

2. we examine least squares $EJ(w) = \frac{1}{2} (t - \phi(w))^T (t - \phi(w))$ and
 $w^* = \phi^T t$

we can also use sequential algorithm $\hat{w} \leftarrow w + \eta [t_n - w^T \phi(w)] \phi(x_n)$

ML EXAMS 18/02/2018 A

Exercise 1

1.

if $C = c_1$ and $B = b_1$ then YES NO ①

if $C = c_1$ and $B = b_2$ then YES ②

if $C = c_2$ and $B = a_2$ then YES ③

if $C = c_2$ and $B = a_2$ and $B = b_1$ ④

if $C = c_2$ and $B = a_2$ and $B = b_1$ ⑤

if $C = c_2$ and $B = a_3$ then NO ⑥

if $C = c_3$ then NO ⑦

2.

is commitment with S_1 because of ①

is commitment with S_2 because of ④

NOT is commitment with S_3 because of ⑦

NOT is commitment with S_4 because of ⑤

Exercise 2

1. the maximum a posteriori hypothesis is $h_{MAP} = \arg\max_{h \in H} \frac{P(D|h)}{P(D)}$

if we assume $P(h_j) = P(h_i)$ we simplify and

we obtain $h_{ML} = \arg\max_{h \in H} P(D|h)$

2. Bayes is an optimal classifier that returns the most probable hypothesis prediction:

$$V_{NB} = \arg\max_{v \in V} \sum_{h \in H} P(v|x, D) p(h_i|x, D) = \sum_{h \in H} P(v|x, h_i) P(h_i|D)$$

3. we can use it if we have analytical solution or if hypothesis model is not too large otherwise it's not practical

Question 3

1. in linear regression we have a target function $f: \mathbf{x} \rightarrow \mathbf{y}$ with $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}$ and the model is:

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=2}^m w_j \phi_j(\mathbf{x}) \text{ in linear in } \mathbf{w}$$

we want find $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{arg\,min}} E_D(\mathbf{w})$ where $E_D(\mathbf{w})$ is error function

given a dataset $D = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$ target value t_n is given by $y(\mathbf{x}, \mathbf{w})$ affected by additive noise ϵ

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon$$

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

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

assume observations iid

$$P(t_1, \dots, t_m | \mathbf{x}_1, \dots, \mathbf{x}_n, \beta) = \prod_{n=1}^m N(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

so maximum likelihood (zero-mean Gaussian prior)

$$\underset{\mathbf{w}}{\operatorname{arg\,max}} P(t_1, \dots, t_m | \mathbf{x}_1, \dots, \mathbf{x}_n, \beta)$$

conventional least mean square error

$$\underset{\mathbf{w}}{\operatorname{arg\,min}} E_D(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{arg\,min}} \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2$$

