

Demande 5

Risposta non
ancora data

Punteggio max:
6,00

P Contrassegna
domanda

ASSOCIATION RULES

Given the list of transactions reported below compute list the closed itemset using a support threshold of 0.4 (that is 40%). $\Rightarrow \text{min sup} = 2$

- a1,a2,a4,a5
 - a2,a3,a4,a5
 - a3,a5
 - a2
 - a2,a3,a4

List all the closed itemsets using the template included in the answer box. Each closed itemset must be reported on a separate line with its support.

Closed Itemset #1: (write the closed itemset as a1, a2, a3, ...)

Support #1:

Closed Itemset #2: (write the closed itemset as a1, a2, a3, ...)

Support #2:

Closed Itemset #3: (write the closed itemset as a1, a2, a3, ...)

Support #3:

Closed Itemset #3: (write the closed itemset as a1, a2, a3, ...)

Support #3:

Closed Itemset #4: (write the closed itemset as a1, a2, a3, ...)

Support #4:

Closed Itemset #5: (write the closed itemset as a1, a2, a3, ...)

Support #5:

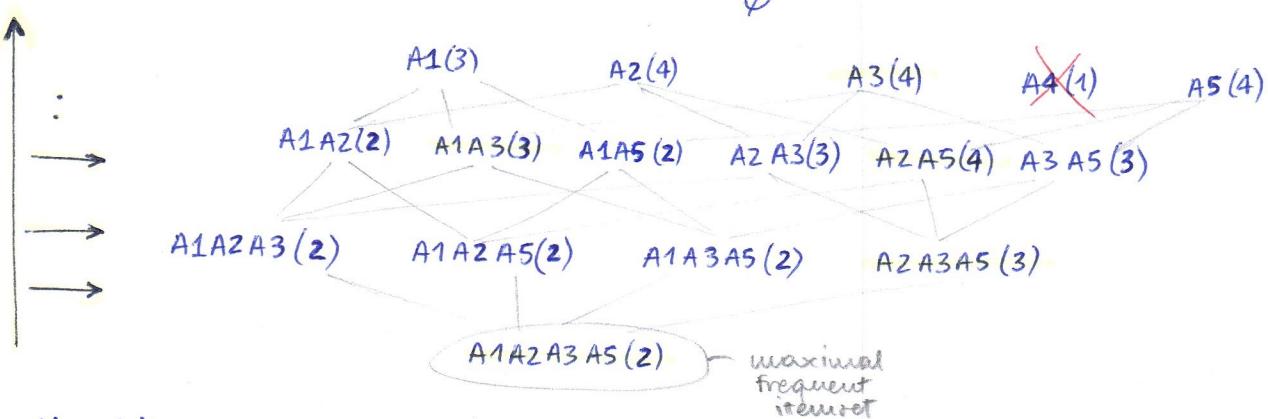
(insert additional entry if needed)

Version 2 :	1	A1	A3	A4	
	2	A2	A3	A5	
	3	A1	A2	A3	A5
	4	A2	A5		
	5	A1	A2	A3	A5

min dup = 2

$$\min \sup = 2$$

How we proceed
the search:



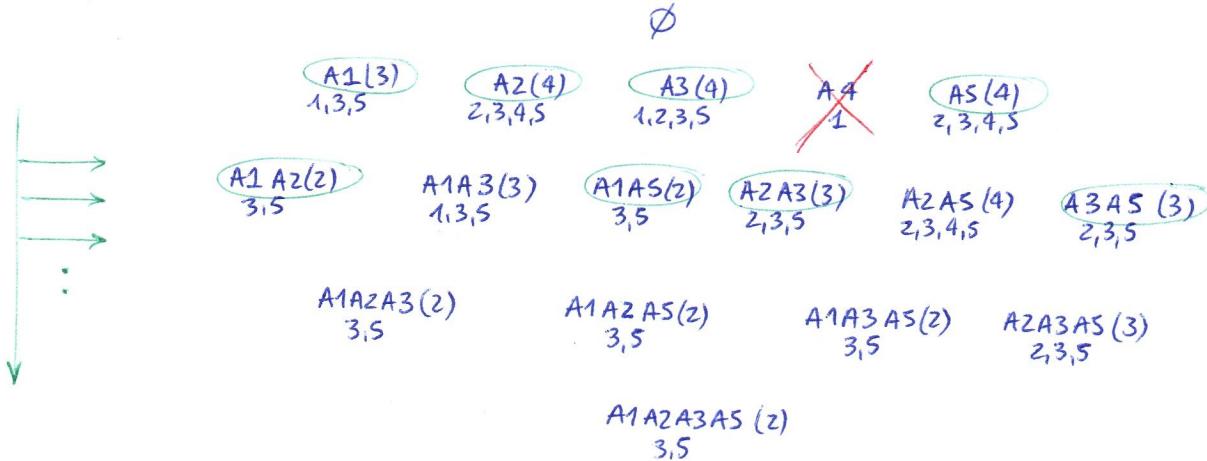
Closed itemset: itemset s.t. all its supersets have less support
We start from the bottom and we consider one level at the time.
→ Closed itemsets in yellow

What if we want to apply Eclat?

