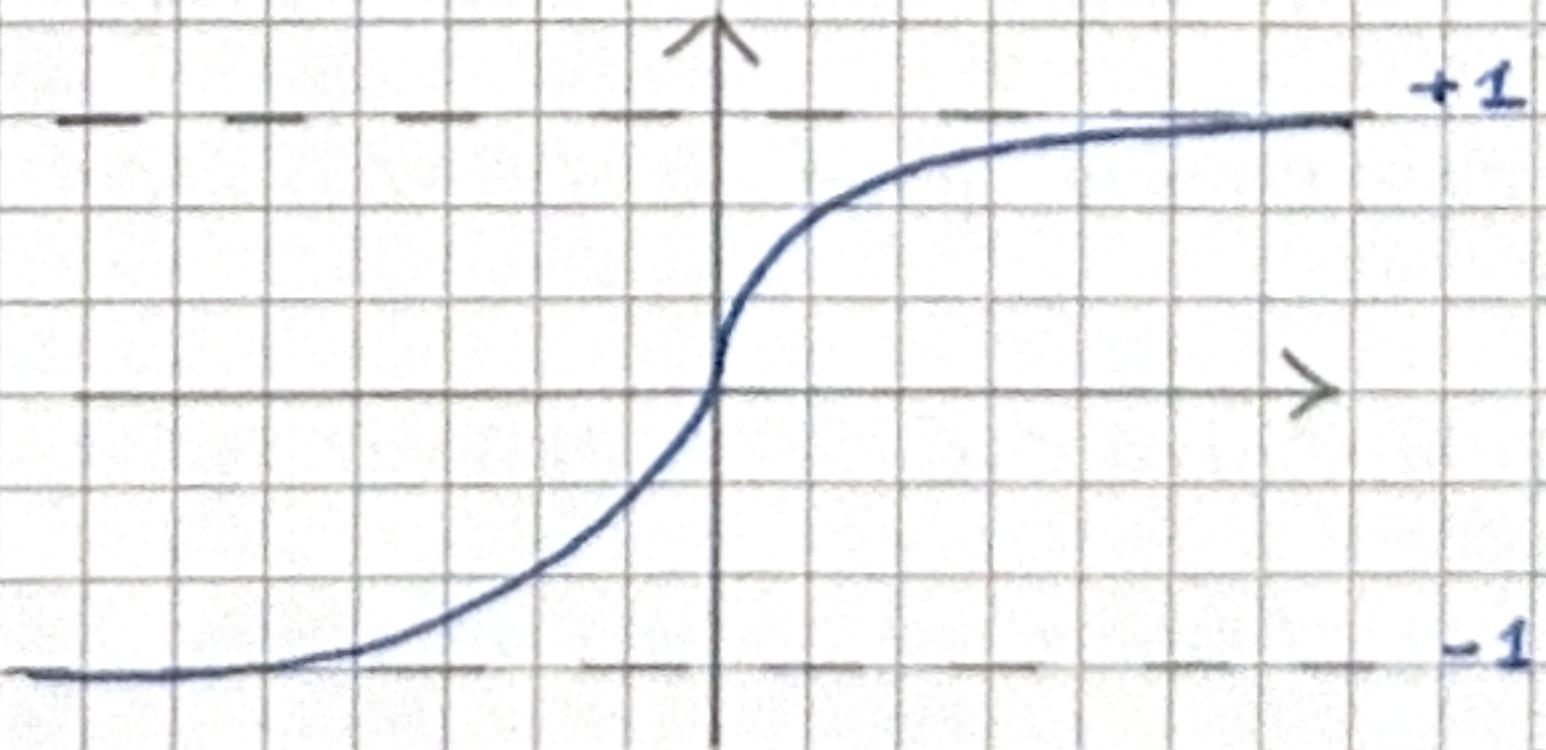
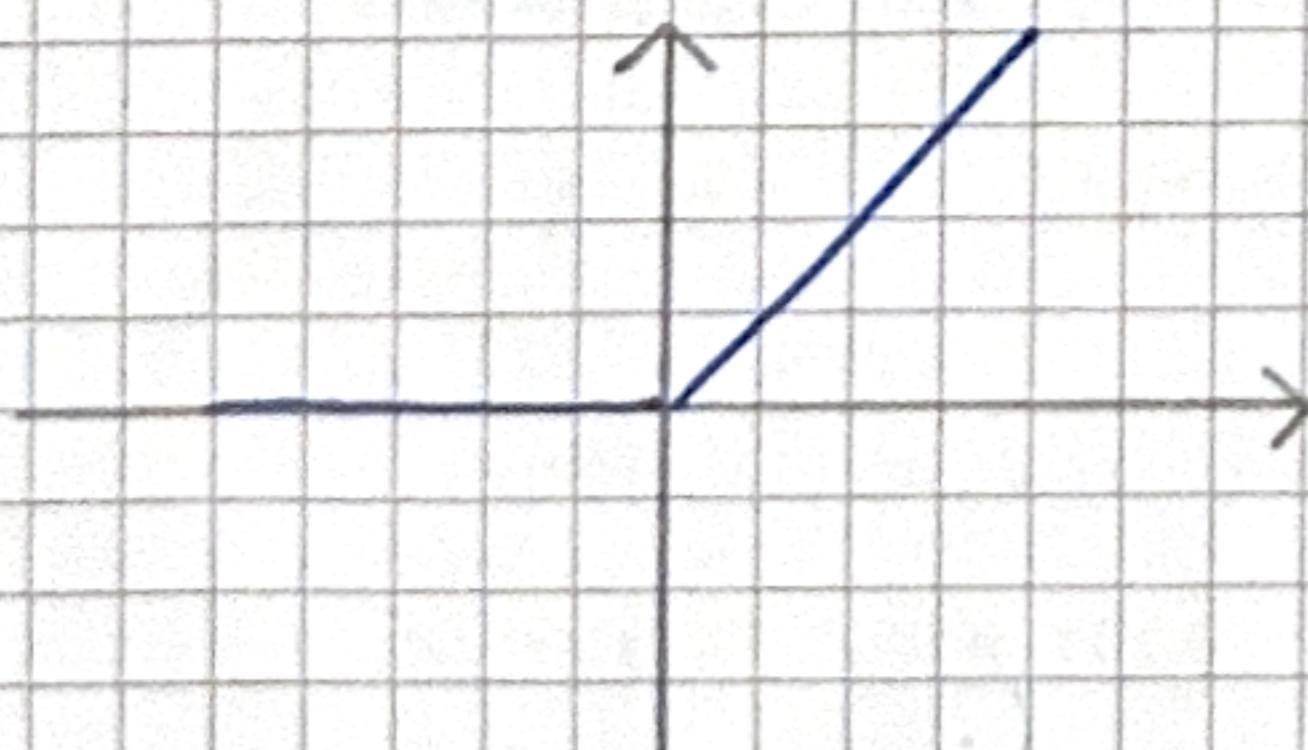
as a regression task, for the activation function of the output layer

