

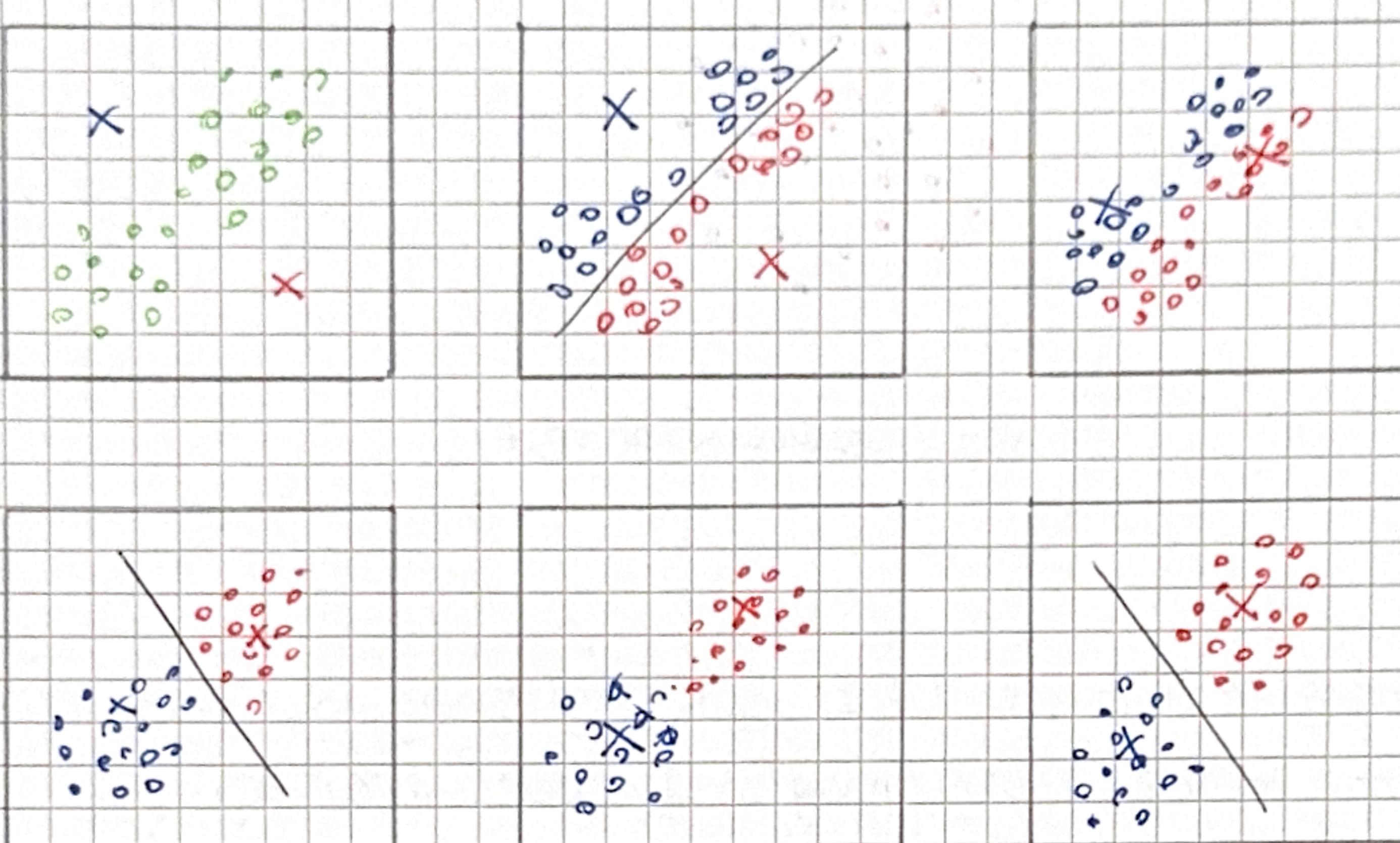
the meanest centroid moves and recompute the centroids

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 $\text{error}_S(h) < \text{error}_D(h^*)$ 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 solution 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 possible node. greedily remove the one that most improves validation set accuracy