2. BATCH MODE: we consider all dataset \rightarrow require high computational power

$$\Delta \mathbf{w}_i = M \sum_{(x, t) \in D} (t - o(x)) \mathbf{x}_i$$

MINI-BATCH MODE: we choose a small subset SCD

$$\Delta \mathbf{w}_i = M \sum_{(x, t) \in S} (t - o(x)) \mathbf{x}_i$$

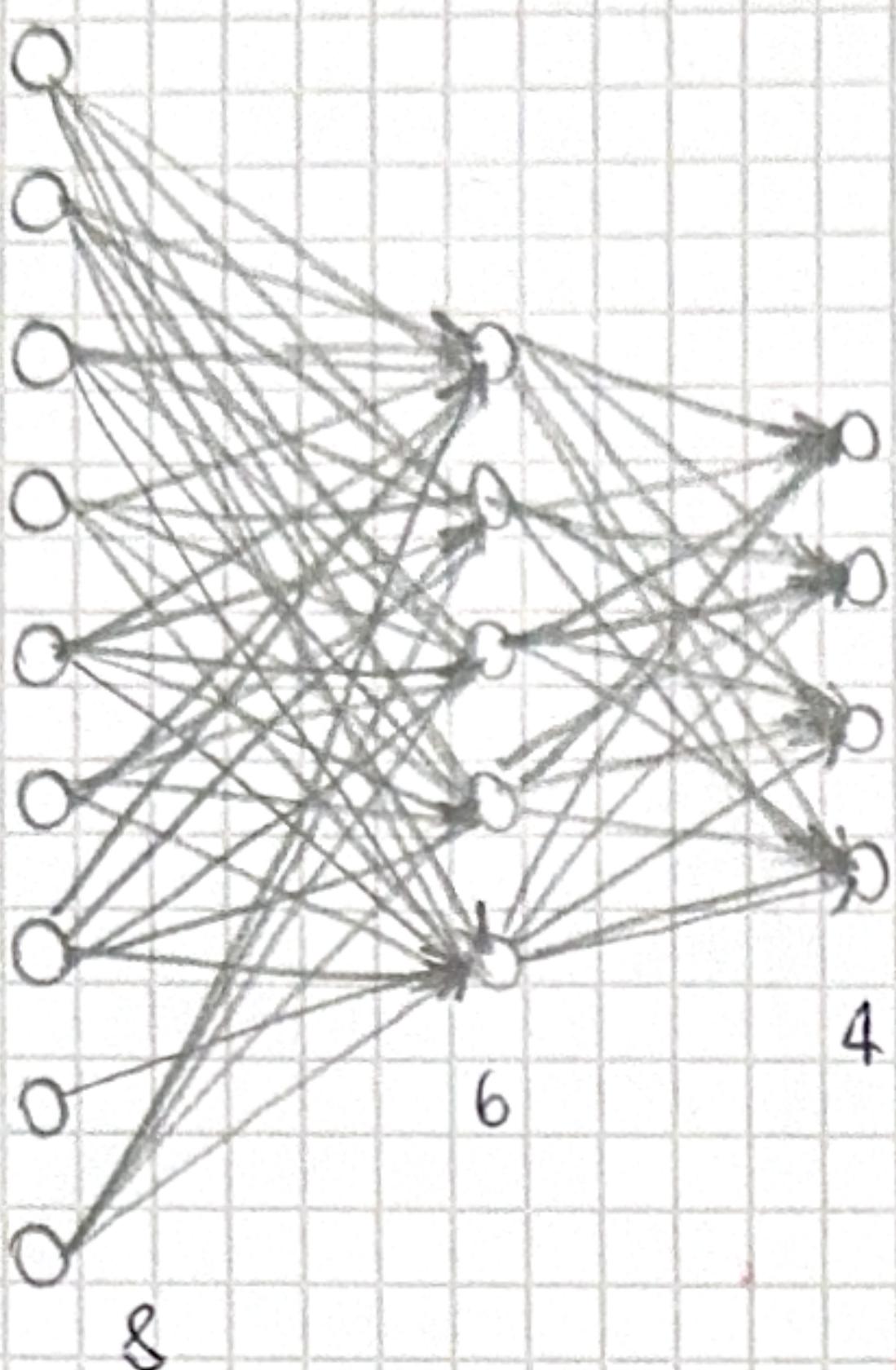
INCREMENTAL MODE: choose one sample $(x, t) \in D$ we have high variance

$$\Delta \mathbf{w}_i = M (t - o(x)) \mathbf{x}_i$$

$o(x) = \mathbf{w}^T \mathbf{x}$ for unthresholded $o(x) = \operatorname{sgn}(\mathbf{w}^T \mathbf{x})$ for thresholded

Exercise 4

1.



$$\# \text{parameters} = (8+1) \cdot 5 + (5+1) \cdot 4 = 63$$

2. and 3.

