

3. Use in an ~~optimization~~ if the hypothesis space is not too large and it have an analytical solutions, otherwise is not practical ~~and better~~  
~~use~~ use Naive Bayes classifier

#### Ex. 6

once we augment more in known the evolution of a dynamic system doesn't depend on previous more actions and observations. the current state contains all informations needed to predict the future future state are conditionally independent of past states and past observations statistically independent

MDP can be describe  $\langle X, A, f, r \rangle$   $X$ : finit set of states

$A$ : finite set of actions  $f: X \times A \rightarrow X$  transition function and  $r$ : reward function

MM can be describe  $\langle X, Z, P \rangle$  where  $X$  is set of observations and  $P$  is initial distribution

MDP has the property of fully observability

MDP

$$x_t \rightarrow x_{t+1}$$

↓  
at

MM

$$x_0 \longrightarrow x_{t+1} \longrightarrow \dots \longrightarrow x_n$$

↓                  ↓  
z\_0            z\_{t+1}

ML EXAMS 15/10/2019

#### Ex. 1

1. how many times an older person of class C in classifier are in the class ej

2.

	A	B	C
A	40	50	10
B	10	80	10
C	5	5	80

	A	B	C	← medicated
A	40%	50%	10%	
B	10%	80%	10%	
C	5%	5%	80%	

↳ true values

3. the accuracy in a confusion matrix is the sum of elements in the diagonal ~~number~~ divide by the all of elements of matrix

$$\text{accuracy} = (40 + 80 + 80) / 300 = 0.7$$

EX. 2

1. for noise free task we will use linear regression

$$\text{we know that } y(x, w) = \sum_{j=1}^N w_j \phi_j(x)$$

that is a linear model in parameter  $w$ . we know that  $t = y(x, w) + \epsilon$  with  $\epsilon$  additive noise (Gaussian)

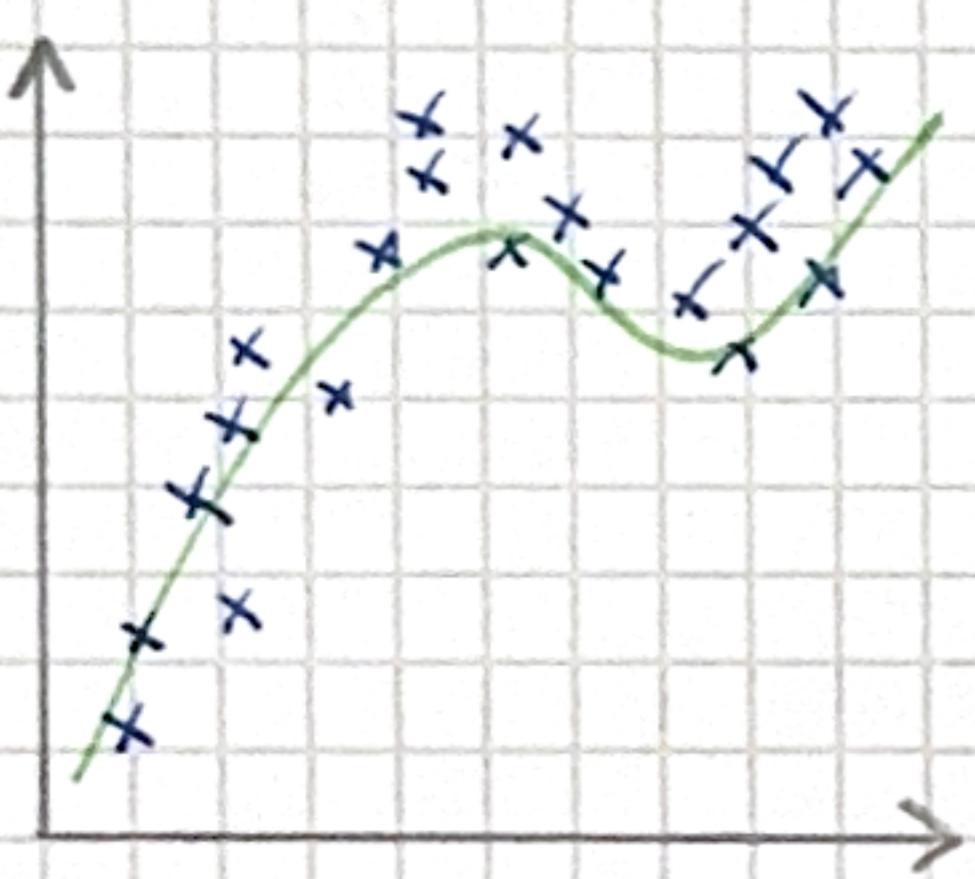
$$\Rightarrow p(\epsilon | \beta) = N(\epsilon | 0, \beta^{-1}) \Rightarrow p(y | x_i, x_m, w, \beta) = N(y | y(x, w), \beta^{-1})$$