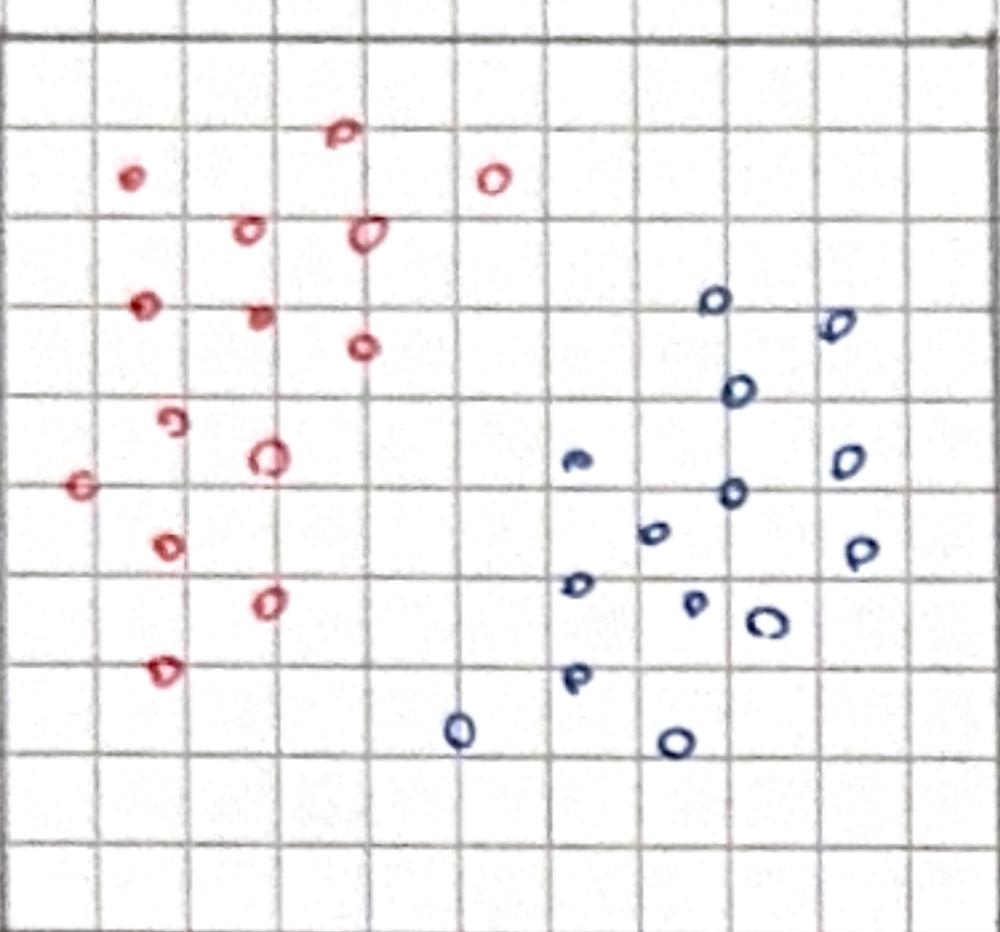
- RULE POST PRUNING (C4.5): if accuracy after pruning is greater than before we have no overfitting

