

Where

Date

Who

Why

- ① You are given the ground truth labels of the data set and the predicted values obtained with a clustering algorithm

$$\text{y-true} = [1, 1, 2, 2, 3, 3]$$

$$\text{y-pred} = [2, 2, 1, 1, 1, 1]$$

compute the Rand Index [RI] between y-true and y-pred

Remember that the Rand Index is given by

$$RI = \frac{a+b}{DP} = \frac{TP + TN}{DP}$$

where:

$TP$  (a) is the number of pairs of elements (without ordering) that are in the same set in y-true and in the same set in y-pred.

$TN$  (b) is the number of pairs of elements (without ordering) that are in the difference set in y-true and in different y-pred.

(DP): is the total number of possible pairs in the Dataset (without ordering)

			TP	TN
AB	1 1	2 2	✓	✗
AC	1 2	2 1	✓	
AD	1 2	2 1	✓	
AE	1 3	2 1	✓	
AF	1 3	2 1	✓	
BC	1 2	2 1	✓	
BD	1 2	2 1	✓	
BE	1 3	2 1	✓	
BF	1 3	2 1	✓	
CD	2 2	1 1	✓	✗
CE	2 3	1 1		✗
CF	2 3	1 1		✗
DE	2 3	1 1		✗
DF	2 3	1 1		✗
EF	3 3	1 1	✓	✗
			3	8

$$TP + TN = \frac{3+8}{15} = \frac{11}{15} = .733$$

$$\binom{6}{2} = \frac{6!}{2!4!} = \frac{6 \times 5 \times 4!}{2!4!} = 15$$

Where Date Who Why

② You are given a dataset containing 300,000 records divided into 3 classes with the following schema.

Id | value\_0 | ... | value\_100 | <sup>Label</sup>  
<sub>Label</sub>

You have to build a classification model to avoid overfitting which is the best strategy to train and test the classifier?

a)  $x_{\text{train}}, x_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train\_test\_split}(x, y, \text{testSize} = 0.2)$

b)  $\text{tscv} = \text{timeSeriesSplit}(n\_splits=4)$

for train, test in tscv.split(x)

c)  $x_{\text{train}}, x_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train\_test\_Split}(x, y, \text{testSize} = 0.2)$

KF =  $Kfold(n\_split = 4)$

for train, test in kf.split(x\_train):

d)  $loo = \text{leaveOneOut}()$

for train\_index, test\_index in loo.split(x)

Where

Date

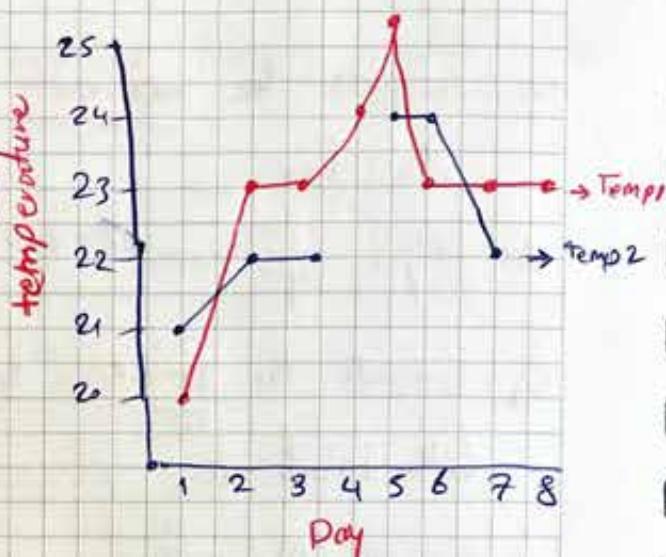
Who

Why

- ③ You have to compare the daily temperature differences between two cities, sample by sample the trends are showed in the picture below.

Question 1: which kind of preprocessing would you choose?

