

EXAM 01/07/2019

EXERCISE N°1

POINT 1

The supervised learning problem is the problem to learn a function $h : X \rightarrow Y$ where :

- X be the domain set, which is the set of all possible objects to make predictions about, where a domain point $\vec{x} \in X$ is called instance and is usually represented by a vector of features
- Y be the label set that defines the set of all possible labels
- H be the hypothesis class

The function \hat{h} that we need to pick from H must be the one with the lowest generalization error i.e.

$L_d(\hat{h}) = E_{z \sim D}[l(\hat{h}, z)]$ where :

- D is the unknown probability distribution over Z from which $(x_i, y_i) \in S$ have been drawn (as independent samples).
- $l : H \times Z \rightarrow R^+$ where $Z = X \times Y$ be the loss function namely a function that given an hypothesis provides a measure of how much we lose by predicting the value $h(\vec{x})$ for \vec{x} instead of the

However because we do not know D ,under certain hypotheses, a good estimate of $L_d(h)$ is given by the training error namely : $L_s(h) = \frac{1}{m} \sum_{i=1}^m l(h, (\vec{x}_i, y_i))$ where $S = ((\vec{x}_1, y_1) \dots (\vec{x}_m, y_m))$ is the training set

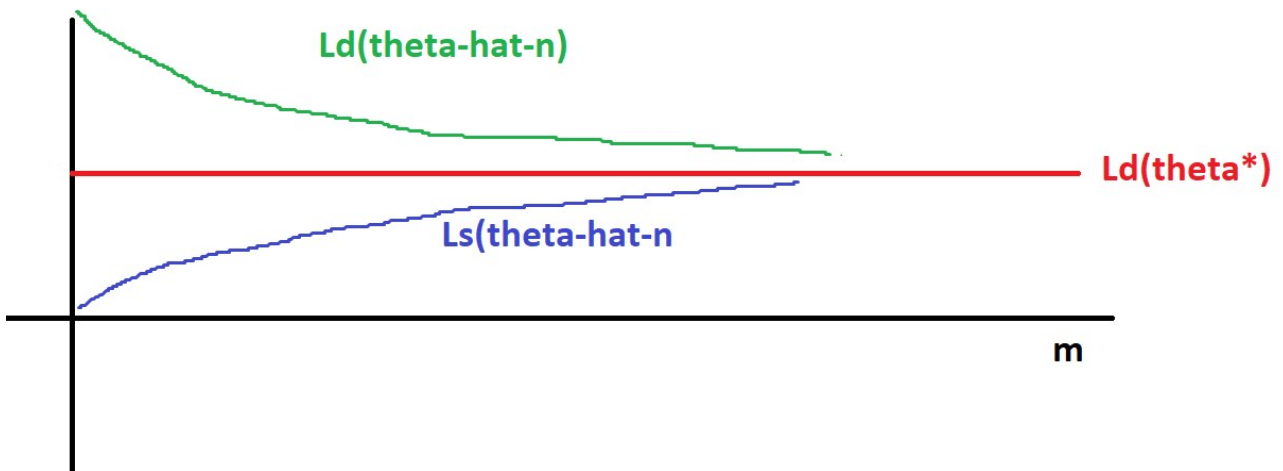
POINT 2

One way to pursue the objective explained in point 1, is ,as already introduce in point 1 as well, to use the so called Empirical Risk Minimization (ERM for short). So given a training set

$S = ((\vec{x}_1, y_1) \dots (\vec{x}_m, y_m))$ with $\vec{x}_i \in X, y_i \in Y \forall i = 1, \dots, m$ to find \hat{h} that has the lowest generalization error we can find the hypothesis that minimizes the training error namely : $L_s(h) = \frac{1}{m} \sum_{i=1}^m l(h, (\vec{x}_i, y_i))$ where $l : H \times Z \rightarrow R^+$ is the loss function .

Notice that this paradigm works only under certain hypotheses regarding the hypothesis class H .