EX. 2

2.1

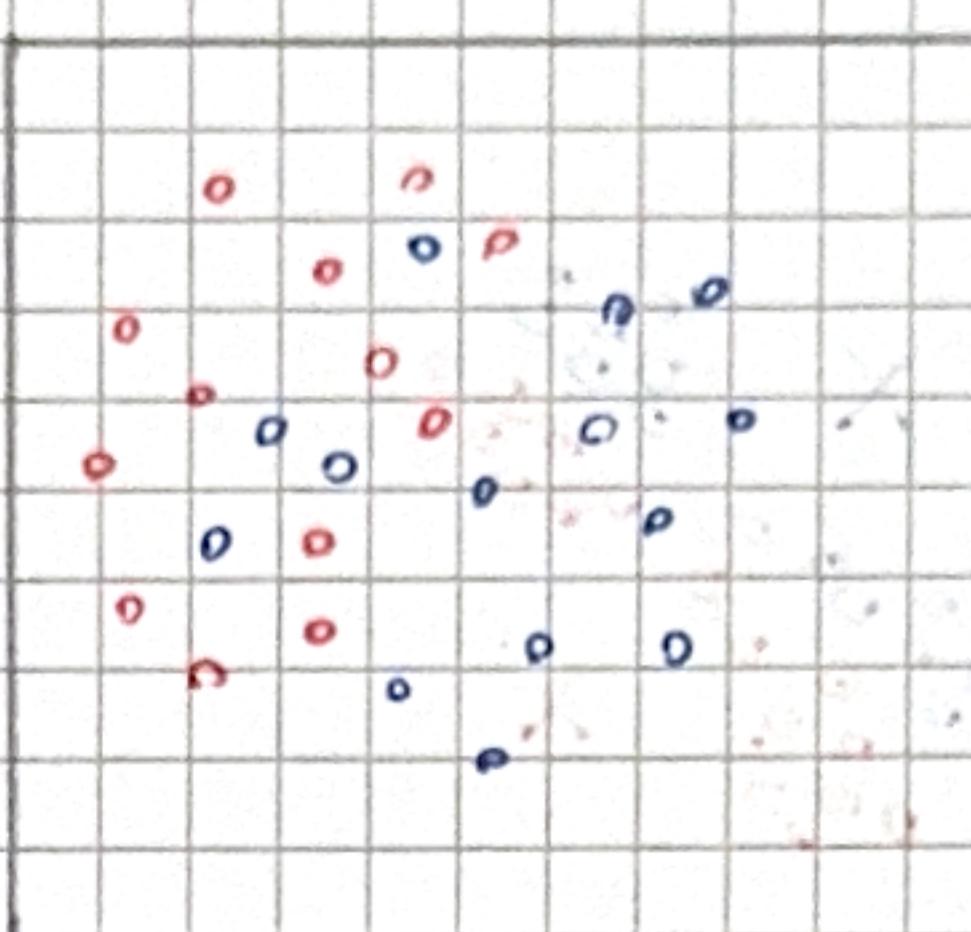
data is linearly separable if exist a surface that split data in two regions such that differently classified instances are separated

2.2



LINEARLY SEPARABLE

2.3



NOT LINEARLY SEPARABLE

2.4

in case of linearly separable data one (2.2) sum can easily solve this problem because aims to maximum margin for better accuracy a margin is defined as $\frac{|g(x_n)|}{\|w\|}$ so the goal is

$$\underset{w}{\operatorname{argmax}} \left(\min_{n} \frac{|g(x_n)|}{\|w\|} \right) = \underset{w}{\operatorname{argmax}} \frac{1}{\|w\|} \min_{n} |g(x_n)| =$$

$$= \underset{w}{\operatorname{argmax}} \frac{1}{\|w\|} \min_{n} [t_n + \frac{1}{\|w\|} [w^T x_n + w_0]]$$

so an optimization problem can be solved using lagrangian multipliers and if also we averaging for all $SV = \{x_k \in D \mid t_k(y_k(x_k)) = 1\}$ we obtain

~~argmax~~

$$w_0^* = \frac{1}{|SV|} \sum_{x_k \in SV} t_k \left(- \sum_{j \neq k} a_j^* t_j x_k^T x_j \right)$$

$$y(x, w) = \operatorname{sign} \left(\sum_{x_k \in SV} a_k^* t_k x_k^T x + w_0^* \right)$$

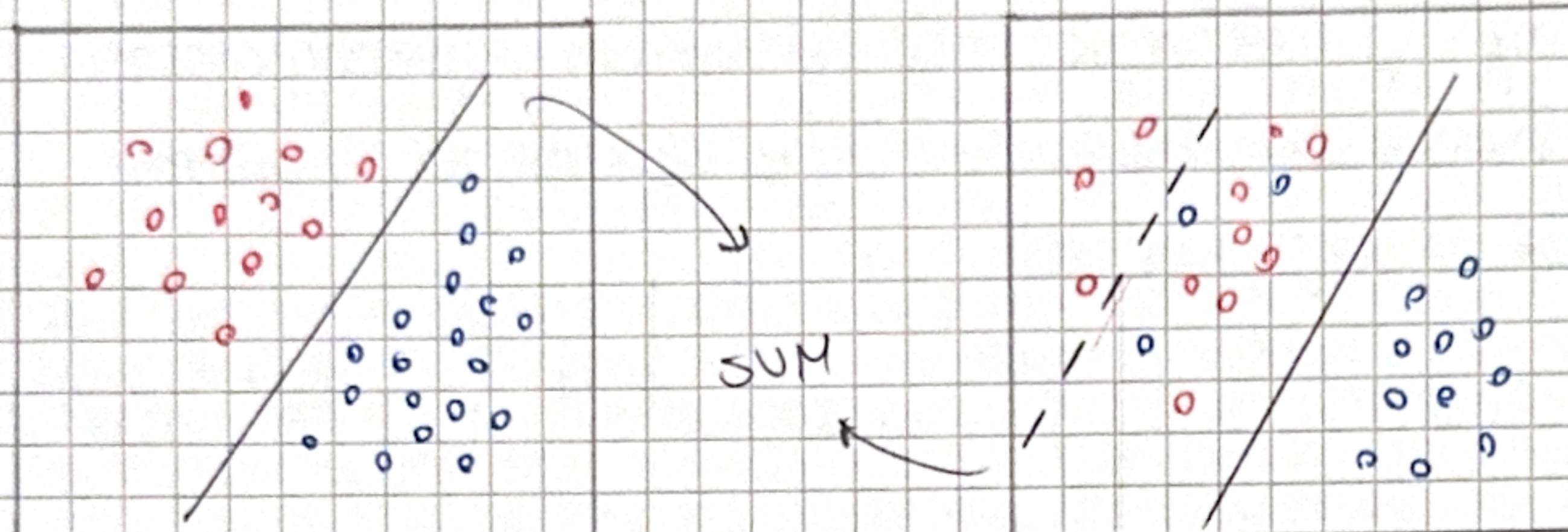
for the second case due the fact the data is not linearly separable we need to use soft margins constraint in detail

