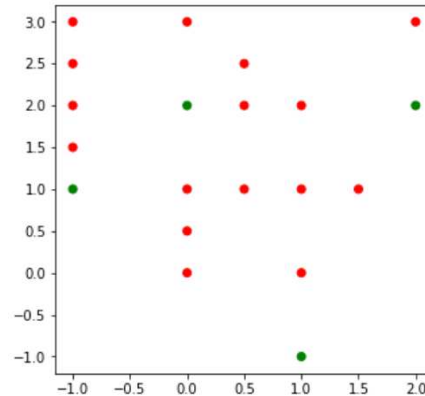


## Final exam. Practical exercises

Total: 5p. Minimum: 1.25p

1. (0.5p) The following dataset is retained in the code variables X and Y.

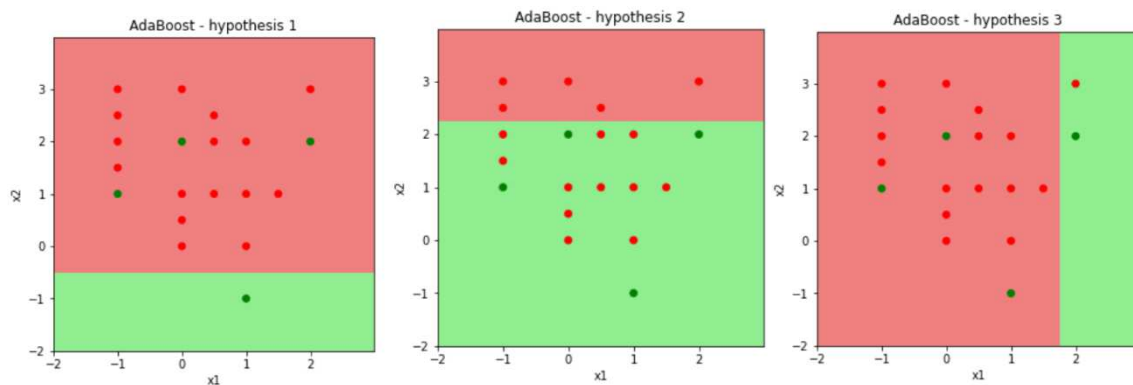


We run the following code

```
# [...] X and Y are initialised as usual with the features
# and the classes, respectively
from sklearn.ensemble import AdaBoostClassifier
ab = AdaBoostClassifier(n_estimators=3, algorithm="SAMME").fit(X, Y)
alphas = ab.estimator_weights_/2
print(alphas)
```

and obtain `[0.86730053 0.36879947 0.44804401]`.

Furthermore, we visualize the decision surfaces of each hypothesis/estimator  $h_t$  in the AdaBoost ensemble:



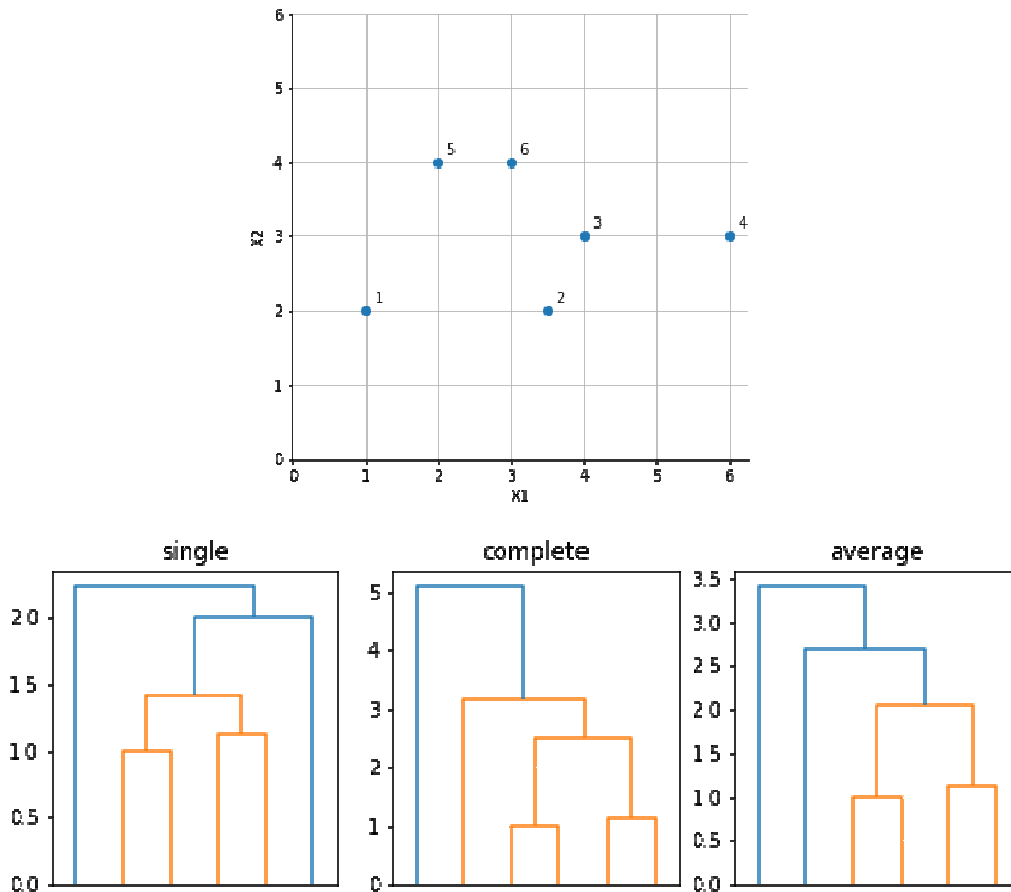
a. What is the output of the following code? The red class is encoded by 1 and the green class by 0. Justify your answer.

```
ab.predict([[2.0, 2.0]])
```