POINT 3

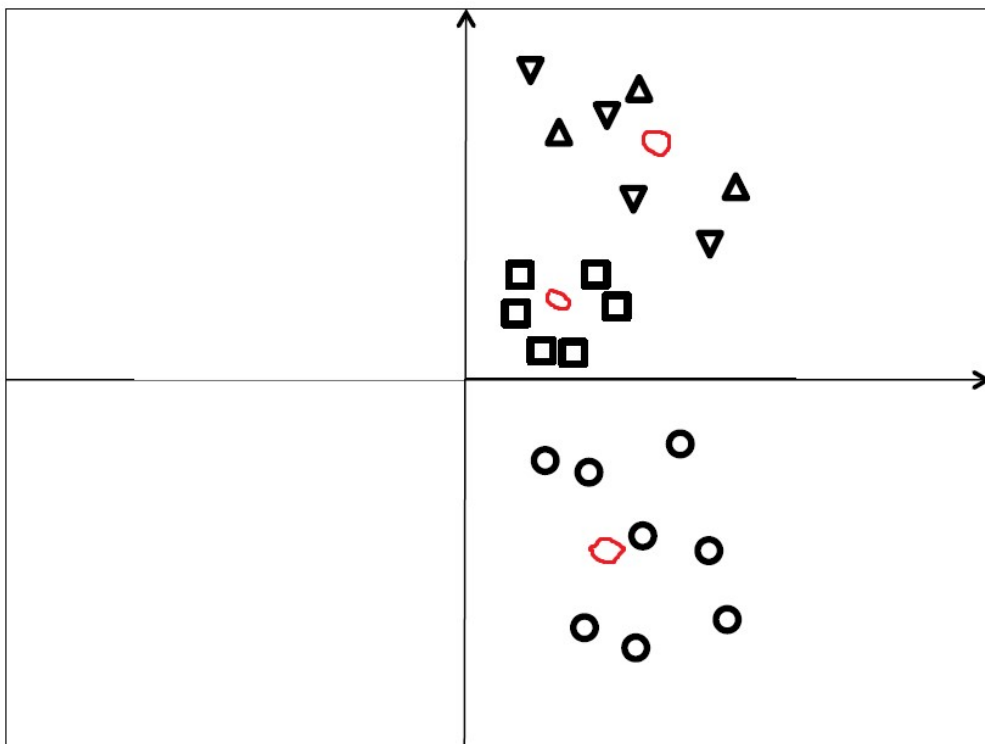


EXERCISE N°2

POINT 1

N°2.1
No there is no way to cluster the to obtain a solution where all the Δ , \square , \circ are
separated because with K-means objective function we minimize the distance between
points and centers.

POINT 2



Rule to obtain such clustering :

Let $\vec{x} \in X$ then :

If $(x_1 < 0 \text{ and } x_2 > 0)$ then $\vec{x}' = [|x_1|, -|x_2|]$

else if $(x_1 < 0 \text{ and } x_2 < 0)$ then $\vec{x}' = [|x_1|, |x_2|]$

else then leave as it is

POINT 3

11.4.3

To identify a good rule for such objective function we must compute the derivative of it with respect to $\vec{\mu}_j$ and then put it equal to $\vec{0}$.

So:

$$\frac{\partial}{\partial \vec{\mu}_j} \left(\sum_{i=1}^K \sum_{\vec{x} \in C_i} d(\vec{x}, \vec{\mu}_i)^3 \right) = \sum_{i=1}^K \frac{\partial}{\partial \vec{\mu}_j} \sum_{\vec{x} \in C_i} d(\vec{x}, \vec{\mu}_i)^3 =$$

$$\frac{\partial}{\partial \vec{\mu}_j} \sum_{\vec{x} \in C_j} d(\vec{x}, \vec{\mu}_j)^3 = \sum_{\vec{x} \in C_j} \frac{\partial}{\partial \vec{\mu}_j} (d(\vec{x}, \vec{\mu}_j)^3) = \sum_{\vec{x} \in C_j} \frac{\partial}{\partial \vec{\mu}_j} \left(\left(\sum_{l=1}^d (x_{lj} - \mu_{lj})^2 \right)^{\frac{3}{2}} \right)$$

$$= \sum_{\vec{x} \in C_j} \frac{\partial}{\partial \vec{\mu}_j} \left(\sum_{l=1}^d (x_{lj} - \mu_{lj})^2 \right) = \sum_{\vec{x} \in C_j} \frac{\partial}{\partial \vec{\mu}_j} d(\vec{x} - \vec{\mu}_j)^2$$