Translation of the database:

A1	A2	A3	A4	A5
1	2	1	1	2
3	3	2		3
5	4	3		4
5	5			5

And then:



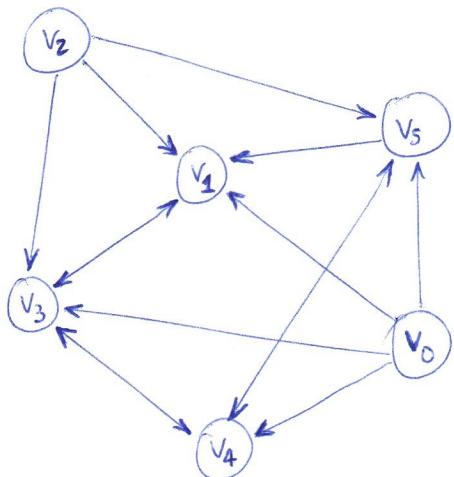
What if we look for minimal generators?

(itemsets that don't have
subsets with the same support)

It's convenient to start from the top and consider one level at the time.

TRAWLING

Consider the graph below, identify the communities using a support threshold of 3.



$$V_0 = \{V_1, V_3, V_4, V_5\}$$

$$V_1 = \{V_3\}$$

$$V_2 = \{V_1, V_3, V_5\}$$

$$V_3 = \{V_1, V_4\}$$

$$V_4 = \{V_3, V_5\}$$

$$V_5 = \{V_1, V_4\}$$

Frequent itemsets (a priori):

\emptyset

$V_2(4)$

~~$V_2(0)$~~

$V_3(4)$

$V_4(3)$

$V_5(3)$

~~$V_2V_3(2)$~~

$V_1V_4(3)$

~~$V_2V_5(2)$~~

~~$V_3V_4(2)$~~

$V_3V_5(3)$

~~$V_4V_5(1)$~~

\Rightarrow

V_1 $s=4$

V_2V_4 $s=3$

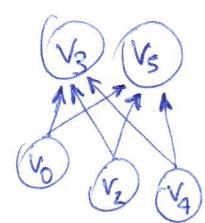
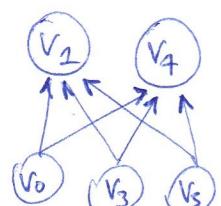
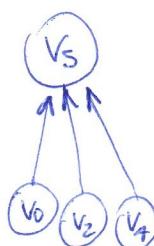
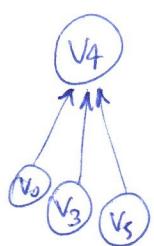
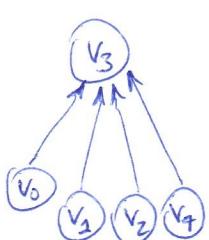
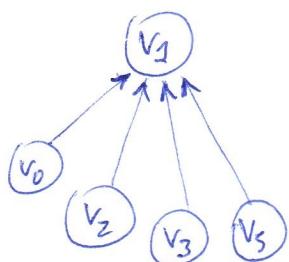
V_3 $s=4$

V_3V_5 $s=3$

V_4 $s=3$

V_5 $s=3$

Communities :



24/06/2020 - CLUSTERING

Consider the 5 datapoints p_0, \dots, p_4 and the following distance matrix:

	p_0	p_1	p_2	p_3	p_4
p_0	0.0	4	10	4	1
p_1	4	0.0	6	2	5
p_2	10	6	0.0	6	11
p_3	4	2	6	0.0	5
p_4	1	5	11	5	0.0

Apply the first step of hierarchical clustering using single link (MIN) approach to measure distance between clusters.