- Consider the difference between the series when missing values
- Fill missing values with  $\infty$
- Fill missing values with previous value
- Add zero padding after temperature<sup>2</sup>



Question 2: Given the strategy that you have selected, calculate the mean of daily differences calculated with the following numpy function.

$$\text{avg\_daily\_diff} = \text{np.mean}(\text{temp1} - \text{temp2})$$

Day	1	2	3	4	5	6	7	8
temp1	20	23	23	24	25	23	23	23
temp2	21	22	22	22	24	24	22	22
diff	-1	+1	+1	+2	+1	-1	+1	+1

$\Rightarrow 5 \Rightarrow \text{mean} = \frac{5}{8} = .625$

Where

Date

Who

Why

- ④ You are given the confusion matrix in the figure below obtained after the validation phase of classifier

$Q_1$  = Compute the accuracy of the classifier

$$\text{Accuracy: } \frac{30}{43} = 0.6977$$

$Q_2$  = Compute the precision for class with Label A.

$$\text{precision (A)} : \frac{15}{25} \rightarrow \text{predicted as class "A"} \Rightarrow \frac{15}{25} = 0.6$$

$$\text{Recall (A)} = \frac{15}{18} \rightarrow \text{actually is belong class A}$$

$$F = \frac{2PR}{P+R} \Rightarrow$$

- ⑤ You are given a decision tree trained on a dataset with categorical features only.

which statement is true:

- Ⓐ each internal node of a decision tree can have at most " $C$ " children, where " $C$ " is equal to the number of distinct values of the attribute associated with the internal node.

		A	B	
		15	3	18
True label	A	10	15	25
	B			18
		25	18	

predicted Label

?

Where

Date

Who

Why

- ⑥ The matrix in the figure below represents the distances between data points, while the table on the right lists the ground truth labels. Points a, b, c, d represent training sample while  $x_1, x_2$  are test samples.

	a	b	c	d	$x_1$	$x_2$
a	0	1	2	3	2	2
b	1	0	2	1	2	5
c	2	2	0	3	4	2
d	3	1	3	0	1	4
$x_1$	2	2	4	1	0	4
$x_2$	2	5	2	4	4	0

Train      Test

Sample      ground-truth

a	True
b	True
c	False
d	False

- Q1 if the value of K is set to 2, what is the number of distances computed by the K-nearest neighbor algorithm to obtain the label of  $x_1$   $\Rightarrow$  The number of computed distance is 4

Q2  $\Rightarrow$

Where

Date

Who

Why

- ⑦ Given the transactional dataset shown in the figure below. Apply the Apriori Algorithm to extract all frequent itemsets. The value of minsup is 2.

An itemset is considered to be frequent if its support count is equal to or higher than the minsup.