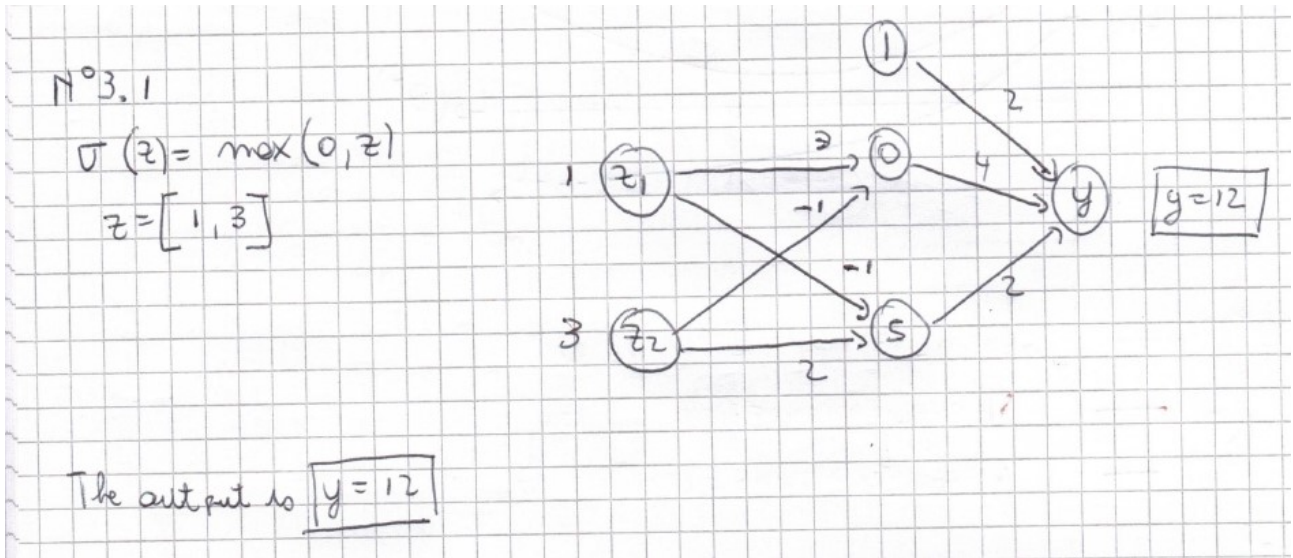
$$= \sum_{\vec{x} \in C_j} -2\vec{x} + 2\vec{\mu}_j = 2|C_j|\vec{\mu}_j - 2 \sum_{\vec{x} \in C_j} \vec{x}$$

At the optimum: $2|C_j|\vec{\mu}_j - 2 \sum_{\vec{x} \in C_j} \vec{x} = 0$

$$\vec{\mu}_j = \frac{1}{|C_j|} \sum_{\vec{x} \in C_j} \vec{x}$$

EXERCISE N°3

POINT 1



POINT 2

The back propagation algorithm is the algorithm to train a neural network based on stochastic gradient descent. The main structure of the algorithm is the following :

Input : Training data $(\vec{x}_1, y_1) \dots (\vec{x}_m, y_m)$ and a neural network with no weights

Output : a neural network with weights $w_{ij}^{(t)} \forall i, j, t$

And it works as follows :

- Randomly initialize the weights $w_{i,j}^{(t)} \forall i, j, t$
- Until convergence is reached then repeat the following steps:
 1. Pick a point \vec{x}_k, y_k randomly from the training data
 2. Apply the forward propagation algorithm for such data and so compute $v_{t,j} \forall j, t$
 3. Compute the sensitivity vectors to make the updates of the weights namely : $\delta_j^{(t)} \forall j, t$
 4. Update the weights with the following rule : $w_{i,j}^{(t+1)} = w_{i,j}^{(t)} - \eta v_{t-1,j} \delta_j^{(t)} \forall i, j, t$
- If convergence is reached then return all the weights $w_{i,j}^{(t)} \forall i, j, t$

POINT 3

The main reasons why the ReLU function is used are the followings :

- It is simpler to compute
- It is as powerful as others
- It does not suffer of the vanishing gradient problem which is common problem in deep learning when using activation functions like the sigmoid or the hyperbolic tangent. Such problem consists of having a roughly 0 value as value of the delta terms to update the weights, on deep levels when applying the back propagation algorithm, and therefore such value will propagates to the shallow levels and so weights are not updated.