$E = 0$ if point is on or inside the correct margin boundary

$0 < \xi_n \leq 1$ if point is inside the margin but in correct boundary
 $\xi_n > 1$ if point is on wrong side of boundary

In this case the problem can be solved:

$$w^*, w_0^* = \text{argmin}_w \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n$$



\hookrightarrow SVM with soft margins constraint

EX.3

3.1

Given a dataset $D = \{(x_n, t_n)\}_{n=1}^N$ we can define a linear model for noise regression problem in this way:

$$y(x, w) = \sum_{i=1}^M w_i \phi(x_i) \text{ still linear in } w$$

The number of trainable parameters is $M+1$ (1 in bias term)

3.2

The goal of linear regression target value ~~can~~ affected by additive noise E can be written

$$y(t) = y(x, w) + E$$

If we assume Gaussian noise $P(E|\beta) = N(\epsilon|0, \beta^{-1})$ we have

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

If we assume iid observation we have

$$\begin{aligned} P(t_1, \dots, t_N | x_1, \dots, x_N, w, \beta) &= \prod_{i=1}^N P(t_i | x_i, w, \beta) = \\ &= -\beta \cdot \frac{1}{2} \sum_{n=1}^N [t_n - w^\top \phi(x_n)]^2 - \frac{N}{2} \ln(2\pi\beta^{-1}) \end{aligned}$$

maximum likelihood

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

correspond least squares

$$\underset{w}{\text{argmin}} \quad E_0(w) = \underset{w}{\text{argmin}} \quad \frac{1}{2} \sum_{n=1}^N [t_n - w^\top \phi(x_n)]^2$$

EX. 4

4.1

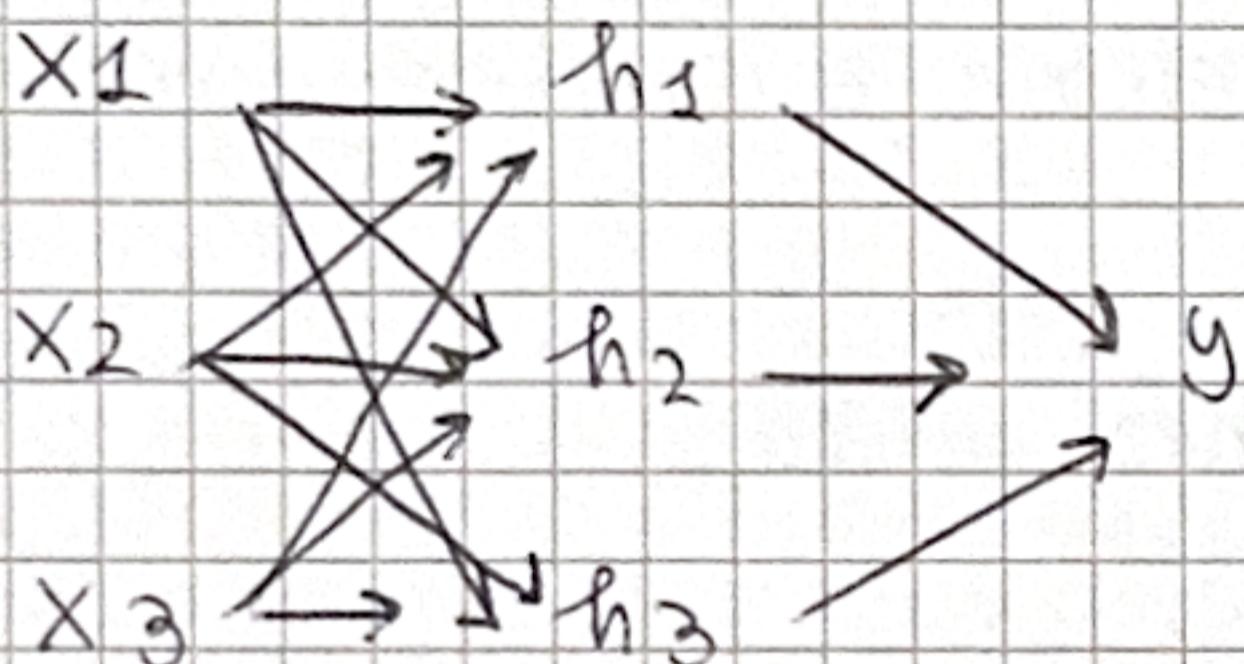
feed forward network \rightarrow information flow from input to output without any loop and network is composition of elementary functions in an acyclic graph.