I will use IDENTITY ACTIVATION FUNCTION  $y = w^T h + b$

and an loss function MSE mean square error  $MSE = \frac{1}{N} \sum_{n=1}^N (t_n - x_n)^2$

5.2

we can use an activation function ReLU ( $y(\alpha) = \max(0, \alpha)$ ) or sigmoid and hyperbolic tangent  $y(\alpha) = \sigma(\alpha)$  or  $y(\alpha) = \tanh(\alpha)$



(7)

5.3

negative learning rate  $\eta > 0$

negative initial values of  $\theta^{(k)}$

$k \leftarrow 1$

while stopping criterion not met do:

sample a random subset (minibatch)  $\{x^{(1)}, \dots, x^{(m)}\}$

compute the gradient estimate  $\hat{g} = \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}, \theta^{(k)}), t^{(i)})$

apply update:  $\theta^{(k+1)} \leftarrow \theta^{(k)} - \eta \hat{g}$

$k \leftarrow k + 1$

end while

$\nabla_{\theta} L(f(x^{(i)}, \theta^{(k)}), t^{(i)})$  is obtain with BACKPROPAGATION

hyperparameter  $(m)$  → number of batch

$\eta$  → learning rate

if  $\eta$  too small could not converge  $M$  too large could diverge

EX. 6

6.1

K-MEANS is an unsupervised learning algorithm that try to compute  $\pi_k$ ,  $\mu_k$  and  $\Sigma_k$  for compute

$$p(x) = \sum_{k=1}^K \pi_k N(x; \mu_k, \Sigma_k)$$

where  $\pi_k$  is prior probability  $\mu_k$  the mean and  $\Sigma_k$  is covariance matrix

each instance  $x_n$  is yemaining by

1. choosing  $K$  gaussian according to  $\pi_1, \dots, \pi_K$
2. generating an instances at random using according to that gaussian  $(\mu_k, \Sigma_k)$

PSEUDOCODE of KMEANS

1. begin a decision on the value  $K =$  number of cluster
2. put any portion of data into  $K$  cluster that can be done systematically or randomly
  - I. select  $K$  sample as centroid of cluster
  - II. assign the remain  $N-K$  samples to the nearest centroid and recompute the centroid
3. take each sample in sequence and compute its distance from centroid of each cluster. if same samples in more than