Backprop is an algorithm used to compute the gradient and the gradient will be propagate throughout all the network is not affected by overfitting and local minima

Exercise 5

1. the dataset is separable because exist an hyperplane that divided our dataset into two regions such that different classified instances are separated

2. in this case I use a polynomial kernel function $k(x, x') = (\beta x^T x' + \gamma)^d$

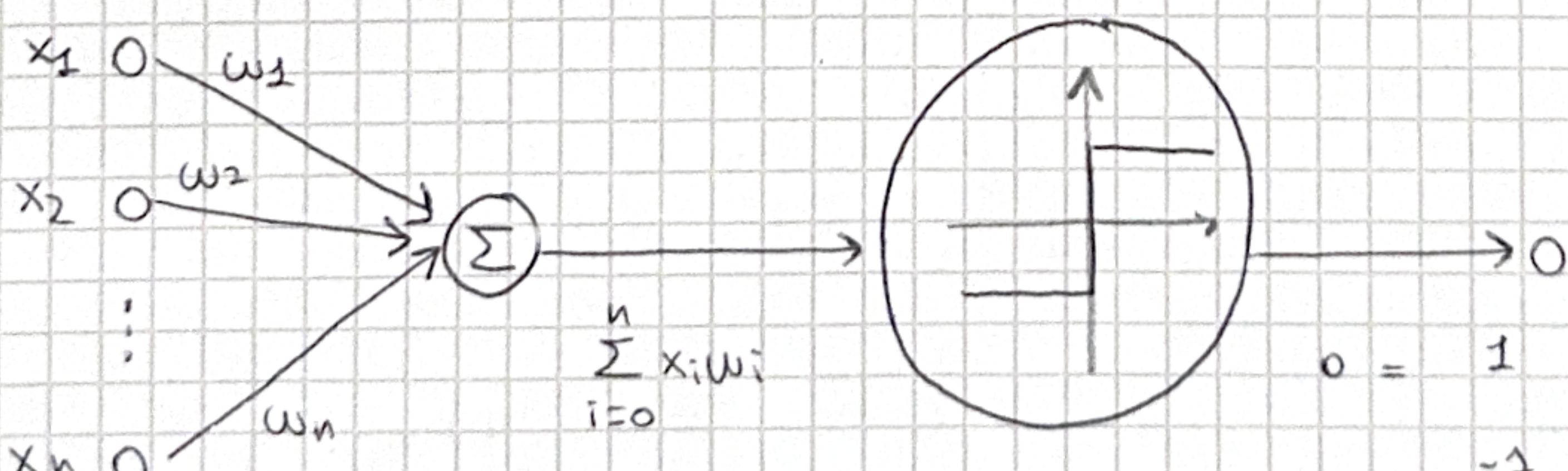
3. $y(x, \alpha) = \text{sign}(w_0 + \sum_{n=1}^N \alpha_n x_n^T x)$ applying kernel trick

$$y(x, \alpha) = \text{sign}(w_0 + \sum_{n=1}^N \alpha_n k(x_n, x))$$

$$w_0 = \frac{1}{|SV|} \sum_{x_i \in SV} (t_i - \sum_{x_j \in SV} \alpha_j t_j k(x_i, x_j))$$

Exercise 6

1.



$$o = \begin{cases} 1 & \text{if } \sum_{i=0}^n w_i x_i \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

we want to find w that minimizes error function (squared error)

