

2) In MDP we have  $\langle X, A, S, r \rangle$

$X$ : finite set of states  $A$ : finite set of actions

$S: X \times A \rightarrow X$  is the transition function and  $r: X \times A \rightarrow \mathbb{R}$ : reward function

In POMDP we have  $\langle X, A, Z, d, R, O \rangle$  where

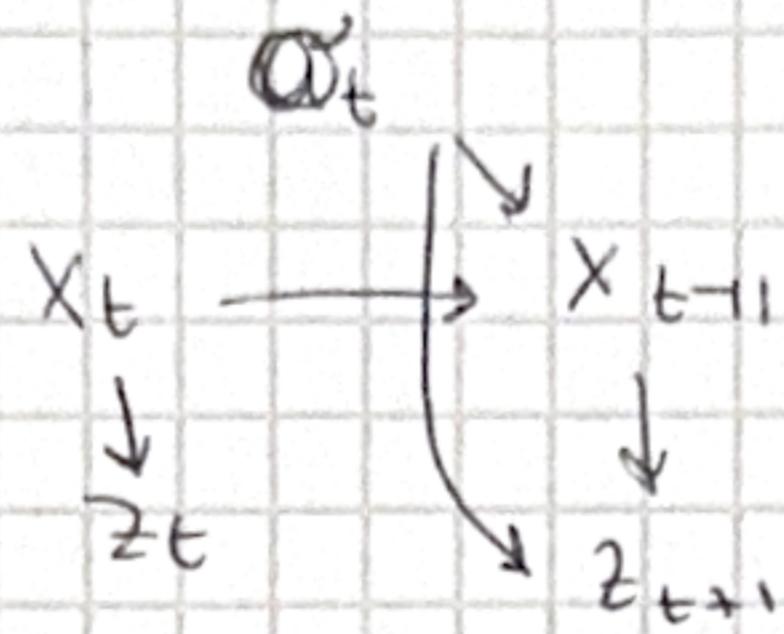
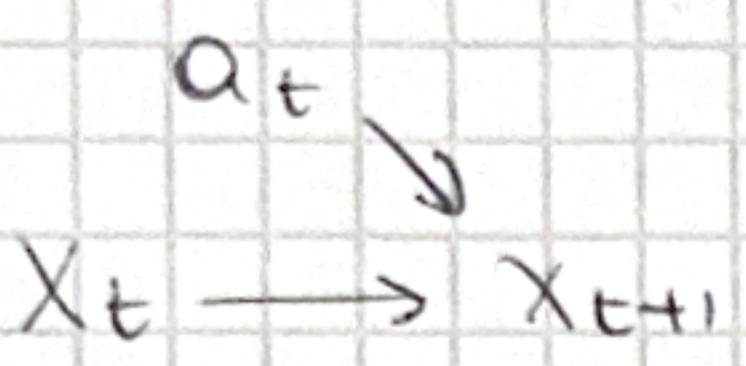
$X$ : finite set of states  $A$ : finite set of actions  $Z$ : finite set of

$d: X \times A \rightarrow X$  transition function  $R: X \times A \rightarrow \mathbb{R}$  is the observation

reward function and  $O$  is the probability of observation and

$P(x_0)$  is the probability of the initial state

Ex 3.



Ex. 6

$$1. \dim(w_1) = 128 - 10 = 1280 \quad \dim(w_2) = 50 \cdot 10 = 500$$

$$2. \text{use } y(x) = w^T \max(0, w^T x + c) + b \quad h = y(w^T x + c)$$

$$g(\alpha) = \max(0, \alpha) \rightarrow \text{ReLU}$$

ML EXAMS 18/01/2018

Ex. 1

input image  $= 1242 \times 378 \times 3$  (RGB) kernel  $= 5 \times 5$  pad  $= 2$  stride  $= 1$

$$w_{out} = \frac{w_{in} - w_k + 2p}{s} + 1 = 1242$$

$$d_{out} = \frac{d_{in} - (w_k + 2p) + s}{s} + 1 = 378$$

number of trainable parameters:  $5 \times 5 \times 3 \times 64$

$$w_{out} = \frac{w_{in} - dp}{2} + 1 = 621 \quad \left\{ \text{apply pool 1} \right.$$

$$d_{out} = \frac{d_{in} - dp}{2} + 1 = 183$$

$$w_{out} = 310 \quad h_{out} = 84 \quad C_1 \rightarrow 1242 \times 378 \times 69$$

$$w_{out} = 78 \quad h_{out} = 24 \quad P_1 \rightarrow 621 \times 189 \times 64$$

$$C_2 \rightarrow 310 \times 84 \times 128 \quad P_2 \rightarrow 78 \times 24 \times 128$$

$$2) 1 \text{ param} = 3 \times 5 \times 5 + 64 = 416$$

$$2 \text{ param} = 64 \times 3 \times 3 \times 128 + 121 = 78884 \quad 73856$$

## EX.2

1. In a synchrotron imager of  $12 \times 12$  pixels no the dimensionality of data map is  $2^{12 \times 12}$ , the intrinsic dimensionality is 3  $\rightarrow$  two coordinates for the center one for rotation