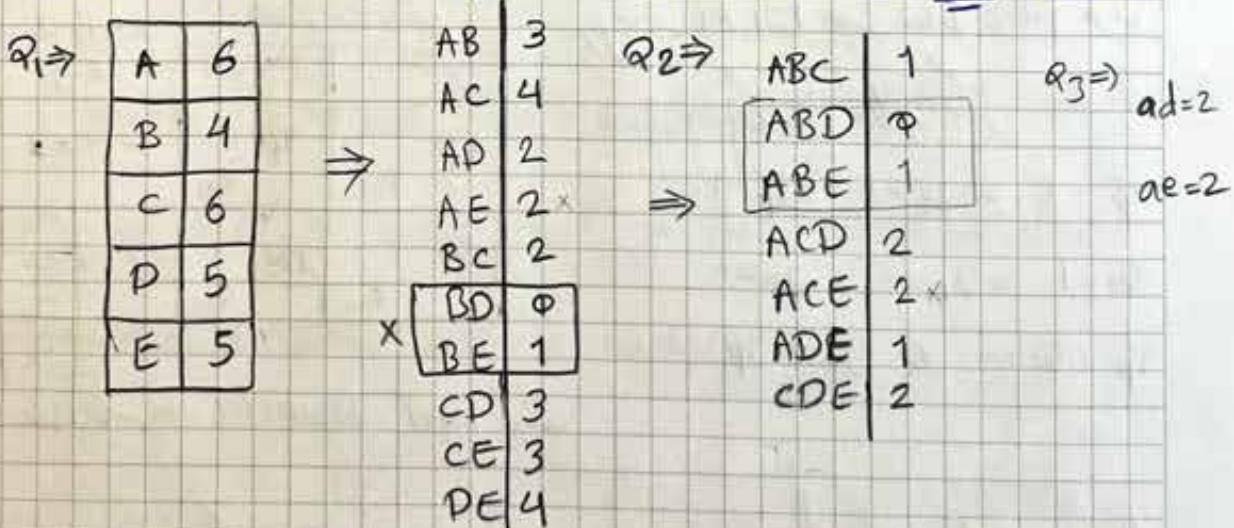
### Transactions

ab
cde
bc
ab
acd
abce
acde
de
ac
de

Question 1) list all of frequent items having length 2,

Question 2) list all itemset of length 3 that have been generated by Apriori after the join and prune steps before counting their support on the database.

Question 3) list all frequent itemsets that are not closed



### Where

Date

$A \rightarrow B$

Who

### Why

## Support

$$\text{Sup} = \frac{\text{bItemset}}{\text{نادر مثل سطرها}} = \frac{\#\{N, P\}}{|T|}$$

$\sup >, \min \sup$

## confidence

$$\text{Conf} = \frac{\text{Lo itemset}}{\text{تعداد مجموعه}} = \frac{\text{Sup}(A, B)}{\text{Sup}(A)}$$

$$\text{conf} > \text{min conf}$$

\* Frequent itemset : is a itemset whose support is greater than or equal to minSup threshold

## Rand Index

$$\text{Rand Index} = \frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}} = \frac{T_N + T_P}{T_N + F_P + T_P + F_N}$$

بازیابی RI اینجا جای مالام سایر هر دار  
نادر شوی خود را در آن سعی نماید پس از آن دوسر

$$T_p = 1 \Rightarrow j_t = 1, j_p = 1$$

$$T_N=1 \Rightarrow y_f=0, y_P=0$$

$$T_P = T_N \Rightarrow \exists t \neq y_P$$

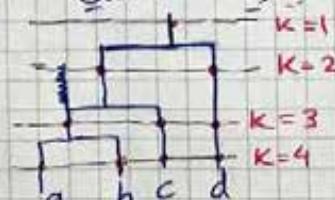
دسته بندی و آنچه		پیش خود گزینی بازرسی	
FP	TP	دسته خوب	پیش خود گزینی
$P_{\text{F}} = \frac{\text{FP}}{\text{FP} + \text{TN}}$	$P_{\text{TP}} = \frac{\text{TP}}{\text{TP} + \text{FN}}$	دسته خوب	پیش خود گزینی
$P_{\text{TN}} = \frac{\text{TN}}{\text{FP} + \text{TN}}$	$P_{\text{FN}} = \frac{\text{FN}}{\text{TP} + \text{FN}}$	در درج گزینی	در درج گزینی

نیت hyper Parameter بجای RI،

## Hierarchical clustering:

## Single linkage (min)

complete linkage (max)



در چارکل دنیاں کترن خدار میں ہے ریاستی نامہ

در فریز و بیال پیشتر مقرر بیس - پیشتر ماحصله با دور زدی

Where

Date

Who

Why

## Naive Bayes classifier

$$x \in R^n \quad x = (x_1, x_2, \dots, x_n)$$

$$P(\text{class} | x) = \underbrace{P(x_1 | \text{class})}_{\substack{\text{خوارج} \\ \Delta \\ x}} * \underbrace{P(x_2 | \text{class})}_{\substack{\text{خوارج} \\ \Delta \\ x}} * \dots * P(\text{class})$$

مجموع كل خوارج موزع على

## Manhattan distance

$$\text{dist}(P_1, P_2) = |x_1 - x_2| + |y_1 - y_2|$$

## DBSCAN :

core : Points with at least minpoints neighbors within a distance

border: at least a core point in the neighbor hood

noise: neither a core nor a border

## accuracy score in classification

$$\text{Accuracy} = \frac{\text{نواتي صحيح}}{\text{نواتي صحيح + غير صحيح}} = \frac{\text{نواتي صحيح}}{\text{نواتي صحيح + نواتي خاطئ}}$$

