

EXERCISE 1

- Provide the definition of *Confusion matrix* for a multi-class classification problem, formally explain the content of component $C_{i,j}$ of the matrix.
- Provide a numerical example of a confusion matrix for a 3-classes classification problem with a balanced data set including 200 samples for each class (600 samples in total) and an average accuracy around 70% for class 1, 80% for class 2 and 90% for class 3. The matrix must not be symmetric. Show the confusion matrix in two formats: with absolute values and with the corresponding percentage values. (Hint: use simple numerical values, so that you do not need to make complex calculations.)
- Compute the accuracy of the classifier for the numerical example provided above.

the confusion matrix is the matrix C where $C_{ij} = \text{number of objects of class } i \text{ classified as class } j$. The trace of C is the number of correct classification.

Example:

	CLASS		
	1	2	3
1	140	20	15
2	30	160	5
3	30	20	180

in perc.

$$\Rightarrow$$

	CLASS		
	1	2	3
1	23.3	3.3	2.5
2	5	26.6	0.8
3	5	3.3	30

the accuracy is

$$\frac{\text{trace}(C)}{\sum C_{ij}} = 0.8 = 80\%$$

EXERCISE 2

Consider a binary classification problem $X \rightarrow \{T, F\}$, with $X = \{T, F\}^3$, i.e. $(x_1, x_2, x_3) \in X$ and $x_i \in \{T, F\}$, and the dataset $D = \{(F, F, F), F\}, \{(F, T, T), T\}, \{(T, T, F), T\}, \{(T, F, T), T\}\}$. Consider the two hypothesis $h_1 = (x_1 \wedge \neg x_2 \wedge x_3) \vee x_2$ and $h_2 = (\neg x_1 \wedge x_2 \wedge x_3) \vee x_1$.

- Determine whether h_1 and h_2 are consistent with D , showing all the passages needed to answer.
- Assuming the likelihood probabilities $P(D|h_1) = 0.6$ and $P(D|h_2) = 0.8$ and the prior probabilities $P(h_1) = 0.2$ and $P(h_2) = 0.1$, determine the higher a posteriori hypothesis between h_1 and h_2 .

check for h_1 :

$$\begin{aligned} h_1(x_1) &= F = b_1 \quad \checkmark & h_1(x_3) &= T = b_3 \quad \checkmark \\ h_1(x_2) &= T = b_2 \quad \checkmark & h_1(x_4) &= T = b_4 \quad \checkmark \end{aligned} \Rightarrow h_1 \text{ is cons.}$$

check for h_2 :

$$\begin{aligned} h_2(x_1) &= F = b_1 \quad \checkmark & h_2(x_3) &= T = b_3 \\ h_2(x_2) &= T = b_2 \quad \checkmark & h_2(x_4) &= T = b_4 \end{aligned} \Rightarrow h_2 \text{ is cons.}$$

$$h_{MAP} = \underset{h}{\operatorname{argmax}} P(h|D) = \underset{h}{\operatorname{argmax}} P(D|h)P(h) \stackrel{\substack{P(D|h_1)P(h_1) = 3/25 \\ P(D|h_2)P(h_2) = 2/25}}{\Rightarrow} h_{MAP} = h_2.$$

EXERCISE 3

Consider a regression problem $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with a dataset $D = \{(x_n, t_n)\}_{n=1}^N$, where f is known to be non-linear in x .

- Describe a linear model for this problem and determine the trainable parameters and the size of the model (i.e., number of trainable parameters).
- Describe a solution of the problem in terms of least square error minimization. Define the error function corresponding to the model given above and illustrate a method to find a solution of the optimization problem.

We consider $\gamma(\boldsymbol{x}) = \mathbf{w}^T \phi(\boldsymbol{x}) + w_0$ where $\mathbf{w} \in \mathbb{R}^d$, $w_0 \in \mathbb{R}$ are the parameters

and $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a non linear Function.

We consider the error Function $E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \gamma_i)^2$ where $\gamma_i = \gamma(x_i)$. We have:

$$\frac{\partial E}{\partial w_j} = \frac{\partial}{\partial w_j} \left(\frac{1}{2} \sum_{i=1}^N (t_i - \gamma_i)^2 \right) = - \sum_{i=1}^N (t_i - \gamma_i) \phi(x)_i j$$

We consider $\nabla E = \left(\frac{\partial E}{\partial w_j} \right)_{j=1}^d$

and we apply the gradient descent

Initialize \mathbf{w}
until some termination condition do:

$$\mathbf{w} \leftarrow \mathbf{w} - \gamma \nabla E$$

where γ is a small positive real number.

