

## EX. 1

1. given a dataset  $D$  and an hypothesis space  $H$  a model is overfitting when exist another hypothesis  $h' \in H$  such that :
- $$\underset{D}{\text{error}}(h) > \underset{D}{\text{error}}(h') \quad \text{and} \quad \underset{S}{\text{error}}(h) < \underset{S}{\text{error}}(h')$$
- two solutions are

2. if we have two decision trees and one is deeper respect to other for sure we have overfitting we have two possible methods to avoid

F. REDUCE error pruning : using validation dataset, we prune the node that doesn't improve validation accuracy

II. naive post-prune : if the accuracy after pruning is greater respect before we still have overfitting

## Ex 2

1. naive bayes classifier uses conditional independence to approximate the solution:

$x$  is conditionally independent from  $y$  given  $z$

$$P(x \neq y | z) = P(x|y, z) P(y|z) = P(x|z) P(y|z)$$

assume target function  $f: x \rightarrow V$  were each inten of  $x$  are describe by  $(a_1, \dots, a_n)$

$$\operatorname{argmax}_{v_j \in V} P(v_j | x, D) = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, \dots, a_n, D) = \operatorname{VMAP}$$

with bayes assumption

$$P(a_1, \dots, a_n | v_j, D) = \prod_j P(a_i | v_j, D)$$

$$VNB = \operatorname{argmax}_{v_j \in V} P(v_j | D) \prod_i P(a_i | v_j, D)$$

2. in a classification task: the target function  $f: x \rightarrow y$  with  $x = \{ \text{title} \times \text{author} \times \text{obs} \times \text{rate} \}$  and  $y = \{ \text{ML}, \text{KRR}, \text{PL} \}$  the dataset is  $D = \{(d_i, y_i)\}_{i=1}^N$

$$VNB = \operatorname{argmax}_{v_j \in V} P(v_j | D) \prod_i P(d_i | v_j, D)$$

we know  $P(d_i | V, D) = \prod_{i=2}^L P(d_i = w_i | V, D)$  where  $L = \text{length}(d_i)$

and  $P(d_i = w_i | V, D)$  probability that the word occurs in the document  $i$

$$\hat{P}(d_i | V, D) = \frac{\text{TF}_{ij} + 2}{\text{TF}_i + 1}$$

### EX.3

- Because we have two classes 0, 1 it consider a logistic regression that is a probabilistic discriminative model based on maximum likelihood likelihoods

$$P(t|\vec{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{(1-t_n)} \text{ with } y_n = P(C_1 | \vec{x}_n)$$