b. Draw the decision boundary of the AdaBoost model/ensemble. You do not need to draw the points. Justify your answer.

2. (0.9p) For the dataset below, the dendrogram for agglomerative clustering has been generated for each type of linkage, but the labels of the dendrogram are missing.

What are the labels for each?



3. (1.2p) What does the following code print? Justify your answer!

```
import pandas as pd
import numpy as np
from sklearn.cluster import KMeans

d = pd.DataFrame({
    'X': [0.0, 4.0, 6.0, 8.0, 10.0],
})

init_centroids = np.array([[ -2.0], [20.0]])
model = KMeans(init=init_centroids, n_clusters=2, max_iter=100)
clusters = model.fit_predict(d)
print(clusters)
```

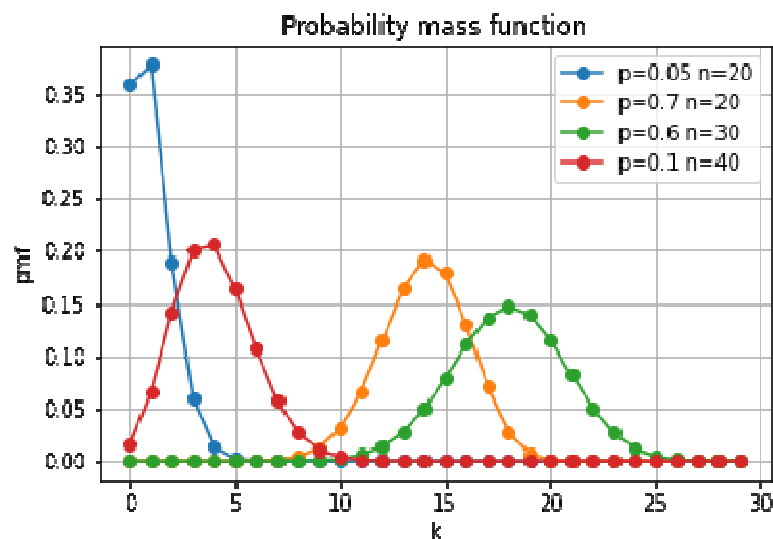
4. (1p) You are using the standard KMeans algorithm to cluster a group of students based on the scores they obtained over four different tests. During the execution of the algorithm, you have the following clusters and centroids. What is the value of the inertia (or J criterion) at this point?

Cluster 1				Cluster 2			
T1	T2	T3	T4	T1	T2	T3	T4
6	9	10	7	8	6	9	9
5	9	9	7	7	7	9	7
7	9	8	8	8	7	9	8

Centroid 1				Centroid 2			
T1	T2	T3	T4	T1	T2	T3	T4
6	9	9	7	8	7	9	8

5. (0.5p) A website has roughly the same number of  $n$  visitors every day and the same conversion rate  $p$ . You learn that this is a binomial distribution of parameters  $p$  and  $n$ , where a success means that a visitor has made a purchase, thus becoming a customer. If the website is seeing 10 customers in one day and 15 in the next, which of the following four sets of parameters is most likely to apply to this website? *(Please include calculations and reasoning!)*