		Actual		
		a	b	c
predict	a	1	0	0
	b	0	1	0
	c	0	0	1
	d	0	0	0

$$\text{Recall}(A) = \frac{\text{نواتي صحيح}}{\text{نواتي صحيح + غير صحيح}} = \frac{\text{نواتي صحيح}}{\text{نواتي صحيح + نواتي خاطئ}} \Rightarrow P_{\text{truth}}$$

$$\text{Precision}(A) = \frac{\text{نواتي صحيح}}{\text{نواتي صحيح + غير صحيح}} = \frac{\text{نواتي صحيح}}{\text{نواتي صحيح + نواتي خاطئ}} \rightarrow P_{\text{predict}}$$

$$F1-\text{Score} = \frac{2P \times R}{P+R}$$

measures of node impurity

Where Date

classification:

\* Decision Tree:

- 1) Hunt's Who Algorithm
- 2) CART
- 3) ID3, C4.5, C5.0
- SLIQ, SPRINT

Why

- Gini index
- Entropy
- misclassification error

(X) interpretable

Gini  $\Rightarrow$   
split

$$gini_{split} = \sum_{i=1}^K \frac{n_i}{n} Gini(i)$$

نطر کو رکار خار فریز نہیں

نطر کو رکار خار نہیں

$$Gini(i) = 1 - \sum [p(j/t)]^2$$

نطر کو رکار خار فریز نہیں

\* Rule-based classifier

\* i) mutually exclusive :

1) دو فریزانہ حوال درست نہیں

2) حکم کو دو بیرونی اور دست بھیں نہیں مکانیں (Rule)

\* ii) exhaustive rules :

حریکیں ملکے از محاذیں مولیں خار فریز کیا جائے

حریکیں ملکے از محاذیں مولیں خار فریز کیا جائے

\* Associative classification

\* K-nearest neighbor

(Y) incremental

$$\text{Euclidean distance} \Rightarrow d(p, q) = \sqrt{\sum (p_i - q_i)^2} \quad \text{not (X)}$$

$$w = \frac{1}{d^2}$$

Terms occurring frequently in a single document but rarely in the whole collection are preferred.

Where

Date

Who

Why

### \* Bayesian classification

Fully incremental  
not (X)

{ C: any class label

{ X:  $\langle x_1, \dots, x_k \rangle$  record to be classified

① Compute  $\underline{P(C|X)}$  for all classes  $\rightarrow \underline{(P(C|X))}$

② Assign  $\underline{X}$  to the class with maximal  $\underline{P(c|X)}$

$$P(c|X) = \underbrace{P(x_1|c) \cdot P(x_2|c) \cdots P(x_k|c)}_{N \times N} / \underbrace{P(X)}_{\text{نحو عکس}} \cdot \underbrace{P(c)}_{N \times N}$$

$$P(x_1, x_2, \dots, x_k | c) = P(x_1|c) \cdot P(x_2|c) \cdots \cdot P(x_k|c)$$

$$\Rightarrow P(x_k | c) = \frac{|x_k|}{N_c} \quad \begin{array}{l} \text{one instance} \\ \text{in } x_k \text{ has } c \\ \text{in doc } c \text{ has } K \text{ terms} \end{array}$$

Probability distribution

### \* Support vector machines (SVM)

Not (Y)  
Not (X)

### \* Artificial neural networks

Not (Y)  
Not (X)

### tf-idf

tf-idf of term  $t$  in document  $d$  of collection D (consist of  $m$  documents)

$$\text{tf-idf}(t) = \text{freq}(t, d) \times \log\left(\frac{m}{\text{freq}(t, D)}\right) \quad \text{② A collection of heterogeneous doc}$$

Suitable for: ① A single doc consisting of many sections or sub-sections

Where

Date

Who

Why

## clustering Algorithms:

- 1) K-means and its variants  $\Rightarrow$  closeness  $\geq$
- 2) Hierarchical
- 3) Density-based

Euclidean distance  
cosine  
correlation  
similarity

tf-df, suitable for: ① Single Documents or parts of document with **homogeneous** content. ② A collection of documents **varying** over the same topic.