EXERCISE 4

1. Describe the principle of soft margins used by SVM classifiers. Illustrate the concept also with a geometric example.
2. Draw a dataset for binary classification of 2D samples and show two solutions based on SVM with and without soft margin constraints. Choose a proper example that illustrates well the concept, i.e., in which the two solutions are significantly different.

We consider as margin, the distance between the linear discriminant and the closest sample to it: $\min_{x_i \in D} \frac{|\mathcal{E}(x_i)|}{\|w\|}$. The SVM aims to find the w that maximizes such margin, while classifying correctly all the samples:

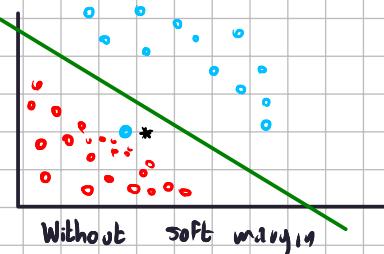
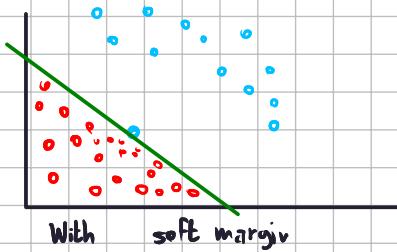
$$D = \{(x_i, t_i)\}_{i=1}^N, t_i \in \{+1, -1\} \Rightarrow \mathcal{E}(x_i)t_i \geq 0 \quad \forall i.$$

We can add a relaxation in the optimization problem:

$$\begin{cases} w^* = \arg \min_w \left(\min_{x_i \in D} \frac{|\mathcal{E}(x_i)|}{\|w\|} \right) + C \sum_{i=1}^N \xi_i^+ \\ \text{s.t. } t_i \mathcal{E}(x_i) \geq 0 \quad \forall i \\ \xi_i \geq 0 \quad \forall i \end{cases}$$

\Rightarrow we allow some samples to overcome the margin

$$\begin{cases} \xi_i = 0 \text{ in his side} \\ \xi_i \in [0, 1] \text{ violates margin} \\ \xi_i > 1 \text{ misclassified} \end{cases}$$



* may be a very noisy sample.

EXERCISE 5

Let D be a dataset containing the following input values $X = \{(3.3, 1.6), (7.5, 48.2), \dots, (98.3, 43.5), (87.2, 92.4)\}$ and target values $T = \{0, 2, \dots, 4, 3\}$.

Consider designing a Feedforward Neural Network for learning the function $t = f(x)$.

1. Explain what is a valid choice for the activation function of the output layer and for the loss function.
2. Provide some valid options for the activation functions of the hidden units.
3. Formally describe the Stochastic Gradient Descent (SGD) algorithm and illustrate its hyper-parameters.

Since the target values are not limited in $[0, 1]$, a linear activation for the output layer is mandatory. For the hidden units, a good choice might be the sigmoid, we observe how the samples have different scales for the attributes, so a normalization procedure of X would benefit the training process.

Given the parametric net $\gamma(w, x)$ with w parameters, the SGD follows:

SGD {

```

    initialize w with small random values
    while (!termination condition) {
        X_x = sample K points from X (randomly)
        g = 0
        for each w_i in w {
            g = g + ∂E / ∂w_i * e_i // e_i = [0, 0, ..., 1, 0, ...] ← i-th pos.
        }
        w = w - γ g
    }
  
```

$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ std basis vector

the hyper parameters are: the choice of the loss function $E(w)$, the choice of the termination condition, the choice of the batch size $K \in \mathbb{N}$, the choice of the step size $\gamma \in \mathbb{R}^+$.

EXERCISE 6

1. Describe the K-means algorithm in a formal way (i.e., with precise mathematical formulas and equations), including: input and output of the algorithm, its main steps, and the termination condition.
2. Draw a suitable 2-D data set for K-means.
3. Qualitatively simulate the execution of K-means in such 2-D data, showing at least three steps of the algorithm and the final output.

The K-means is an algorithm that assumes that the samples belongs to K classes and have a gaussian distribution $P(x | c_i) = \mathcal{N}(x, \mu_i, I)$ with unitary variance matrix (a circle). The algorithm aims (given K) to estimates the K means μ_i .

K-means {

```

let  $m_1 \dots m_K$ , K random points called MEANS
while (termin. condition){}
    define  $X_1 \dots X_K$  as K empty sets
    For  $i=1 \dots K$ {}
         $D_i = \{ \text{point whose closest MEAN is } m_i \}$ 
         $m_i = \frac{1}{|D_i|} \sum_{x \in D_i} x$ 
    }
return  $m_1 \dots m_K$ 
```