- The parameters that the model has to learn are  $a_i$  and  $n_i$

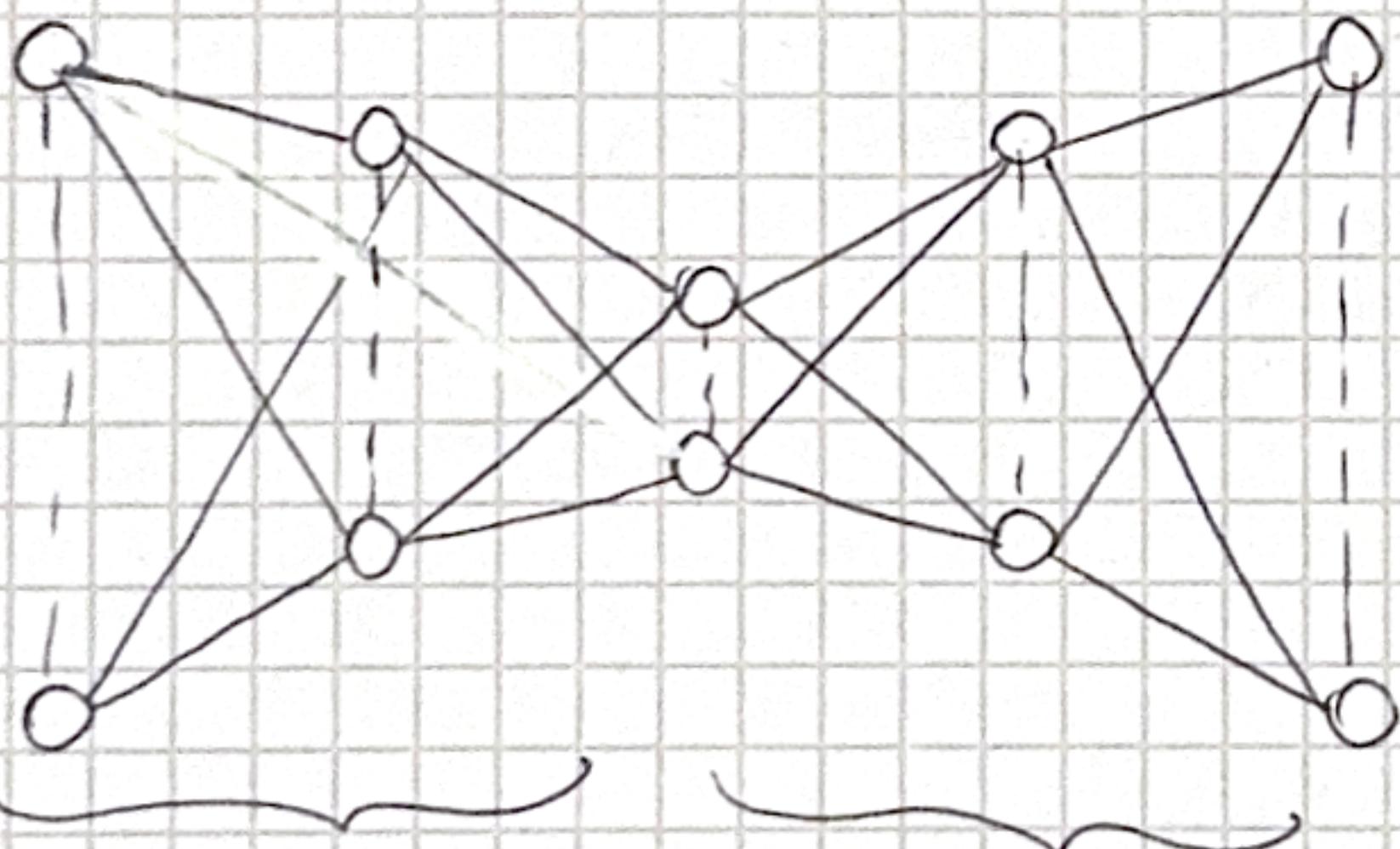
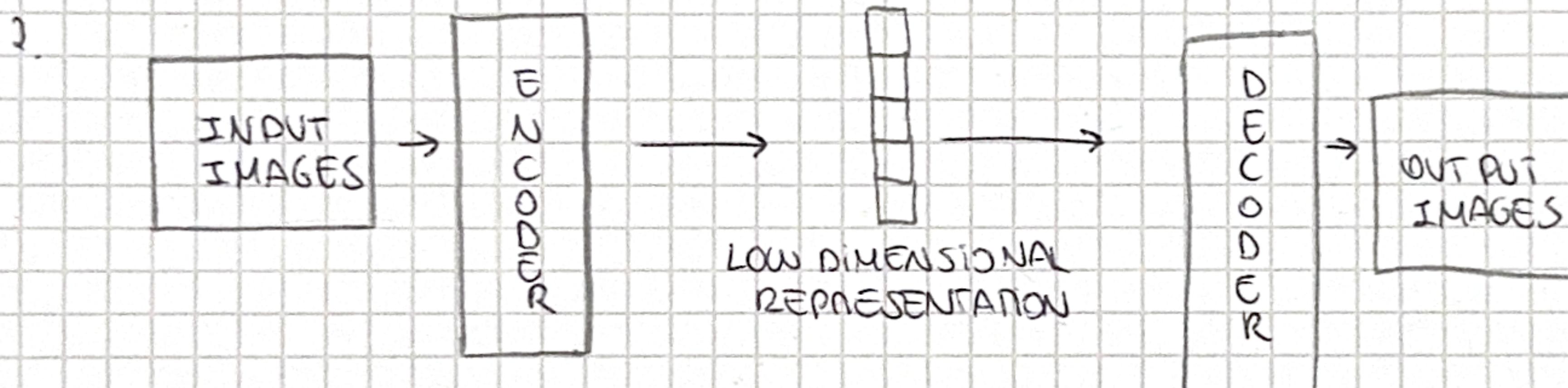
- A suitable error function will be cross entropy loss

$$E(\vec{w}) = -\ln(P(t|\vec{w})) = -\sum_{n=1}^N [t_n \ln(y_n) + (1-t_n) \ln(1-y_n)]$$

### EX.4

- An autoencoder is combination of two ANN a encoder and decoder the training is based on reconstruction loss.