- Report the points involved in the first merge and the corresponding distances
- Report the resulting matrix showing one row after another

Let's do it completely!

p_0	-			
p_1	4	-		
p_2	10	6	-	
p_3	4	2	6	-
p_4	1	5	11	5
	p_0	p_1	p_2	p_3



Note: to determine the first merge ($\{p_0, p_4\}$) we consider the points with smaller distance (1 linkage).
first merge: p_0 and p_4 ($d=1$)
we use: single linkage

p_0, p_4	-		
p_1	4	-	
p_2	10	6	-
p_3	4	2	6
	p_0, p_4	p_1	p_2

Remember: the matrix distance decides which elements are going to merge, the linkage decides how we'll update the matrix distance

$$d(\{p_0, p_4\}, p_1) = \min \{ d(p_0, p_1), d(p_4, p_1) \} \\ \vdash \min \{ 4, 5 \} = 4$$

$$d(\{p_0, p_4\}, p_2) = \min \{ 10, 11 \} = 10$$

$$d(\{p_0, p_4\}, p_3) = \min \{ 4, 5 \} = 4$$

Second merge: p_1 and p_3 ($d=2$)

p_0, p_4	-		
p_2, p_3	4	-	
p_2	10	6	-
	p_0, p_4	p_2, p_3	p_2



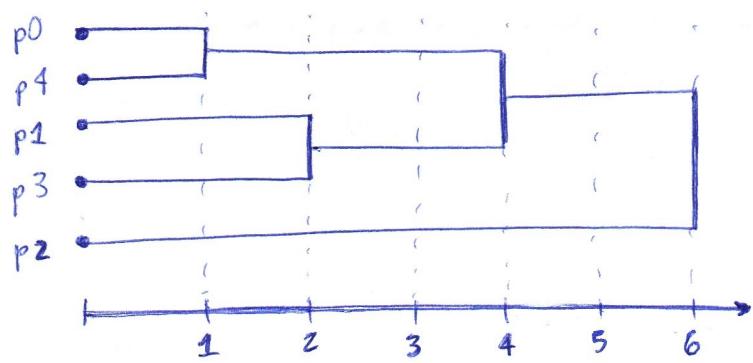
$$d(\{p_1, p_3\}, p_2) = \min \{ 6, 6 \} = 6 \\ d(\{p_0, p_4\}, \{p_1, p_3\}) = \\ = \min \{ d(\{p_0, p_4\}, p_1), d(\{p_0, p_4\}, p_3) \} \\ = \min \{ 4, 4 \} = 4$$

p_0, p_1, p_3, p_4	-	
p_2	6	-
	p_0, p_1, p_3, p_4	p_2

$$d(\{ \cdot \cdot \cdot \}, p_2) = \min \{ 10, 6 \} = 6$$

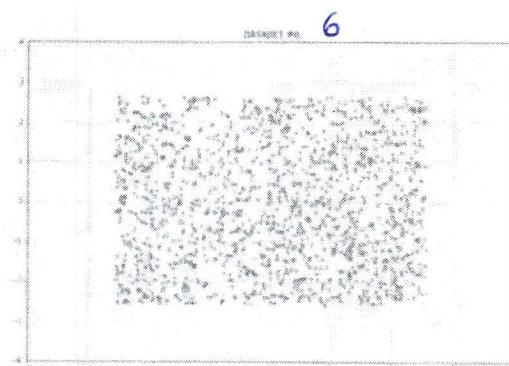
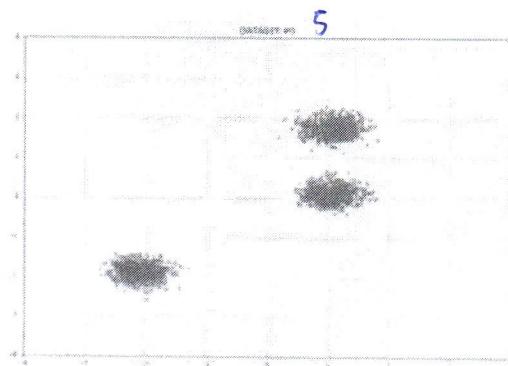
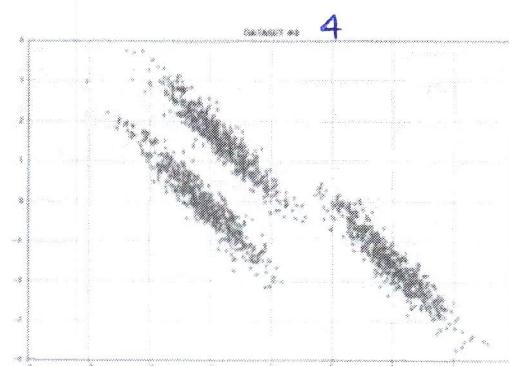
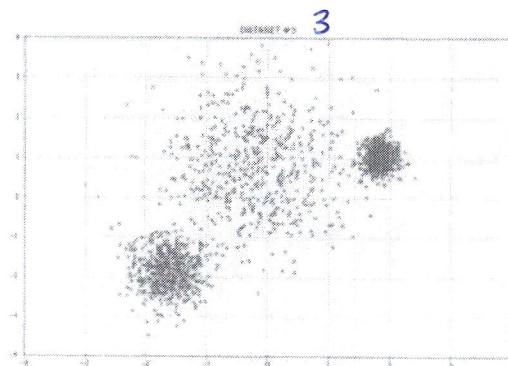
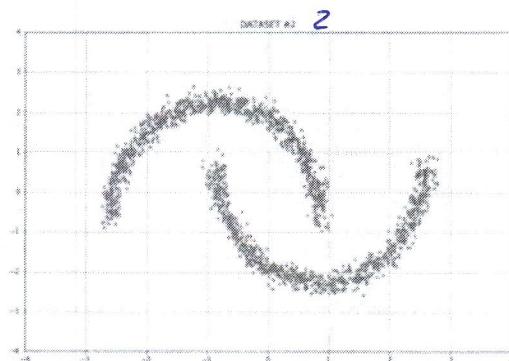
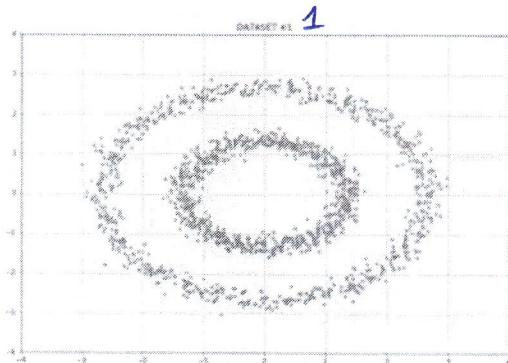
The final merge is at distance 6.