EXERCISE N°4

POINT 1

- Classifier C_1 is associated to CURVE A because C_1 is the exact representation of a classifier that suffers of underfitting, namely, the model produced is too vague and has an high training error.
- Classifier C_3 is associated to CURVE C because C_3 is the exact representation of a classifier that suffers of overfitting, namely it is classifying every point of the training set correctly and therefore has a 0 training error at the end of the training iterations.
- Classifier C_2 is associated to CURVE B because it is the perfect balance between a classifier that predicts anything (in the T.S) correct and one that is too vague.

POINT 2

- Classifier C_1 is associated to curve D, namely it is UNDERFITTING both validation and training error are high.
- Classifier C_2 is associated to curve E, namely it has a quite reasonable training and also validation error.
- Classifier C_3 is associated to curve F, namely it is OVERFITTING because it has an high validation error but a small training error.

POINT 3

One strategy to reduce overfitting issue is to use cross-validation to choose the classifier C_i and then, once C_i is chosen we train C_i over all the data. Therefore recalling i, ϑ , we can apply the cross validation method.

To be more precise, cross validation consists in a way to select a value of a parameter ϑ by understanding what is the best value that such parameter can assume with the data we have.

To do so, we split the dataset into k different folds of size m/k (this quantity is supposed to be integer).

Then for each value of the parameter ϑ we find the best model for all the possible $k-1$ folds that we can select. In addition to that, we compute the average error of each parameter as the average of the errors on the folds left out, and when we have repeated such procedure for all the values of the parameters, then we select the optimal one by choosing the one that minimizes the error computed for each parameter. To conclude, we train our model on all the dataset with the parameter found.

The pseudo code is the following :

Input : $S = ((\vec{x}_1, y_1) \dots (\vec{x}_m, y_m))$; set of parameters Θ ; integer k ; learning algorithm A

Split S into S_1, \dots, S_k

foreach $\vartheta \in \Theta$

for $i = 1 \dots k$

$$h_{i,\vartheta} = A(S \setminus S_i; \vartheta)$$

$$error(\vartheta) = \frac{1}{k} \sum_{i=1}^k L_{S_i}(h_{i,\vartheta})$$

Output : $\vartheta^* = \operatorname{argmin}_{\vartheta} (error(\vartheta))$

$$h_{\vartheta^*} = A(S; \vartheta^*)$$

Assuming that we have enough data another strategy is to use validation to choose C_i and then, once C_i is chosen, we use all the data to train C_i .

To be more precise validation consists in a way to select a value of a parameter ϑ by understanding what is the best value that such parameter can assume with the data we have.

To do so we split out dataset into 2 parts training set and validation set. Then, for each value of the parameter that we have, we find the best model for every possible values of the parameter on the training set. Then among such models that correspond to a value of the parameter we select the one that minimizes the validation error. The model obtain then is trained on the entire dataset.

Notice that training error is computed as : $L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, (\vec{x}_i y_i))$ where S is the training set, namely, $S = ((\vec{x}_1, y_1) \dots (\vec{x}_m, y_m))$, $l : H \times Z \rightarrow R^+$ is the loss function where $Z = X \times Y$, that, given an hypothesis provides a measure of how much we lose by predicting the value $h(\vec{x})$ for \vec{x} instead of the correct value y .

The validation error instead is computed as : $L_V(h) = \frac{1}{mv} \sum_{i=1}^{mv} l(h, (\vec{x}_i y_i))$ where V is the validation set, namely, $V = ((\vec{x}_1, y_1) \dots (\vec{x}_{mv}, y_{mv}))$, $l : H \times Z \rightarrow R^+$ is the loss function where $Z = X \times Y$, that, given an hypothesis provides a measure of how much we lose by predicting the value $h(\vec{x})$ for \vec{x} instead of the correct value y .