

## EXERCISE 1

- Provide a formal (domain-independent and solution-independent) definition of overfitting.
- Discuss the problem of overfitting in learning with Decision Trees and illustrate possible solutions to it.

Let  $h$  to be an hypothesis, we denote  $\text{error}(h)$  the TRUE error of  $h$ .

$h: X \rightarrow Y$ ,  $\text{error}(h) = \frac{1}{|X|} \sum_{x \in X} \mathbb{I}(h(x) \neq f(x))$  where  $f$ : Target Funct. and  $X$  is finite.

If  $D$  is the dataset, the sample error is  $\text{error}_D(h) = \frac{1}{|D|} \sum_{x \in D} \mathbb{I}(h(x) \neq f(x))$

$h$  is overfitting if  $\exists h'$  such that:

$$\text{error}_D(h) < \text{error}_D(h') \quad \text{and} \quad \text{error}(h) > \text{error}(h')$$

The decision trees overfits when we have a leaf for each sample in the dataset, we can overcome this with pruning:

```

PRUNING(h:tree) {
    until termination condition do {
        n = random node from h
        h' = we cut from h the subtree with root n
        a = accuracy(h)
        a' = accuracy(h')
        if a' > a {
            h = h'
        }
    }
    return h
}

```

A possible term. cond. might be: the accuracy didn't increased in the last  $k$  iterations.

## EXERCISE 2

- Describe the Naive Bayes Classifier and highlight the approximation made with respect to the Bayes Optimal Classifier.

- Provide design and implementation choices for solving the following problem through Naive Bayes Classifier:

Classification of scientific papers in categories according to their main subject. The categories to be considered are: ML (Machine Learning), KR (Knowledge Representation), PL (Planning). Data available for each scientific paper are: title, authors, abstract and publication site (name of the journal and/or of the conference).

Given a dataset  $D = \{(x_i, t_i)\}_{i=1}^N$  the optimal Bayes classifier is given

by considering  $P(c|x, D) = \sum_{h \in H} P(c|x, D, h)P(h|x, D) = \sum_{h \in H} P(c|x, h)P(h|D)$

where  $H$  is the set of all the hypothesis. The naive Bayes classifiers assume that the attributes of the samples are independent:

$$\text{if } x_i = (x_{i1}, \dots, x_{id})^T \quad P(x_{ij}=z | x_{ik}=z') = P(x_{ij}=z) \quad \forall j \neq k.$$

$$\Rightarrow P(c|x, D) = P(c|x_{i1}, \dots, x_{id}, h) = P(x_{i1}, \dots, x_{id} | c, D)P(c|D) \cdot \frac{1}{a}$$

$$\Rightarrow \text{argmax}_c P(c|x, D) = \text{argmax}_c \prod_{j=1}^d P(x_{ij} | c, D)P(c|D) =$$

$$\text{argmax}_c P(c|D) \cdot \prod_{j=1}^d \frac{\text{number of } C \text{ sample}}{\text{number of sample}}$$

$$\downarrow \frac{\text{number of } C \text{ sample with } x_{ij}(j) = x_{ij}}{\text{number of sample}}$$

A row of the dataset is

TITLE AUTHORS ABSTRACT SITE |

$$(T_i, A_i, AB_i, SI_i) \in D \text{ and } C_i \in \{ML, PL, KR\}$$

We have  $P(ML|D) = \frac{|\{(x_i, C_i) \in D : C_i = ML\}|}{|D|}$   $P(KR|D) = \frac{|\{(x_i, C_i) \in D : C_i = KR\}|}{|D|}$

$$P(PL|D) = \frac{|\{(x_i, C_i) \in D : C_i = PL\}|}{|D|}$$

Then, For a given title  $T^* \in \{\text{All possible titles}\}$  we have

$$P(T^* | C, D) = \frac{|\{(x_i, C_i) \in D : C_i = C \wedge T_i = T^*\}|}{|\{(x_i, C_i) \in D : C_i = C\}|}$$

The same follows for the other attributes. A new  $x'$  is classified.

$$c(x') = \arg \max_C P(C|D) \cdot \prod_{j=1}^5 P(x'_{c_j} | C, D)$$

#### EXERCISE 3

Consider a dataset  $D = \{(a_1, s_1), p_1\}, \dots, (a_N, s_N), p_N\}$  containing the number of hours  $a_i$  a student has attended a course, the number of hours  $s_i$  s/he has studied for the course and whether or not s/he has passed the exam  $p_i = \{0, 1\}$ .

1. Define a model based on logistic regression that, given the values of  $a$  and  $s$ , estimates whether a student passes the exam or not.
2. Discuss which are the parameters of the model that have to be learned based on the given data.
3. What is a suitable error function for learning the parameters of the model?

$p = \text{all the labels in } D$

Given  $D$ , we estimate  $P(P|x, D) = \prod_{i=1}^N \pi_i^{p_i} (1-\pi_i)^{1-p_i}$  con  $\pi_i = P(p_i = 1|x_i) = \sigma(w^T x_i)$

$\Rightarrow \sigma(w^T x_i) = (1 - \exp(w^T x_i))^{-1}$ . We want to maximize  $P(P|x, D)$  that depends on  $w$ .