6. (0.9p) Apply one iteration, i.e. the E step and the M step, of the EM algorithm for the mixture of two Bernoulli distributions:

$$\underbrace{\pi \cdot \text{Bernoulli}(\theta_1)}_{\text{Bernoulli}_1} + (1 - \pi) \cdot \underbrace{\text{Bernoulli}(\theta_2)}_{\text{Bernoulli}_2}$$

Pieces of information:

- We are at iteration  $t$ .
- The dataset with which we work is  $D = \{0,0,1\}$ .
- The parameters determined at the previous iteration ( $t-1$ ) are those in the the following mixture:

$$0.6 \cdot \text{Bernoulli}(0.3) + 0.4 \cdot \text{Bernoulli}(0.8)$$

- At the E step you will have to complete the remaining boxes in the following table.

$\gamma_{ij} = \mathbb{E}_{Z X, \pi^{(t-1)}, \theta_1^{(t-1)}, \theta_2^{(t-1)}} [Z_{ij}] = \begin{cases} \mathbb{E}[z_i x_i]^{(t)} & \text{if } j = 1 \\ 1 - \mathbb{E}[z_i x_i]^{(t)} & \text{if } j = 2 \end{cases}$	Bernoulli <sub>1</sub> (j=1)	Bernoulli <sub>2</sub> (j=2)
$x_1 = 0$	0.84	
$x_2 = 0$		
$x_3 = 1$		

- At the M step you will have to **apply** (not to derive) the formulas already known by you. **Do not compute the final result; just put the numbers in the formulas.**

Învățare automată

— Licență, anul III, 2021-2022, examenul parțial II (seria E) —

Nume student:

Grupa:

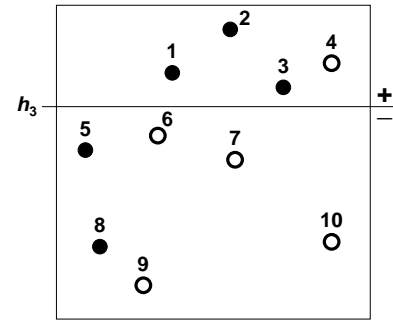
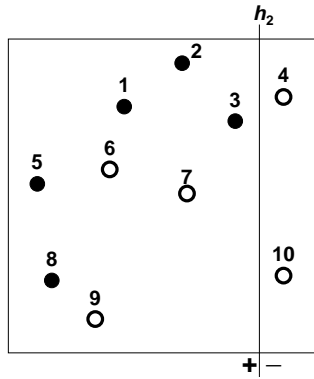
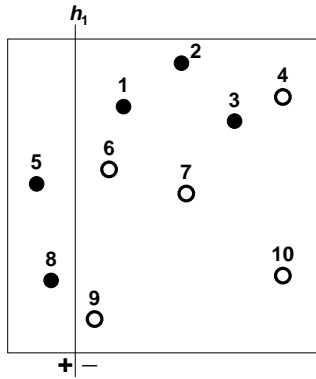
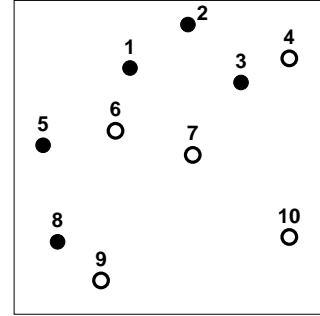
1. [3.5p]

(Algoritmul AdaBoost: aplicare pe un set de date din  $\mathbb{R}^2$ )

Se consideră că aplicăm algoritmul AdaBoost pe dataset-ul din figura alăturată. (Pentru ușurința exprimării la calcule, am notat pe figură indicii instanțelor de antrenament, în imediata apropiere a acestora.)

Folosim convenția noastră obișnuită de notare: simbolul  $\bullet$  desemnează instanțe pozitive, iar simbolul  $\circ$  instanțe negative.

La primele trei iterații ale algoritmului au fost selectați compasii de decizie  $h_1$ ,  $h_2$  și  $h_3$  (în această ordine), așa cum se indică în figurile de mai jos.



*Obiectivul* acestei probleme este să determinăm dacă la sfârșitul celor trei iterații algoritmul AdaBoost reușește să clasifice perfect toate instanțele de antrenament.

a. Calculați

- distribuțiile probabiliste corespunzătoare celor 3 iterații, adică  $D_1(x_i)$ ,  $D_2(x_i)$  și  $D_3(x_i)$ , pentru  $i = 1, \dots, 10$ ;
- pentru fiecare dintre cele 3 iterații ( $t = 1, 2, 3$ ): eroarea ponderată la antrenare ( $\varepsilon_t$ ) produsă de compasul de decizie  $h_t$ , precum și ponderea ( $\alpha_t$ ) asociată ipotezei / compasului de decizie  $h_t$ .