Dendrogram



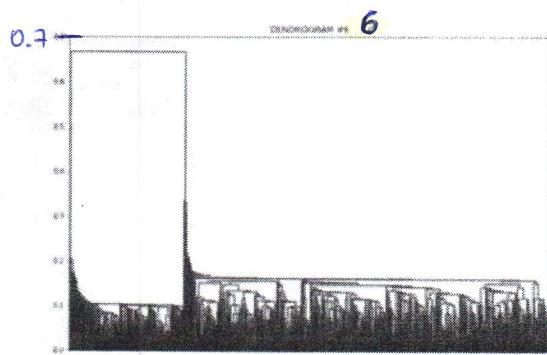
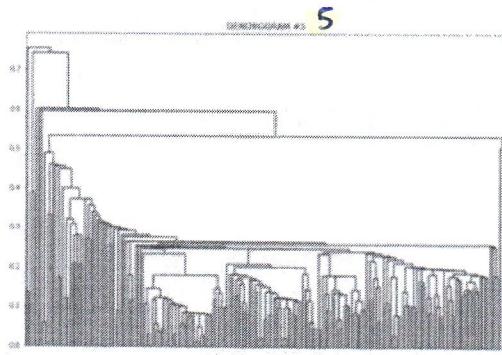
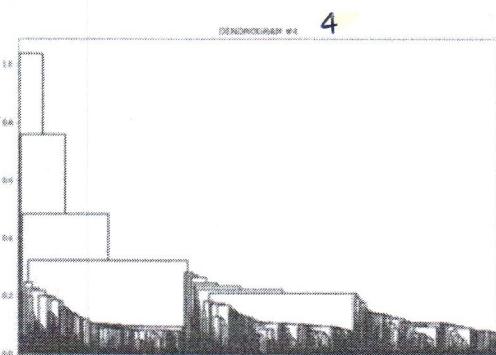
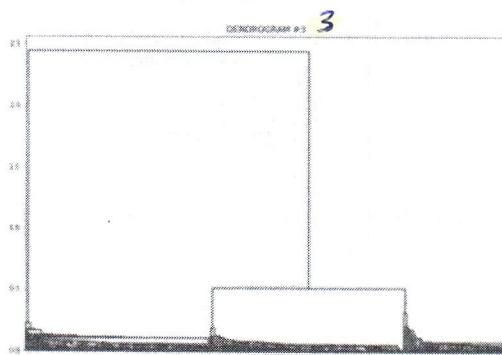
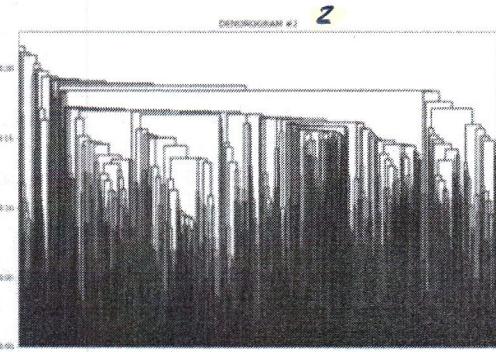
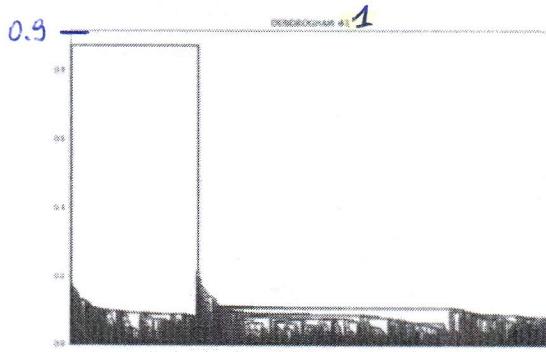
CLUSTERING

