

EXERCISE 1

Given a classification problem for the function $f : A \times B \times C \rightarrow \{YES, NO\}$, with $A = \{a_1, a_2, a_3\}$, $B = \{b_1, b_2\}$, $C = \{c_1, c_2, c_3\}$ and the following decision tree T that is the result of training on some data set:

1. Provide a rule based representation of the tree T .
2. Determine if the tree T is consistent with the following set of samples
 $S \equiv \{(s_1 = (a_2, b_1, c_2), YES), (s_2 = (a_1, b_1, c_1), NO), (s_3 = (a_1, b_2, c_3), YES), (s_4 = (a_3, b_1, c_1), YES), (s_5 = (a_3, b_2, c_2), NO)\}$. Motivate your answer.
3. Compute the accuracy of T with respect to S .

The rule based repr. is:

$$(B = b_2 \wedge C = c_1) \vee (B = b_1 \wedge A = a_2 \wedge C = c_2) \vee (A = a_3 \wedge C = c_2)$$

The tree is not consistent with S since s_3 have $C = c_3$ and "Yes" as the label, but the tree classifies all the inputs with C_3 as "No".

Except s_3 , the other samples are well classified, the accuracy is $\frac{4}{5} = 80\%$

EXERCISE 2

1. Describe the K-nearest neighbors (K-NN) algorithm for classification and discuss its limitations.

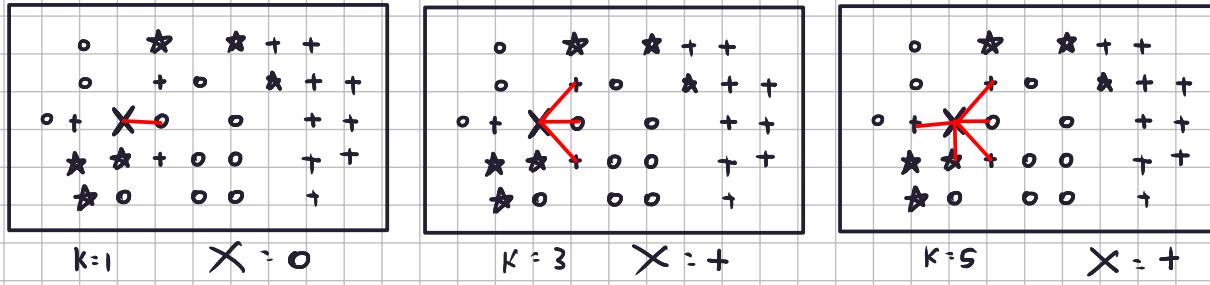
2. Draw a 2D classification dataset for three classes (circle, star, plus), choose one query point close to points in the dataset belonging to different classes, and determine the answer of K-NN for such a query point for $K=1$, $K=3$, and $K=5$, illustrating a situation in which the answer of K-NN depends on K (i.e., you get different solutions for $K=1$, $K=3$, $K=5$). Motivate your answer, showing (with a graphical drawing) which instances contribute to the solution.

The KNN, given a dataset and an input x , classify x with the most popular class among the K closest point to x in D .

$$D_K(x) = \{x_1, \dots, x_K\} \text{ and } \forall x' \notin D_K(x), \exists x_i \in D_K(x) : \|x - x_i\| < \|x - x'\|$$

$$c(x) = \arg \max_C \sum_{x_i \in D_K(x)} \delta(\text{class}(x_i) = C).$$

The main limitation is the fact that we must use and consider all the dataset just to make a single prediction, this is expensive in terms of time and space complexity.



EXERCISE 3

1. Describe the difference between bagging and boosting when combining multiple learners.

2. Consider applying both bagging and boosting to a classification problem and the following models $y_1(x), \dots, y_M(x)$ for bagging and $z_1(x), \dots, z_M(x)$ for boosting as results of the training phase. Write down the formal equations to combine the predictions of the different models for each of the two approaches.

In bagging, we train M models with the same template hypothesis on M different dataset $\{D_i\}$ sampled from the available dataset D .

In boosting, we train sequentially M models on the same dataset, but with weighted samples

-if the model γ_i failed on some points, the model γ_{i+1} will give more importance to that points.