2. We want to maximize the variance after the projection to normal dimension  $u_1$  we know the projected norm is one  $u_1^T x_n$  subject to  $u_1^T u_1 = 1$

$$\text{the variance is } \frac{1}{N} \sum_{n=1}^N (u_1^T x_n - \bar{u}_1^T \bar{x})^2 = u_1^T S u_1$$

$$\text{with } S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$

$$\text{we want maximize } \max_{u_1} u_1^T S u_1$$

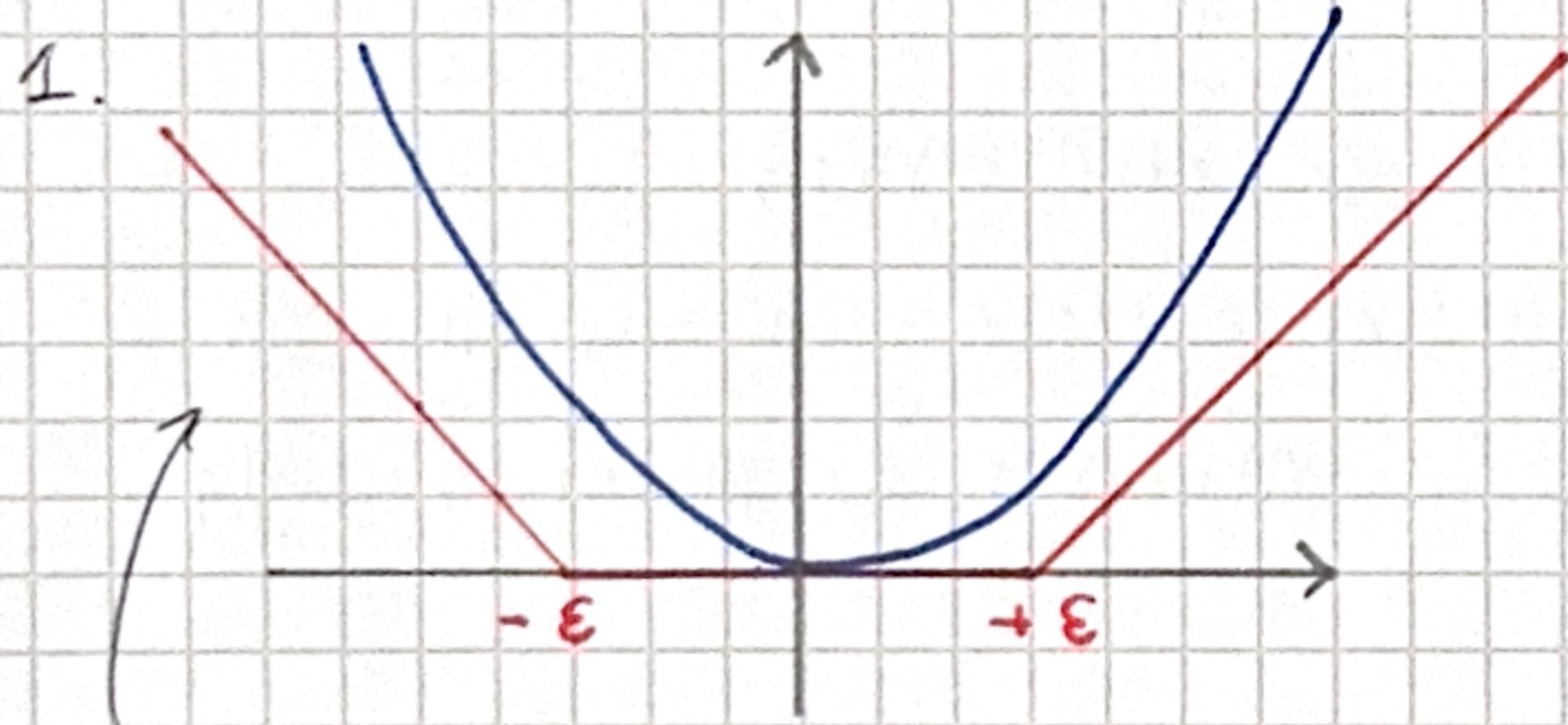
if we keep derivative and we impose that is equal to 0

$$S u_1 = \lambda_1 u_1 \Rightarrow u_1^T S u_1 = \lambda_1$$

$S$  is the eigenvector associated to the largest eigenvalue  $\lambda_1$   
this is called first principal component