14/07/2020 (pt. 1) - Single linkage



CLUSTERING

14/07/2021 (pt. 2)

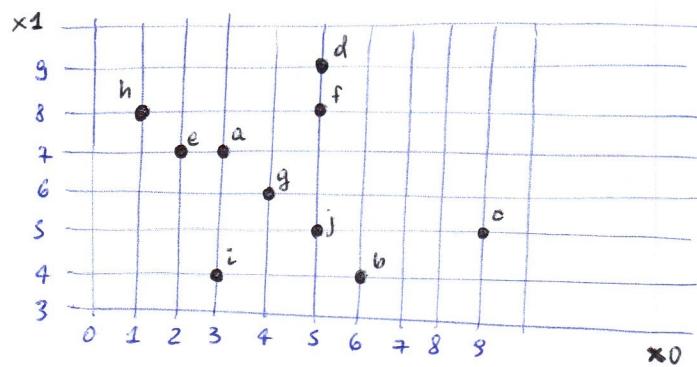


Dendrogram 3 comes from dataset 5
 2
 5
 4
 6
 1

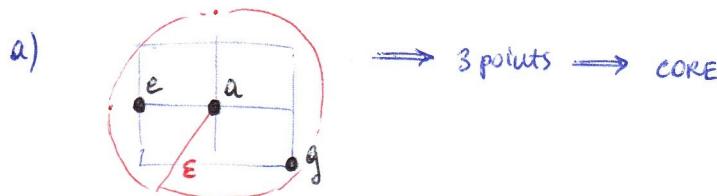
because data are more dense (compact)
 (the points are more closer among each other
 than in the dataset 1)

REPRESENTATIVE-BASED CLUSTERING

Given the dataset reported below, apply DBSCAN with $\epsilon = 1.5$ and $\text{minpts} = 3$ and list all the core, border and noise points.



Core points = points that have at least minpoints in the way of ϵ
The point itself IS PART OF THE COUNT!

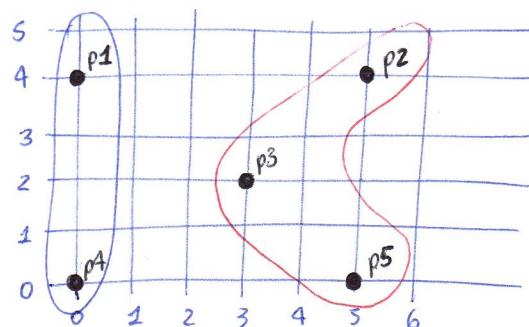


- b) 2 points \Rightarrow BORDER (j) CORE
- c) 1 point \Rightarrow NOISE
- d) 2 points \Rightarrow NOISE
- e) 3 points \Rightarrow CORE
- f) 2 points \Rightarrow NOISE
- j) 3 points \Rightarrow CORE
- h) 2 points \Rightarrow BORDER
- g) 3 points \Rightarrow CORE
- i) 1 point \Rightarrow NOISE

(Taken from 5/02/2020) - REPRESENTATIVE BASED CLUSTERING

Consider the following dataset containing 5 points (p_1, \dots, p_5) in which columns x and y represent the input variables while the column "label" contains the cluster label that has been assigned to each example by applying k-means with $k=2$ and Euclidean distance. Compute for each example the silhouette coefficient using Euclidean distance.