Veți completa tabelele următoare și veți indica succint(!) modul în care ați procedat pentru a ajunge la rezultatele respective.

$i$	1	2	3	4	5	6	7	8	9	10
$D_1(x_i)$										
$D_2(x_i)$	1/6			1/14						
$D_3(x_i)$	7/66			1/22		1/6				

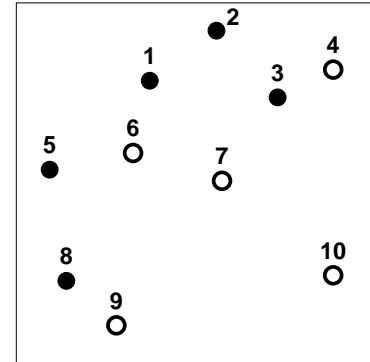
$t$	1	2	3
$\varepsilon_t$			
$\alpha_t$			

**Atenție!** Pentru a vă ușura munca, am completat noi câteva dintre elementele primului tabel. Vă puteți baza pe valorile indicate de noi, ca să vă simplificați raționamentele / calculele!

b. Folosind ipoteza combinată obținută de algoritmul AdaBoost la finalul celei de-a treia iterații, stabiliți

i. eroarea la antrenare produsă (pentru calcularea ei, puteți folosi tabelul de mai jos),

ii. zonele de decizie corespunzătoare acestui clasificator (veți justifica modul în care ați procedat!). Veți indica aceste zone de decizie, precum și granițele de decizie, pe desenul de alăturat.



$t$	$\alpha_t$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$
1	$\alpha_1 = \dots$										
2	$\alpha_2 = \dots$										
3	$\alpha_3 = \dots$										
	$H_3(x_i)$										

**Atenție!** Pentru a vă ușura munca, vă furnizăm noi următoarele valori numerice:  $\ln \sqrt{7/3} \approx 0.423$ ,  $\ln \sqrt{11/3} \approx 0.649$ ,  $\ln \sqrt{15/7} \approx 0.923$ .

2. [2.85p]

(O problemă à la C. Do și S. Batzoglou:  
rezolvarea unei mixturi de doi vectori de distribuții Bernoulli i.i.d.  
identificarea parametrilor și a variabilelor neobservabile;  
scrierea funcției de log-verosimilitate a datelor complete)

Să zicem că avem două monede,  $A$  și  $B$ . Sarcina ta este aceea de a afla (a „estima“)  $\theta_A$  și  $\theta_B$ , probabilitățile de apariție a feței *stemă* (engl., head) pentru fiecare dintre cele două monede. Însă eu sunt cam răutacios și nu-ți dau voie să arunci tu monedele. În schimb, decid să procedez astfel: voi arunca eu însumi monedele și după aceea îți voi comunica ție rezultatele aruncărilor. În mod concret, îți voi spune ceva de forma următoare: am ales una dintre cele două monede (nu-ți spun care anume),<sup>1</sup> am aruncat-o de 10 ori și am obținut în total de 7 ori stema și de 3 ori banul. Apoi am ales din nou una dintre cele două monede (poate aceeași cu cea dinainte, poate nu), am aruncat-o de 10 ori și am obținut în total de 5 ori stema și de 5 ori banul. În total, îți comunic de  $N$  ori câte o astfel de informație. (Așadar, la final vei dispune de rezultatele a  $10N$  aruncări ale monedelor.)

Facem *presupunerea* că la fiecare dintre aceste  $N$  serii cele 10 aruncări sunt independente unele de altele și că am reținut ordinea / succesiunea rezultatelor obținute în urma acestor aruncări și ți-o comunic.

a. Formulează aceasta ca pe o problemă de tip EM (Expectation-Maximization).

i. Care sunt datele observabile?

ii. Care sunt variabilele neobservabile / ascunse / latente (engl., hidden variables)?

iii. Care sunt parametrii modelului (adică parametrii distribuțiilor probabiliste folosite)?

b. Calculează [adică, stabilește pas cu pas care este] expresia funcției de *log-verosimilitate a datelor complete* la iterația  $t$ , pentru această problemă.

(Atenție! Nu ți se cere să rezolvi efectiv problema EM.)

---

<sup>1</sup>Vom considera că moneda  $A$  este aleasă întotdeauna cu probabilitatea  $\pi$ , iar moneda  $B$  cu probabilitatea  $1 - \pi$ .