⑧

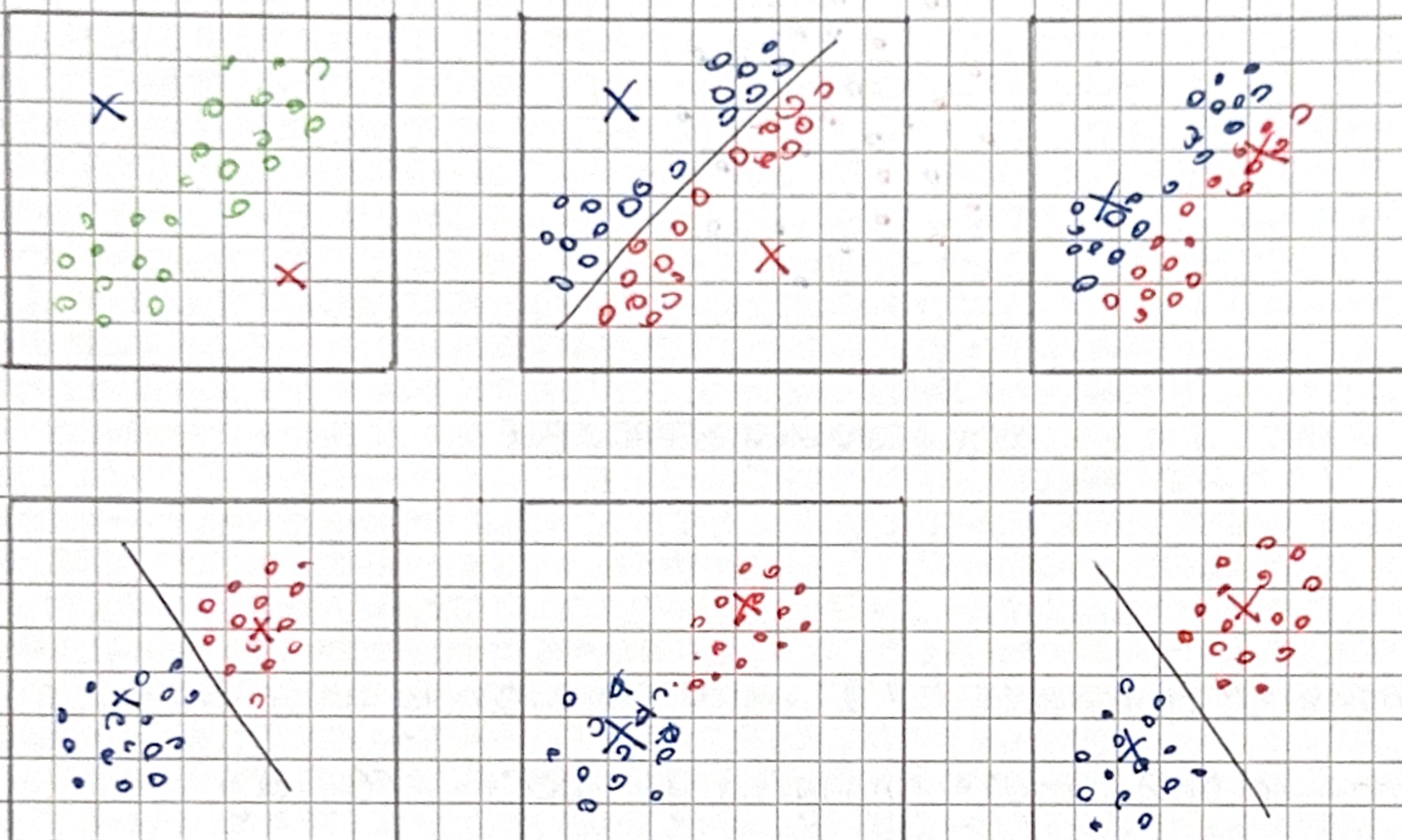
the meanest centroid acquire and recompute the centroid

#### 4. repeat step 3 until convergence

termination: each iteration steps the sum of distance from each training sample is decrease (from centroid)

there are only finitely partitions of training example into k clusters

6.2



ML EXAMS 22/06/2022

EX.1

1.1

we have overfitting when given a dataset  $D$  exist an hypothesis  $h'$  that  $\text{error}_S(h) < \text{error}(h')$  and  $\text{error}_D(h) > \text{error}_D(h')$

where  $S \subseteq D$  and  $h, h' \in \text{Hypothesis space}$

1.2

we are sure that we have overfitting with decision tree when for same problem we have two different trees one more deeper than other

so consider two decision tree  $T$  and  $T'$  obtain with ID3 algorithm

there are two methods for reduce overfitting in decision tree

- REDUCE ERROR PRUNING we must split training data in validation and we evaluate impact on validation set of pruning each terminal node. greedily remove the one that most improves validation set accuracy