An autoencoder takes an input and output the same sample  $x_n$  from the dataset (is used to deal with data far away from Gaussian distribution)



COMPACT REPRESENTATION OF INPUT

HOW TO RECONSTRUCT THE INPUT FROM THE COMPACT REPRESENTATION

### EX.5

- In a dynamic system we have the property of full observability if we can see the state resulting from the execution of action also if we have a non-deterministic action we cannot predict the state before we don't see it

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

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

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

In POMDP we have  $\langle X, A, \delta, \mathcal{O}, r, o \rangle$  where

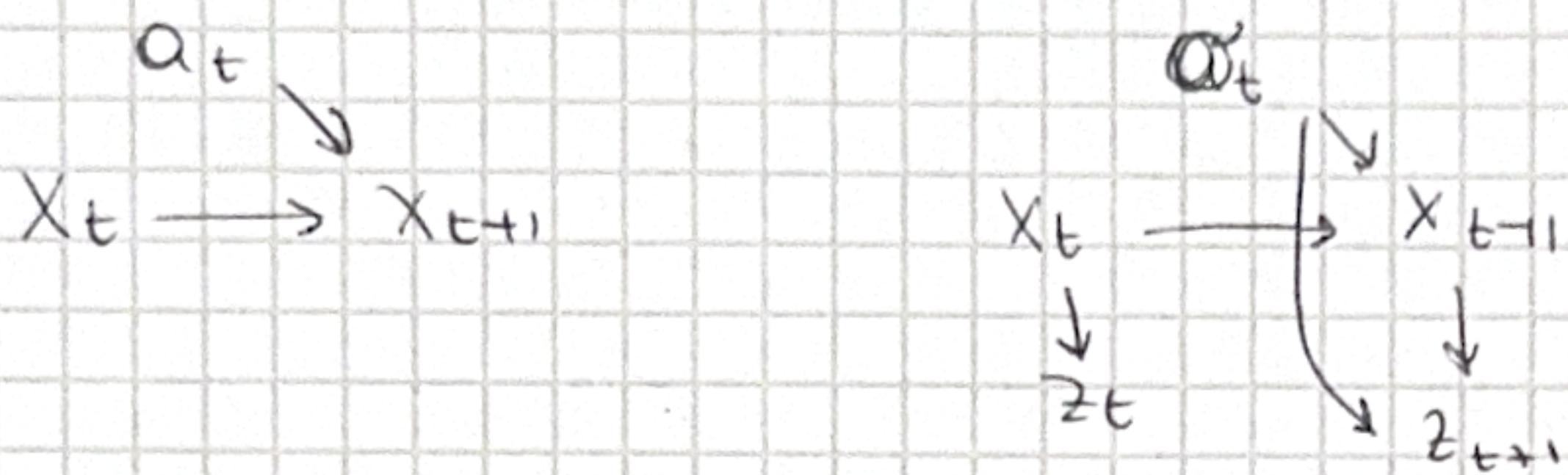
$X$ : finite set of states  $A$ : finite set of actions  $\mathcal{O}$ : finite set of

$\delta: X \times A \rightarrow X$  transition function  $r: X \times A \rightarrow \mathbb{R}$  is the observations

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

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

Ex 3.



Ex. 6

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

$$2. \text{g}_{\theta}(x) = w^T \max(0, w^T x + c) + b \quad h = g(w^T x + c) \\ g(x) = \max(0, x) \rightarrow \text{ReLU}$$

ML EXAMS 18/01/2018

Ex. 1

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

$$w_{out} = \frac{w_{in} - \cancel{w_k} + 2p}{s} + 1 = 1242 \quad \left. \begin{array}{l} \text{after } c_1 \\ \hline \end{array} \right\}$$

$$d_{out} = \frac{d_{in} - \cancel{w_k} + 2p}{s} + 1 = 378 \quad \left. \begin{array}{l} \text{after } c_1 \\ \hline \end{array} \right\}$$

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

$$621 \times 189 \times 64$$

$$310 \times 84$$

$$w_{out} = \frac{w_{in} - dp}{2} + 1 = 621 \quad \left. \begin{array}{l} \text{after pool 1} \\ \hline \end{array} \right\}$$

$$d_{out} = \frac{d_{in}}{2} + 1 = 189$$

$$w_{out} = 310 \quad d_{out} = 84 \quad \left. \begin{array}{l} \text{C1} \rightarrow 1242 \times 378 \times 69 \\ \hline \end{array} \right\}$$

$$w_{out} = 78 \quad d_{out} = 24 \quad \left. \begin{array}{l} \text{P1} \rightarrow 621 \times 189 \times 64 \\ \hline \end{array} \right\}$$

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

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

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