if we use iid hypothesis  $p(y_1, \dots, y_n | x_1, \dots, x_n, w, \beta) =$

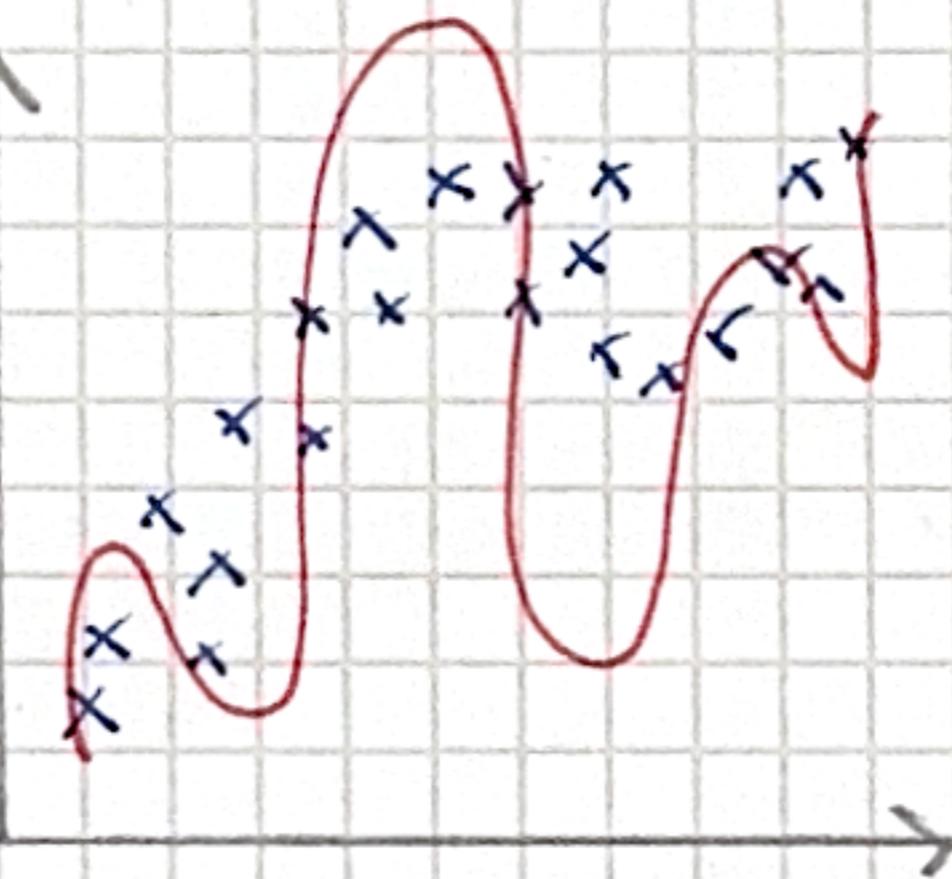
$$= \prod_{n=1}^N N(y_n | w^\top \phi(x_n), \beta^{-1}) = -\frac{\beta}{2} \sum_{n=1}^N (y_n - w^\top \phi(x_n))^2 - \frac{N}{2} \ln(\beta^{-1})$$

the error function is  $E_D(w) = \frac{1}{2} \sum_{n=1}^N (y_n - w^\top \phi(x_n))^2$  and we want find argmin  $E_D(w)$