	x	y	label
p_1	0	4	0
p_2	5	4	1
p_3	3	2	1
p_4	0	0	0
p_5	5	0	1



For every point we need a, b , silhouette:

a = average of the distances with points inside the cluster

b = average of the distances with points of the closest cluster

$$s_i = \frac{b - a}{\max\{b, a\}}$$

	a	b	s_i
p_1	4	5	0.2
p_2	3.4	5.7	0.4
p_3	2.83	3.6	0.21
p_4	4	5	0.2
p_5	3.4	5.7	0.4

$$S_{TOT} = \frac{1}{n} \sum_{i=1}^n s_i$$

$$= 0.282$$

$$p_1) \quad a = 4 = d(p_1, p_4)$$

$$p_2) \quad a = (d(p_2, p_3) + d(p_2, p_5))/2 = (2\sqrt{2} + 4)/2 = \sqrt{2} + 2 = 3.4142 = 3.4$$

$$p_3) \quad a = (d(p_3, p_2) + d(p_3, p_5))/2 = (2\sqrt{2} + 2\sqrt{2})/2 = 2\sqrt{2} = 2.8284 = 2.83$$

$$p_4) \quad = p_1)$$

$$p_5) \quad = p_2)$$

$$p_1) \quad b = (d(p_1, p_2) + d(p_1, p_3) + d(p_1, p_5))/3 = (5 + \sqrt{13} + \sqrt{41})/3 = (5 + 3.6 + 6.4)/3 = 5$$

$$p_2) \quad b = (d(p_2, p_1) + d(p_2, p_4))/2 = (5 + \sqrt{41})/2 = (5 + 6.4)/2 = 5.7$$

$$p_3) \quad b = (d(p_3, p_1) + d(p_3, p_4))/2 = \sqrt{13} = 3.6$$

$$p_4) \quad = p_1)$$

$$p_5) \quad = p_2)$$

$$p_1) \quad (5-4)/5 = s_1 = 1/5 = 0.2$$

$$p_2) \quad s_2 = (5.7 - 3.4)/5.7 = 0.4$$

$$p_3) \quad s_3 = (3.6 - 2.83)/3.6 = 0.21$$

$$p_4) \quad = p_1)$$

$$p_5) \quad = p_2)$$

July 2020 - NAIVE BAYES & BAYESIAN NETWORKS

Consider the dataset below where attribute "class" identify the class attribute. Compute the probability for the class and for the attribute "meat" that are required to compute a Naive Bayes classifier using the Laplace Estimator and report the results in the corresponding boxes below. "cooking" and "side" have been masked using "-" since they're useless for answering the questions below.

meat	cooking	side	class
chicken	fried	fried	dislike
lamb	wanted	veggie	like
veal	-	-	dislike
chicken	-	bacon	like
chicken	-	-	like
veal	-	-	dislike
lamb	-	-	dislike
veal	-	-	dislike
veal	-	-	like
pork	-	-	like
veal	-	-	dislike
lamb	-	-	dislike

Part 1.

1. $P(\text{class}=\text{dislike}) = 7/12 = 0.58$
2. $P(\text{class}=\text{like}) = 5/12 = 0.42$
3. $P(\text{meat}=\text{chicken} | \text{class}=\text{dislike}) = 0.18$
4. $P(\text{lamb} | \text{dislike}) = 0.27$
5. $P(\text{pork} | \text{dislike}) = 0.09$
6. $P(\text{veal} | \text{dislike}) = 0.45$
7. $P(\text{chicken} | \text{like})$
8. $P(\text{lamb} | \text{like})$
9. $P(\text{pork} | \text{like})$
10. $P(\text{veal} | \text{like})$

Meat	like	dislike
chicken	2+1	1+1
lamb	1+1	2+1
veal	1+1	4+1
pork	1+1	0+1
	9	11

Laplace
Estimator

Probabilities		
Meat	like	dislike
chicken	3/9	2/11
lamb	2/9	3/11
veal	2/9	5/11
pork	2/9	1/11

Part 2.

Consider the values of the other probabilities required by the Naive Bayes classifier listed below:

- $P(\text{cooking} = \text{fried} | \text{dislike}) = 0.56$
- $P(\text{roasted} | \text{dislike}) = 0.44$
- $P(\text{fried} | \text{like}) = 0.43$
- $P(\text{roasted} | \text{like}) = 0.57$
- $P(\text{side} = \text{baked} | \text{dislike}) = 0.50$
- $P(\text{fried} | \text{dislike}) = 0.30$
- $P(\text{veggie} | \text{dislike}) = 0.20$
- $P(\text{baked} | \text{like}) = 0.38$
- $P(\text{fried} | \text{like}) = 0.25$
- $P(\text{veggie} | \text{like}) = 0.38$

Compute the probability for Naive Bayes classifier to classify the following 3 examples in the dataset as dislike or like:

1. $P(\text{dislike} | \text{pork, roasted, fried})$
2. $P(\text{like} | \text{pork, roasted, fried})$
3. $P(\text{dislike} | \text{veal, fried, baked})$
4. $P(\text{like} | \text{veal, fried, baked})$
5. $P(\text{dislike} | \text{lamb, fried, baked})$
6. $P(\text{like} | \text{lamb, fried, baked})$

$$1. P(\text{dislike} | \text{pork, roasted, fried}) = P(\text{pork} | \text{dislike}) P(\text{roasted} | \text{dislike}) P(\text{fried} | \text{dislike}) P(\text{dislike}) \\ \downarrow (1/11) \cdot (0.44) \cdot (0.30) \cdot (0.58) = 0,01$$

this still need to be normalized to be a TRUE probability !

$$2. P(\text{like} | \text{pork, roasted, fried}) = \dots = "val"$$

→ the results of the questions must be normalized! "val" + 0,01 ≠ 1, they have to!
→ normalization!

June 2020 - DECISION TREE

Given the following dataset in which "class" represents the class attribute, compute the information gain for attribute "meat".

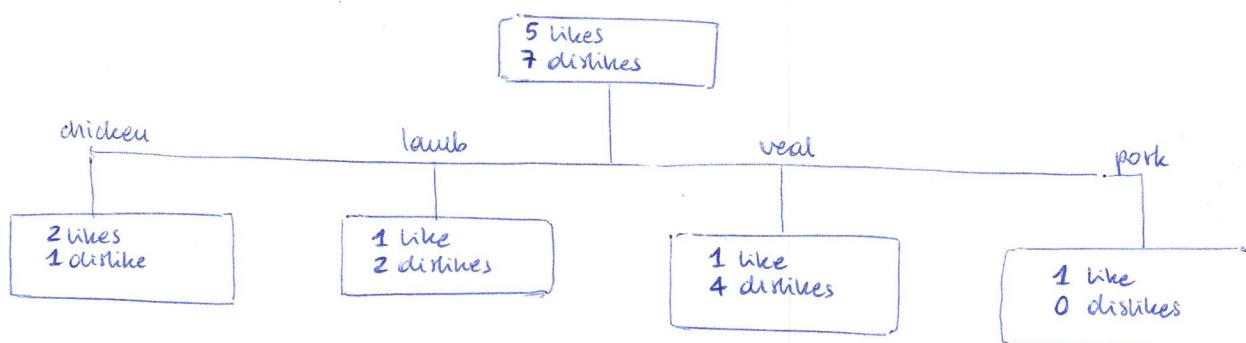
NOTE: the logarithm computation must be performed using base 2. The result must be reported using at least a two digit precision.

	meat	cooking	side	spicyness	class
0	chicken	fried	fried	1.0	dislike
1	lamb	roasted	veggie	1.0	like
2	veal	roasted	veggie	1.0	dislike
3	chicken	fried	baked	1.0	like
4	chicken	fried	baked	2.0	like
5	veal	roasted	fried	2.0	dislike
6	lamb	fried	baked	2.0	dislike
7	veal	roasted	baked	2.0	dislike
8	veal	roasted	veggie	3.0	like
9	pork	roasted	fried	3.0	like
10	veal	fried	baked	4.0	dislike
11	lamb	fried	baked	4.0	dislike

$$\text{Entropy Before Splitting} = 0.98$$

$$\text{Entropy After Splitting} = 0.76$$

$$\text{Information Gain} = 0.22$$



Root:

$$\begin{aligned} \text{info}([5, 7]) &= \text{entropy}([5/12, 7/12]) = -\frac{5}{12} \log_2(5/12) - \frac{7}{12} \log_2(7/12) \\ &= -0.417 \cdot (-1.263) - 0.583 \cdot (-0.778) = 0.98 \end{aligned}$$

chicken:

$$\text{info}([2, 1]) = \text{entropy}([2/3, 1/3]) = -\frac{2}{3} \log_2(2/3) - \frac{1}{3} \log_2(1/3) = 0.918$$

lamb:

$$\text{info}([1, 2]) = \text{info}([2, 1]) = 0.918$$

Pork:

$$\text{info}([1, 0]) = 0 \quad (\text{0 confusion (entropy) when we have 100% of one class})$$

Veal:

$$\text{info}([1, 4]) = \text{entropy}([1/5, 4/5]) = -\frac{1}{5} \log_2(1/5) - \frac{4}{5} \log_2(4/5) = 0.722$$

The overall entropy after the split:

$$= \frac{3}{12} \cdot 0.918 + \frac{3}{12} \cdot 0.918 + \frac{5}{12} \cdot 0.722 + \frac{1}{12} \cdot 0$$

$$= 0.76 = \text{info(meat)}$$

$$\text{gain(meat)} = \text{info}(D) - \text{info(meat}) = 0.98 - 0.76 = 0.22$$

July 2020 - DECISION TREE

Given the following dataset in which "class" represents the class attribute, compute the Information Gain for the numerical attribute "spiciness".