In our case considering $f: \mathbb{R}^3 \rightarrow \mathbb{R}$ the parametric model can be written ~~flexible and expressive where each node approximates a function and θ are parameters of function~~

$$f(x, \theta) = f^{(L)}(f^{(L-1)}(\dots f^{(2)}(x_1, x_2, x_3, \theta^{(1)}); \dots; \theta^{(L-1)}); \theta^{(L)})$$

where L are the layers of network, θ are parameters and x_i are the input

considering an example a two layers of networks we will have



~~weights~~

where h_i are hidden units

$h = g(w^T x + b)$ with g the activation function

$$\text{output } y = w^T h + b$$

4.2

a suitable loss depends on task we are trying to solve in this case due the fact $f: \mathbb{R}^3 \rightarrow \mathbb{R}$ I suppose in a regression task so I will use

RMSE loss function $\frac{1}{N} \sum_{n=1}^N (t_n - x_n)^2$

4.3

I will use stochastic gradient descent:

Require learning rate $\eta > 0$

Require initial values of $\theta^{(1)}$

$$k \leftarrow 1$$

while stopping criterion not met do

sample a subset (mini-batch) $\{x^{(1)}, \dots, x^{(m)}\}$ from dataset

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

$$\text{update } \theta^{(k+1)} \leftarrow \theta^{(k)} - \eta \nabla g$$

$$k \leftarrow k + 1$$

end while

(1) $L(f(x, \theta^{(k)}), t^{(i)})$ is compute with BACK PROPAGATION

(2)

EX.5

S.1

$$MDP = \langle X, A, f, r \rangle$$

$X = \text{room}$

$f: X \times A \rightarrow \mathbb{R}$ is transition function

$A = \text{clean and move}$

$r: \text{reward function}$

S.2

considering $N=3$ we can derive the optimal policy

$N=1$

$N=2$

$N=3$

$\pi \rightarrow$ policy have to maximize the cumulative reward in particular

if $s \in \text{clean} \Rightarrow r(s, a) = +1$ for clean room

$r(s, a) = -2$ move to room

$r(s, a) = 0$ otherwise

$\pi = \{\langle (N=1, \text{clean}), \text{move to } N_2 \rangle, \langle (N_1, \text{dirty}), \text{clean} \rangle, \langle (N_2, \text{dirty}), \text{clean} \rangle$

$\langle (N=2, \text{clean}), \text{move to } N_3 \rangle, \langle (N=3, \text{dirty}), \text{clean} \rangle, \langle (N=3, \text{clean}), \text{stop} \rangle\}$

initial state $= N=1$

EX.6

6.1

$$p(x) = \sum_{n=1}^3 \pi_n N(x, \mu_n, \Sigma_n)$$

x : input data μ_n : mean Σ : covariance matrix

weight of component ~~size~~ π_k $k=3$ but $\sum_k \pi_k = 1 \Rightarrow k-1 = 2$

μ_n is a vector in two dimensions so ~~no~~ #params of $\mu_n = 6$

Σ_k is a symmetric covariance matrix $d=2$ no #params $\frac{d(d+1)}{2} = 3$