### Minkowski Distance

$$d(x, y) = \left( \sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

where n is number of dimensions.

If  $r=1 \Rightarrow$  Manhattan Distance  $\Rightarrow \text{dis}(x_1, y_1) = |x_1 - y_1| + |y_1 - y_2|$

If  $r=2 \Rightarrow$  Euclidean distance  $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$   
dissimilarity  $\rightarrow$

If  $r \rightarrow \infty \Rightarrow$  supremum  $\Rightarrow$  maximum difference between any component of the vectors.

### tf-df

tf-df of term t in document d of collection D

$$\text{tf-df}(t) = \text{freq}(t, d) + \log(\text{freq}(t, D))$$

Terms occurring frequently both in a single document and in the whole collection are preferred.

1) Where

Fp-growth

Date

- ① A-CPB
- ② A-CHT
- ③ AD-CPB

Who

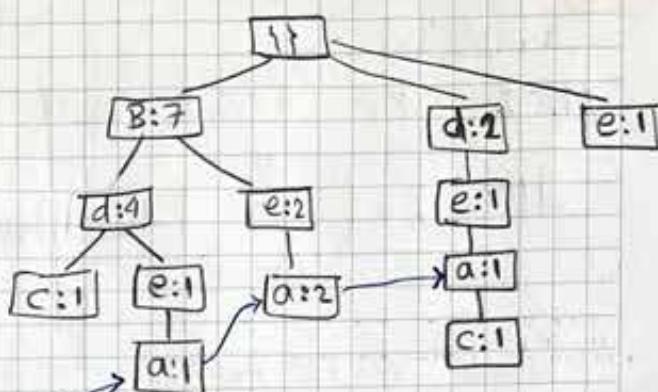
Why

BCD	BDC
ACDE	DEAC
BD	✓
B	✓
ABE	BEA
E	✓
ABE	BEA
BD	✓
D	✓
ABDE	BDEA

A	4
B	7
C	2
D	6
E	5

SICU B

B	7
D	6
E	5
A	4
C	2



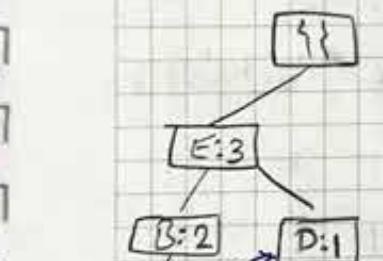
A-CHT

Bde	1	Count B,E B,E
Be	2	
de	1	

B	3
D	2
E	4
P	2

SICU  $\Downarrow$   
A-CPB

EBD	1
EB	2
ED	1



D:2	$\rightarrow$	D:1
	$\Rightarrow$	
EB	1	
E	1	
	$\Rightarrow$	
E	2	
B	1	

(3)