NOTE:

- The split value is computed as the average of the nearby values for example, the split between 0.7 and 0.8 will be 0.75.
- Entropy must be computed using the logarithm with base 2.
- The result must be reported using at least a two digit precision.

meat	cooking	side	spiciness	class
pork	roasted	fried	1.0	dislike
veal	roasted	veggie	1.0	like
chicken	fried	fried	2.0	dislike
pork	fried	baked	3.0	like
chicken	fried	fried	3.0	dislike
veal	fried	baked	3.0	dislike
pork	roasted	baked	3.0	like
chicken	fried	baked	4.0	dislike
chicken	roasted	fried	4.0	like
chicken	roasted	fried	4.0	like
pork	fried	veggie	4.0	like
veal	fried	baked	4.0	dislike

1. Entropy Before Splitting = 1

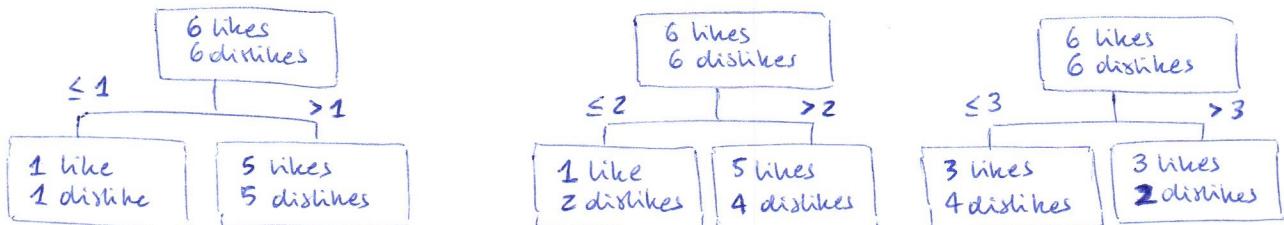
2. Splitting Threshold =

3. Information Gain =

Possible splits:

1.0	1.0	2.0	3.0	3.0	3.0	3.0	4.0	4.0	4.0	4.0	4.0	4.0
dislike	like	dislike	like	dislike	dislike	like	dislike	like	like	like	like	dislike

$$\text{Entropy}(D) = \text{info}(D) = \text{info}([6,6]) = \text{entropy}([\frac{1}{2}, \frac{1}{2}]) = 1 \quad \text{maximal confusion (50%, 50%)}$$



Split 1:

(≤ 1) we got 50%, 50%, maximum confusion \Rightarrow entropy, info = 1

(> 1) entropy = 1

$$\text{info}(\text{split } 1) = \frac{2}{12} \cdot 1 + \frac{10}{12} \cdot 1 = 1 \quad \Rightarrow \text{no gain}$$

Split 2:

$$\text{info}(\leq 2) = \text{info}([1,2]) = \text{entropy}([\frac{1}{3}, \frac{2}{3}]) = -\frac{1}{3} \log_2(\frac{1}{3}) - \frac{2}{3} \log_2(\frac{2}{3}) = 0.918$$

$$\text{info}(> 2) = \text{info}([5,7]) = \text{entropy}([\frac{5}{12}, \frac{7}{12}]) = -\frac{5}{12} \log_2(\frac{5}{12}) - \frac{7}{12} \log_2(\frac{7}{12}) = 0.991$$

$$\text{info}(\text{split } 2) = \frac{3}{12} \cdot 0.918 + \frac{9}{12} \cdot 0.991 = 0.973$$

$$\text{gain}(\text{spiciness}, \leq 2, > 2) = \text{info}(D) - \text{info}(\text{spiciness}, \leq 2, > 2) = 1 - 0.973 = 0.027$$

Split 3 :

$$\text{info}(\leq 3) = \text{info}([3, 4]) = \text{entropy}([3/7, 4/7]) = [\dots]$$

$$\text{info}(> 3) = \text{info}([3, 2]) = \text{entropy}([\dots]) = [\dots]$$

$$\text{info}(\text{split } 3) = [\dots]$$

$$\text{grain}(\text{spiciness}, \leq 3, > 3) = [\dots]$$