3. No, usually  $N$  is greater than the number of intrinsic dimensions because the PCA are not latent variable

## EX.3



is difficult to solve  
because is not  
differentiable

doesn't penalize  
too much outliers

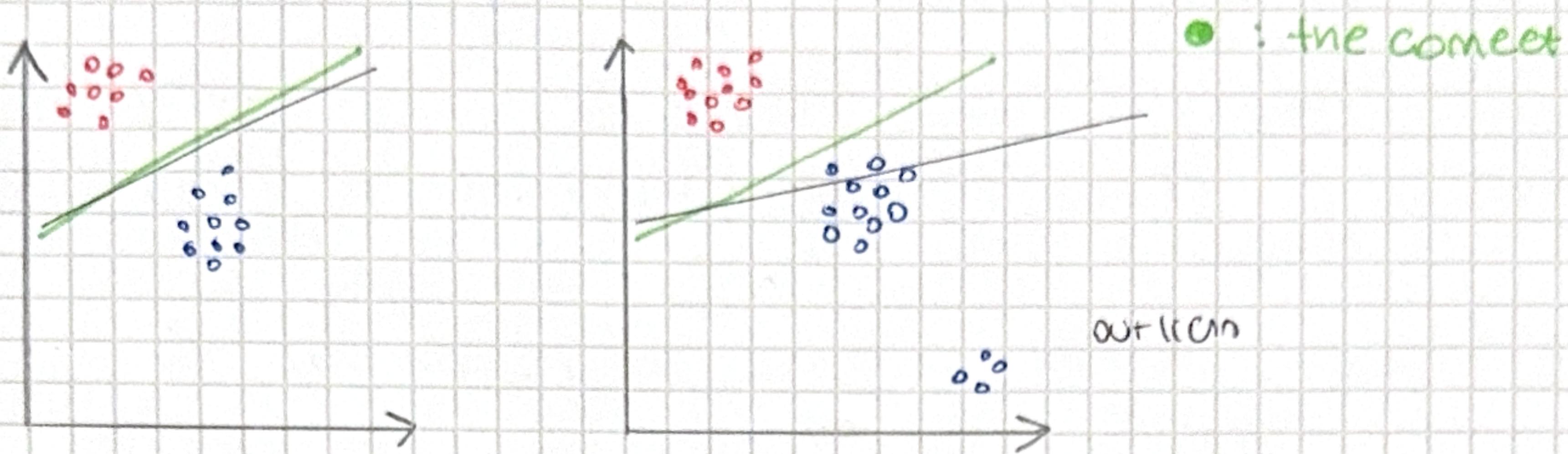
- 2.
- $\epsilon < 0$  is on or inside the correct margin boundary
  - $0 < \epsilon_n < 1$  if points inside margin but correct side
  - $\epsilon_n > 1$  if points in the incorrect boundary

optimization with slack variables become:

$$w^*, w_0^* = \underset{w, w_0}{\operatorname{argmin}} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \epsilon_n$$

the above variables are introduced to allow minimization or margin violations in soft margin classifications

EX. 4



LEAST SQUARES

given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$  find the line of minimization  
 $y(x) = \tilde{w}^\top \tilde{x}$

the goal is minimize the sum of squared error function:

$$E(\tilde{w}) = \frac{1}{2} \text{Tr} \{ (\tilde{x}\tilde{w} - \tau)^\top (\tilde{x}\tilde{w} - \tau) \}$$

closed form  $\tilde{w} = (\alpha \tilde{x}^\top \tilde{x})^{-1} \tilde{x}^\top \tau$

learned model  $y(x) = \tilde{w}^\top \tilde{x} = \tau^\top (\tilde{x}^\top) \tilde{x}$

$$k = \underset{j \in \{1, \dots, k\}}{\operatorname{argmax}} |y_j(x)|$$

is not robust because it simply based on distances

is called in this way because the trace operation is simply an sum and the product is between one matrix to other in only the square.

EX. 3

1.

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h|D) = \underset{h \in H}{\operatorname{argmax}} \frac{P(D|h) P(h)}{P(D)}$$

$$P(d_{hi}) = P(d_{hj}) \quad h_{ML} = \underset{h \in H}{\operatorname{argmax}} (P(D|h) P(h))^{1/2}$$

2. BCE is an optimal estimator that returns always the optimal solution given a target function  $f: X \rightarrow \mathbb{R}^Y$

$$P(v|x, D) = \sum_{h \in H} P(v|x, D, h) P(h|x, D) = \sum_{h \in H} P(v|_{\text{node}}, h) P(h|D)$$

$$v_{ce} = \underset{h \in H}{\operatorname{argmax}} \sum_{h \in H} P(v|x, h) P(h|D)$$

3. Use Naive Bayes classification if the hypothesis space is not too large and it have an analytical solution, otherwise is not practical ~~and~~ but  
Note: use Naive Bayes classifier

EX. 6

Once we augment more information the evolution of a dynamic system doesn't depend on previous state, actions and observations. The current state contains all information needed to predict the future. Future state are conditionally independent of past states and past observations statistically independent.

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

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

HMM can be describe  $\langle X, Z, P_0 \rangle$  where  $X$  is set of observations and  $Z$  is initial distribution

MDP has the property of fully observability

MDP

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

$\downarrow$

at

HMM

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

$\downarrow$

$z_t \quad z_{t+1}$

$\downarrow$

$z_n$

ML EXAMS 15/10/2019

EX. 1

1. how many times an illegal insertion of class C in classifiers are in the class ej

2.

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

	A	B	C	→ predicted
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 ~~correct~~ divide by the all of elements of matrix

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