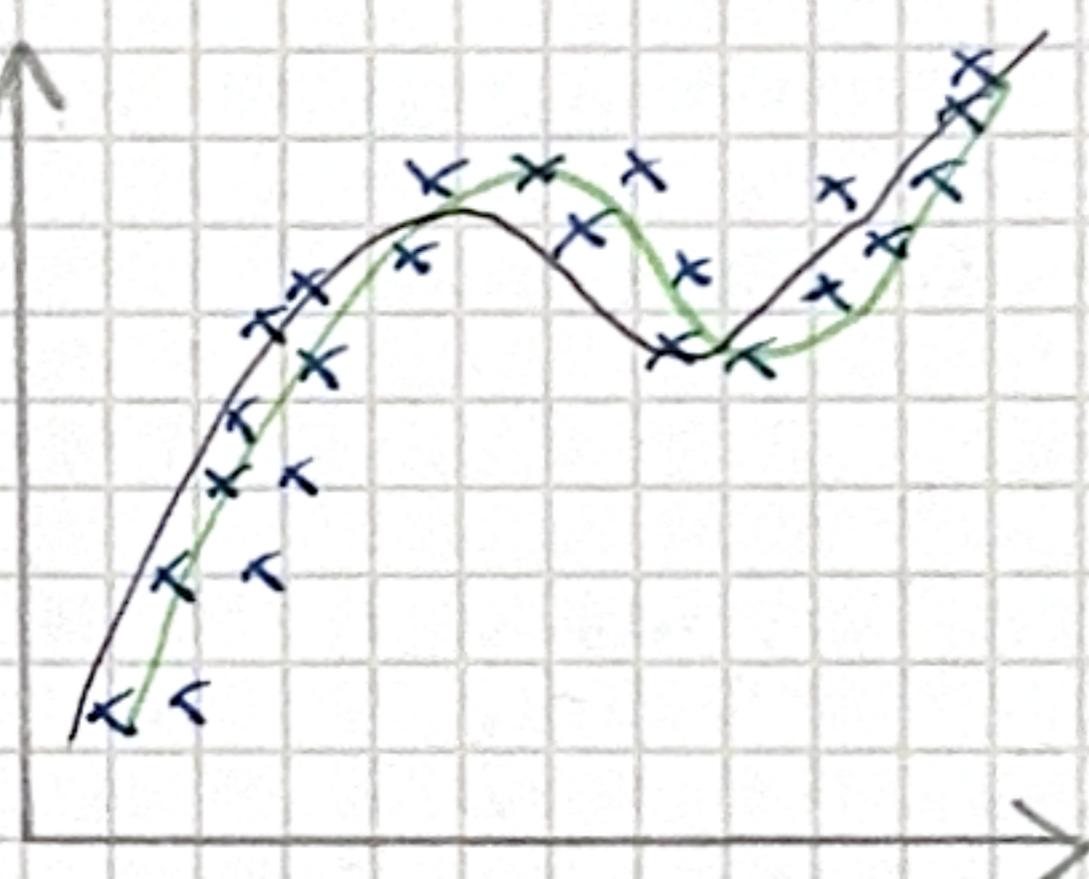
2. we can apply regularization  $w^* = \text{argmin } E_D(w) + \lambda E_w(w)$  with  $\lambda > 0$  and  $E_w(w) = \frac{1}{2} w^\top w$