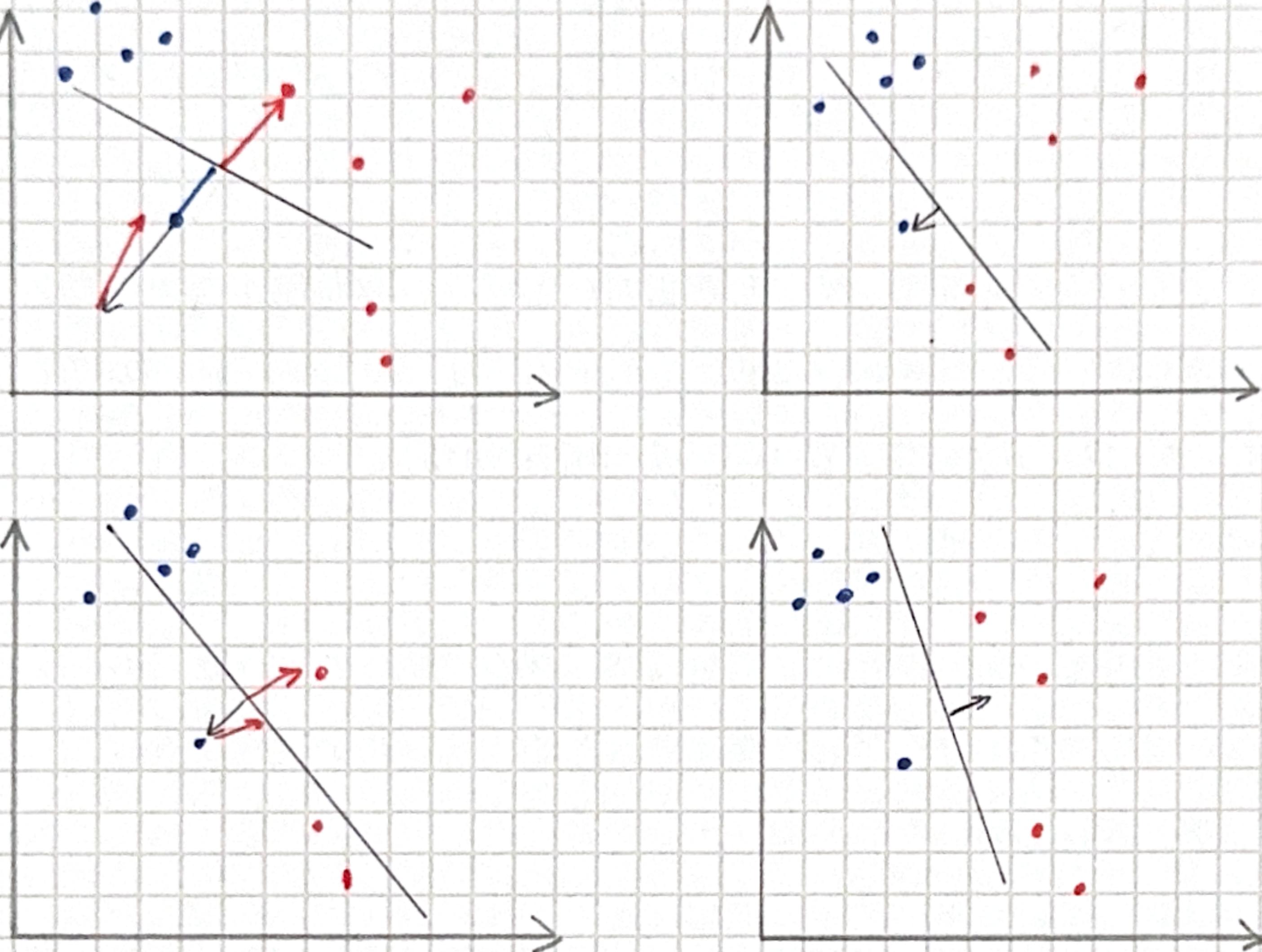
$$E_0(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top x_n)^2 = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top x_n)^2$$

$$\frac{\partial E_0(w)}{\partial w} = \sum_{n=1}^N (t_n - w^\top x_n) (-x_{n,m})$$

we find w^* in a sequential way by applying

$$w_i \leftarrow w_i + \Delta w_i \quad \text{where} \quad \Delta w_i = -\eta \frac{\partial E}{\partial w_i} = \eta \sum_{n=1}^N (t_n - w^\top x_n) x_{n,i}$$

2.



ML EXAMS 12/02/2018 A

EX.1

1. In a Reinforcement Learning problem we have a dataset $D = \{(x_0, a_1, r_1, x_1, \dots, a_n, r_n, x_n)\}_{n=1}^N$ and we want learning behavior $\pi: S \rightarrow A$, meaning we want to find the optimal policy function that maximizes the reward. In RL problems we don't have output and input is triple tuple action, reward, state

2. for each x, a initialize table entry $\hat{Q}_{0,1}(x, a) \leftarrow 0$
observe current state x
for each time $t=1, \dots, T$ (until termination condition)
 - choose an action a
 - execute action a
 - observe the new state x'