In bagging the final prediction is:

$$\hat{y}(x) = \text{sign} \left(\frac{1}{M} \sum_{i=1}^M y_i(x) \right)$$

In boosting:

$$y(x) = \text{sign} \left(\sum_{i=1}^M \alpha_i z_i(x) \right)$$

Where $\{\alpha_i\}$ are M coefficients used in the training process:

$$J_i = \sum_j^{(i)} w_j (t_j - \hat{y}_i(x_j)) \text{ error for } \hat{y}_i$$

α_j = a measure for the reliability of \hat{y}_i (depends on J_i)

$$w_j^{(i+1)} = w_j^{(i)} \exp(\alpha_j \delta(y_i(x_j) \neq t_j))$$

EXERCISE 4

1. Describe the principle of maximum margin used by SVM classifiers through its formal mathematical definition.
2. Draw a linearly separable dataset for 2D binary classification. Draw a possible solution obtained by SVM and highlights the margin and the support vectors.
3. Discuss why the maximum margin solution is preferred for the classification problem.

Let $D = \{(x_i, t_i)\}_{i=1}^N$ to be the dataset ($t_i \in \{-1, +1\}$). If is linearly separable, then there exists w, w_0 such that $t_i(w^T x_i + w_0) \geq 0 \quad \forall i$. We consider the closest sample to the hyperplane defined by w, w_0 :

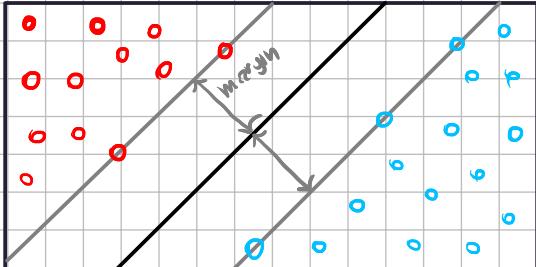
$\min_{x_i} \frac{|y(x_i)|}{\|w\|}$ ← This is the margin. We can set the opt. problem to

Find the w that maximizes that margin:

$$w^*, w_0^* = \underset{w, w_0}{\operatorname{argmax}} \left(\min_{x_i} \frac{|y(x_i)|}{\|w\|} \right) = \underset{w, w_0}{\operatorname{argmax}} \left(\min_{x_i} \frac{t_i(w^T x_i + w_0)}{\|w\|} \right) \begin{array}{l} \text{subj. } t_i y(x_i) \geq 0 \\ \text{to. } \forall i \end{array}$$

We can scale the points s.t. the closest to w, w_0 is at distance 1, so the problem becomes:

$$w^*, w_0^* = \underset{w, w_0}{\operatorname{argmax}} \frac{1}{\|w\|} = \underset{w, w_0}{\operatorname{argmin}} \|w\|^2 \quad \text{s.t. } t_i y(x_i) \geq 1 \quad \forall i$$



The SVM are preferred since it decreases the probability to missclassify new samples (thanks to the margin), respects to other methods (like perceptrons) where the discriminant found as solution is often near to the samples.

EXERCISE 5

Consider a dataset D for the classification problem $f : \mathbb{R}^5 \mapsto \{C_1, C_2, C_3\}$.

1. Describe a probabilistic generative model for such a classification problem, assuming Gaussian distributions.
2. Identify the parameters of the model and determine the size of the model (i.e., the number of independent parameters).

In the probabilistic generative model, we assume that there is a normal distr. on the samples: $P(x|C_i) = N(x; \mu_i, \Sigma)$, and there are the a priori prob. for each classes p_1, p_2, p_3 with $\sum p_i = 1$. So we have $P(x) = \sum_{i=1}^3 p_i N(x; \mu_i, \Sigma)$. The goal is to estimate the parameters of the distributions. For each class we have:

2 mean $\mu_i = 5$ parameter. $p_i = 1$ parameter.
 Then the covariance matrix Σ have 25 prn. } \Rightarrow Total size: 43

EXERCISE 6

Consider a two-layers ANN which receives in input real-valued vectors x of dimension 100 and produces in output real-valued vectors 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.

1. Compute the number of trainable parameters including the bias terms, motivating the answer (i.e., show how to compute this value)
2. Provide the formula explicitly stating how the values of y are computed given an input vector x in terms of the weight matrices and the activation functions (including the bias terms).
3. Provide a suitable loss function for training the network. Motivate your answer, answers with no motivation will not be considered.

We have: rows \downarrow cols
 $W^{(1)} \in \text{Mat}(50 \times 100)$ $b^{(1)} \in \mathbb{R}^{50}$, given the input $x \in \mathbb{R}^{100}$ the output is:
 $W^{(2)} \in \text{Mat}(10 \times 50)$ $b^{(2)} \in \mathbb{R}^{10}$ $e(x) = W^{(2)} \text{ReLU}(W^{(1)}x + b^{(1)}) + b^{(2)}$

So the total number of parameters are: 5560.

The chosen loss Function is the sum of square, since we know that,
 $E = \sum_{i=1}^N \|t_i - e(x_i)\|^2$ For regression, this is equivalent to find the maximum likelihood solution, assuming zero mean Gaussian noise on the samples.