true target function

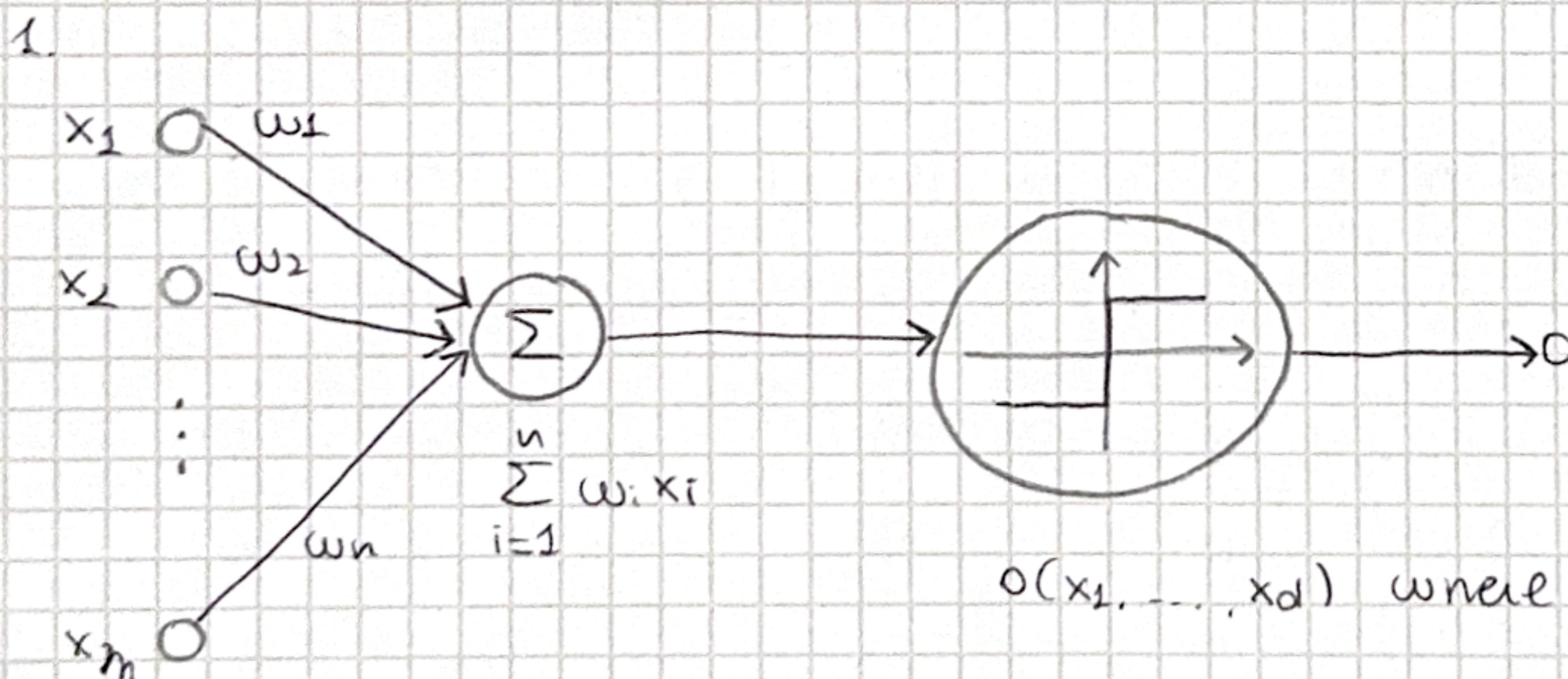


OVERFITTING



regularization

EX. 3



$$o(x_1, \dots, x_d) \text{ where } o(x_1, \dots, x_d) = \begin{cases} 1 & \text{if } \sum_{i=1}^d w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

We want minimize the function:

$$E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - o(n))^2 = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \mathbf{w})^2$$

$$\frac{\partial E(w)}{\partial w} = \sum_{n=1}^N (t_n - w^T x_n) (-x_n)_m = -\sum_{n=1}^N (t_n - \text{sign}(w^T x_n)) (x_n)_m$$

2) we want to find  $w^*$  numerically:

$$\hat{w} \leftarrow \hat{w} + \eta \Delta w_i \quad \text{when } \Delta w_i = -M \sum_{n=1}^N (t_n - \text{sign}(w^T x_n)) (x_n)_m$$

$\Delta w_i$  can be update with

batch minimization or sequential mode

3) if we use an large  $\eta$  is possible that ~~the~~ will diverge otherwise  
if  $\eta$  is too small the algorithm will never converge  
but the algorithm is robust to local ~~non~~ minima