$\ln P(P|x, D) = \sum_{i=1}^N p_i \ln(\pi_i) + (1-p_i) \ln(1-\pi_i)$  recall that  $\pi_i$  depends on  $w$ .

so:  $E(w) = - \sum_{i=1}^N p_i \ln(\pi_i) + (1-p_i) \ln(1-\pi_i)$  and we can find the optimal  $w^*$  by using a method based on  $\nabla E$  and  $H(w) = \nabla \nabla E$  (Hessian).

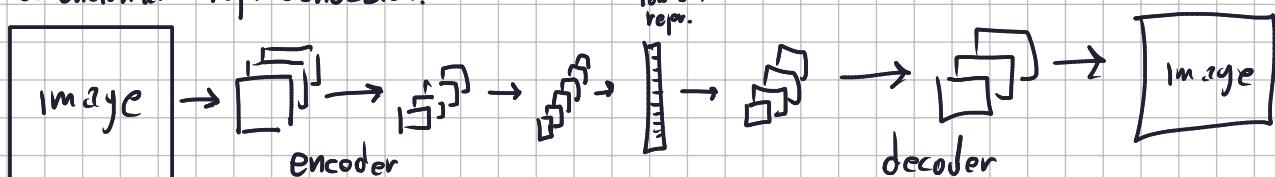
#### EXERCISE 4

1. Briefly describe what is the architecture of an autoencoder and its purpose.
2. Draw an example of autoencoder.

An autoencoder is a NN that is composed by two NN:

1) the first one transform the input (like an image) in a low dimension representation, usually a 1-dim. tensor.

2) the second, from the given low-dimensional encoding, reconstruct the high-dimensional representation.



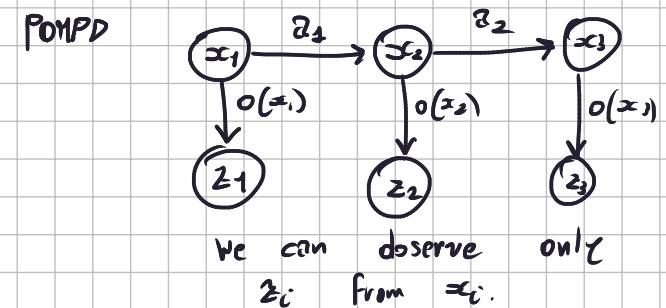
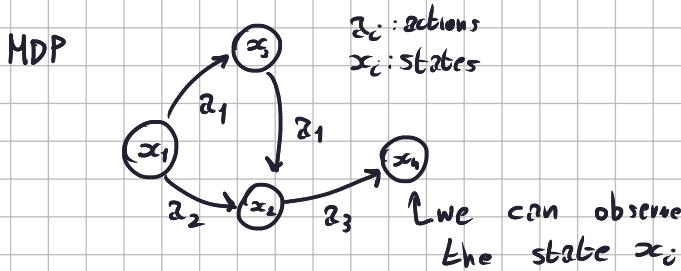
## EXERCISE 5

1. Describe the concept of full observability in models representing dynamic systems.
2. Describe the difference between a Markov Decision Process (MDP) and a Partially Observable Markov Decision Process (POMDP), referring to their formal models.
3. Draw and explain the graphical models of MDP and POMDP.

A dynamic system is Fully observable if we can read his entire state vector  $\mathbf{x}_t$ .

A POMDP is not Fully observable, From  $\mathbf{x}_t$  we can get  $\mathbf{z}_t$  such that :  $\mathbf{z}_t$  is a restriction of  $\mathbf{x}_t$

OR there is a distribution  $D(\mathbf{z}, \mathbf{x}_t)$  over  $\mathbf{z}$  that gives the probability of observing a given  $\mathbf{z}$ :  $P(\mathbf{z}_t = \mathbf{z}^* | \mathbf{x}_t) = D(\mathbf{z}, \mathbf{x}_t)$



## EXERCISE 6

Consider a two-layers ANN which receives in input vectors  $\mathbf{x}$  of dimension 128 and produces output vectors  $\mathbf{y}$  of dimension 10. The hidden layer of the ANN is composed of 50 units which use the ReLU activation function. The output units use a linear activation function. The weight matrices of the hidden and output layers are denoted  $W_1$  and  $W_2$ , respectively.

1. Provide the dimensions of the weight matrices  $W_1$  and  $W_2$ .

2. Provide the formula explicitly stating how the values of  $\mathbf{y}$  are computed given an input vector  $\mathbf{x}$  in terms of the weight matrices and the activation functions (you can ignore the bias terms).

Col: 10  
Row: out

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 4 \\ 3 & 4 & 5 \\ 4 & 5 & 6 \end{bmatrix} \left| \begin{array}{c} \text{ReLU} \\ \text{Linear} \end{array} \right| \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_{10} \end{bmatrix}$$

$W_1$  is a map from  $\mathbb{R}^{128}$  to  $\mathbb{R}^{50}$  so  $W_1$  have 128 columns and 50 rows.

$W_2$  is a map from  $\mathbb{R}^{50}$  to  $\mathbb{R}^{10}$  so  $W_2$  have 50 columns and 10 rows.

$$h_1(\mathbf{x}) = \text{ReLU}(W_1 \mathbf{x}), \quad h_2 = W_2 h_1 \Rightarrow \mathbf{y}(\mathbf{x}) = W_2 \cdot \text{ReLU}(W_1 \cdot \mathbf{x})$$