EX. 4

1. SVM aims at maximum margin providing for better accuracy

the margin is  $|y(x_n)| / \|w\|$  (minimum distance between  $x_k$  and  $h$ )

given a dataset  $D$  and hyperplane  $h$  the margin is computed

as

$$\min_{n=1, \dots, m} \frac{|y(x_n)|}{\|w\|} = \dots = \frac{1}{\|w\|} \min [t_n (w^T x_n + w_0)]$$

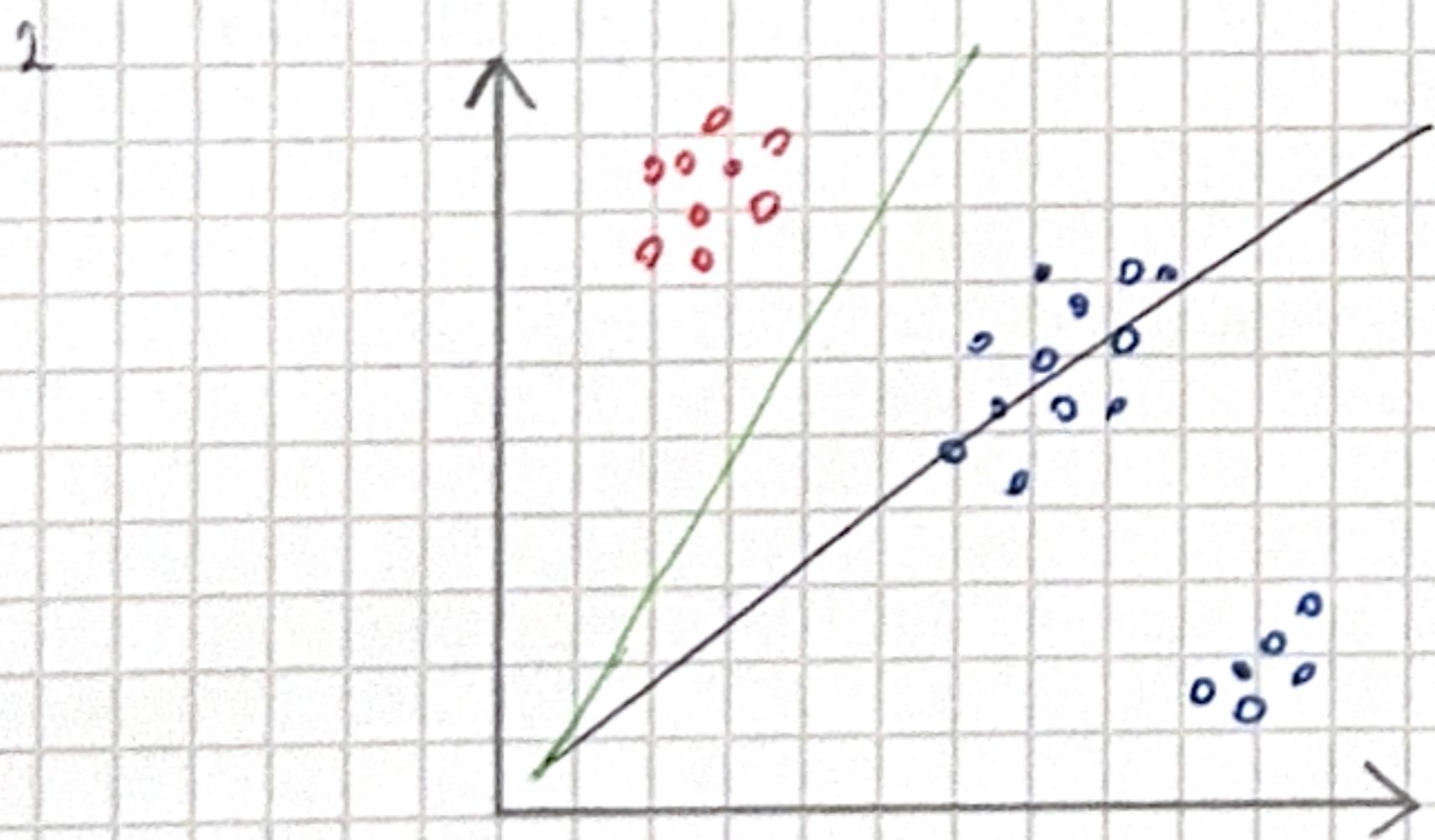
so the maximum margin is computed as:

$$w^*, w_0^* = \underset{w, w_0}{\operatorname{argmax}} \frac{1}{\|w\|} \min [t_n (w^T x_n + w_0)]$$

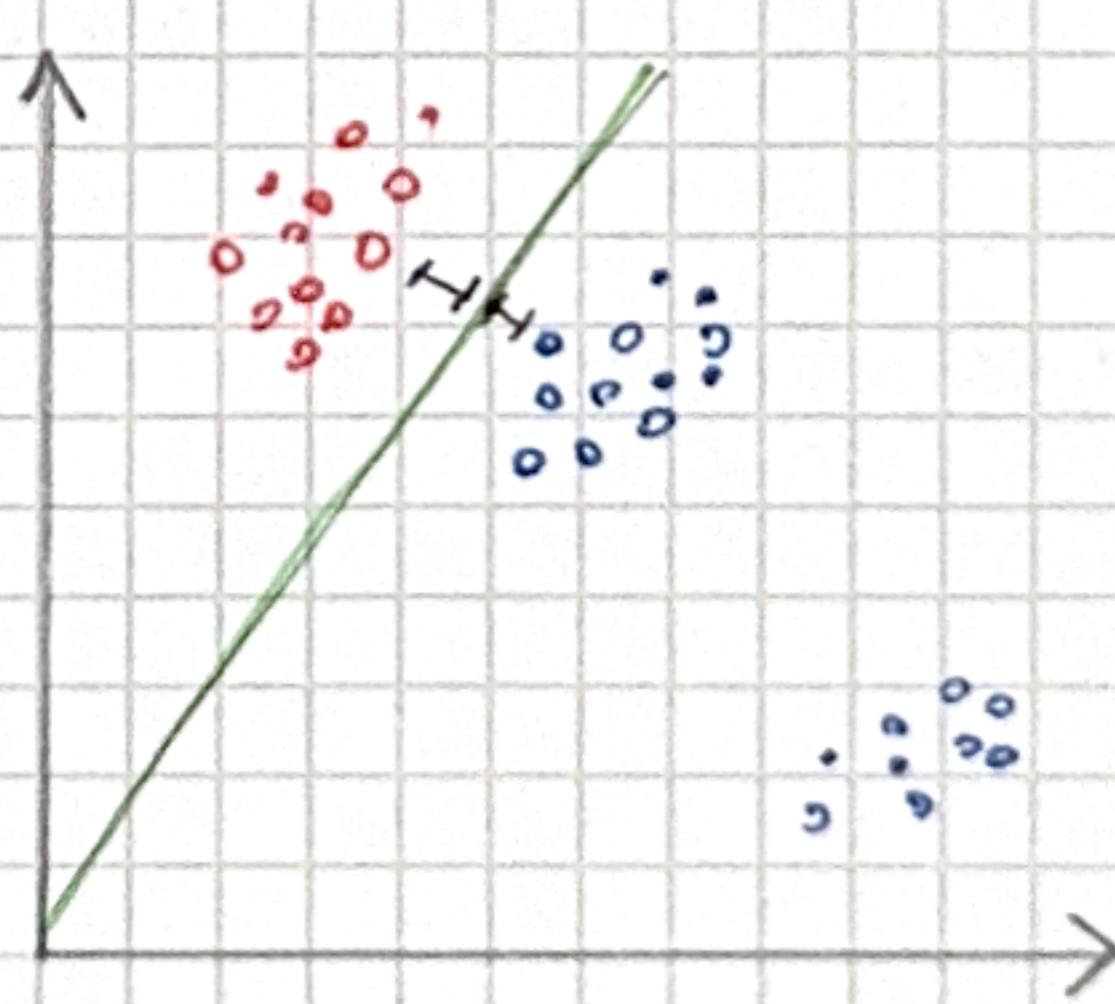
a more stable version is computed by averaging over all support vectors

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

where  $\alpha_j^*$  are lagrangian optimizers results of optimization problem



least squares



SVM  
also robust to outliers

### EX. 5

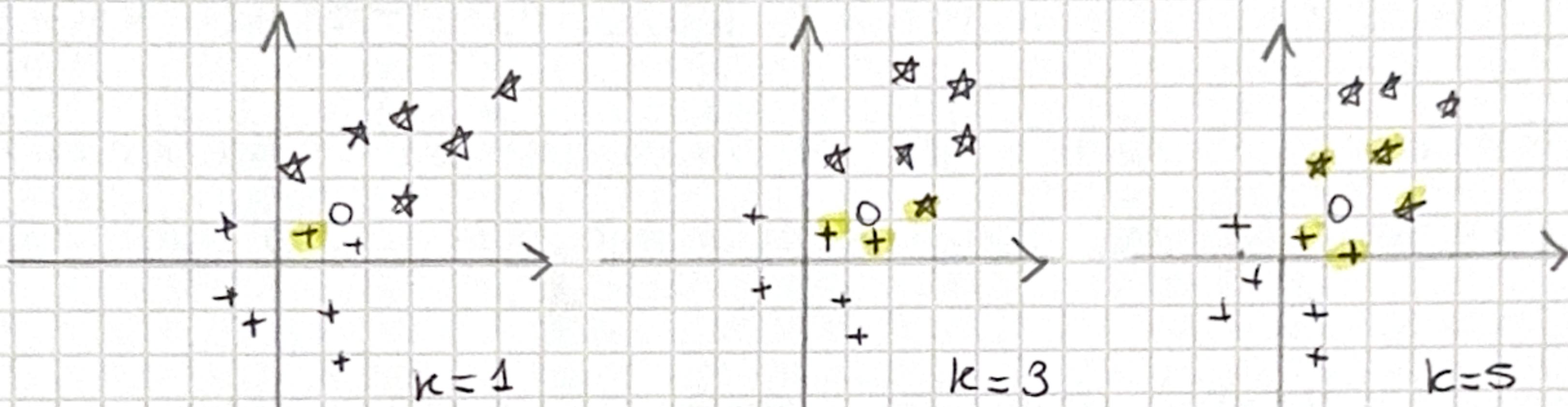
1. given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$  and a target function  $f: X \rightarrow Y$   
 kNN is an instance based model:
1. find  $k$  nearest neighbor of new instance  $x$
  2. assign to  $x$  the most common label among the majority of neighbors

likelihood of class  $c$  of new instance  $x$  is

$$P(c|x, D, x) = \frac{1}{k} \sum_{x_n \in N_k(x, D)} \mathbb{I}(t_n = c)$$

where  $\mathbb{I}(e) = \begin{cases} 1 & \text{if } e = \text{true} \\ 0 & \text{otherwise} \end{cases}$

2.



+ : 1     $k=1$

+ : 2     $x : 1$

$$P(+|x, D, 1) = \frac{1}{1} (1) = 1 \quad P(+|x, D, 3) = \frac{1}{3} (2) \quad P(+|x, D, 5) = \frac{1}{3} (2) =$$

classified as +

$$P(x|x, D, 3) = \frac{1}{3} (1)$$

classified as +

$$P(x|x, D, 5) = \frac{1}{5} (3) =$$

classified as \*

### EX. 6

1. Boosting is an ensemble method. the main idea is instead of train one complex simple classifier we train different classifiers and we combine them result  
 in boosting we are training models in sequence each model is based on the error of the previous model

### 2. AdaBoost

I. minimizing error function:

$$\mathcal{J}_m = \sum_{n=1}^N w_n^{(m)} f(y_m(x_n) \neq t_n)$$