abd  $\Rightarrow$  abd  
 abde  $\Rightarrow$  abde

D-CPB = ? {ab:1, ae:2, aeb:2}

D-CHT = ? {a:5, e:4, b:3}

DB-CPB = ? {a:1, ae:2}

ce

abe

acde

abde

abc

{ }

{ }

{ }

{ }

{ }

{ }

{ }

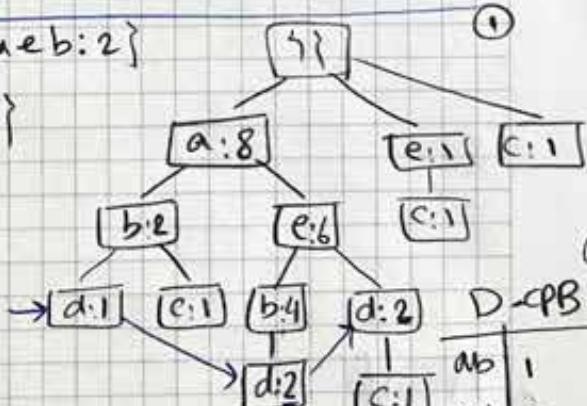
{ }

DB-CPB

ade

abc

a	1
ae	2



1

2

D-CPB

ab	1
aeb	2
ae	2

a	5
e	4
b	3

(ab(1), aeb(2))	1+2=3
-----------------	-------

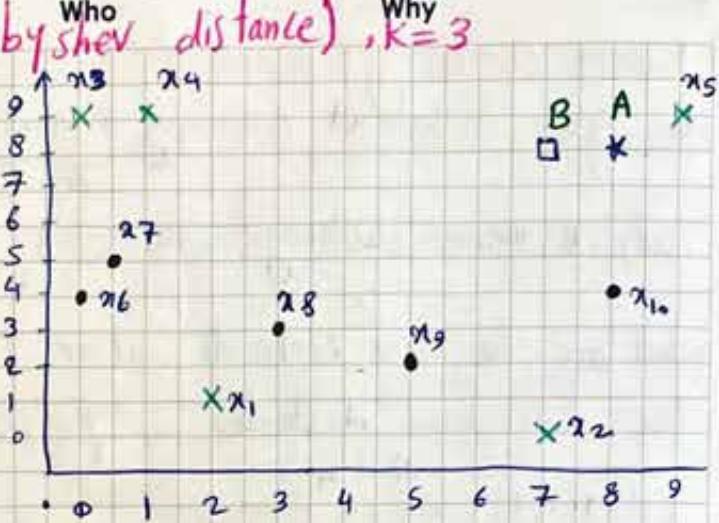
10 جولی ۲۰۲۱

2) ۲۰۲۱/۸/۶ (۱)

$$\rightarrow \max_{i=1}^n |x_i - y_i| \text{ where } x, y \text{ is two points}$$

Where KNN Algorithm (with chebychev distance), k=3

	A	B
d(x_i)	(8,8)	(7,8)
x_1 (2,1)	7	7
x_2 (7,0)	8	8
x_3 (0,9)	8	7
x_4 (1,9)	7	6
x_5 (9,9)	1 ✓	2 ✓
x_6 (0,4)	8	7
x_7 (1,5)	7	6
x_8 (3,3)	5 ✓	5 ✓
x_9 (5,2)	6	6
x_{10} (8,4)	4 ✓	4 ✓



١) مجموعه داده کو اسکیم (نمایش) کریں

٢) صدی درود را اس فرمات کر کو طلاق کریں

Question 1)  $\left\{ \begin{array}{l} A \Rightarrow 20 \\ 1 \times \end{array} \right. \Rightarrow \text{class}(A) = \text{circle} \quad (\forall i, B), (\forall i, A) \text{ را بینو} \quad ③$

$$\max_{i=1}^n |x_i - y_i| \text{ را بینو}$$

↓  $B \Rightarrow \left\{ \begin{array}{l} 20 \\ 1 \times \end{array} \right. \Rightarrow \text{class}(B) = \text{circle} \quad (k=3) \text{ را بینو} \quad ④$

what is the label assigned to A, B when neighbors' votes are weighted uniformly?  $\rightarrow$  when neighbors' votes are weighted uniformly?

Question 2) What are the class probabilities for A, B, if the vote cast by

each neighbor is  $w = \frac{1}{1 + \text{distance}}$

$$\frac{\frac{1}{2}}{\frac{11}{30} + \frac{1}{2}} = \frac{\frac{1}{2}}{\frac{26}{30}} = \frac{15}{26}$$

A	w	Class Label
x_5	1	x
x_8	5	0
x_{10}	4	0

$$p(X|A) =$$

$$\frac{1}{2}$$

$$\frac{\frac{1}{2}}{\frac{11}{30} + \frac{1}{2}} = \frac{\frac{1}{2}}{\frac{26}{30}} = \frac{15}{26}$$

$$\frac{\frac{1}{2}}{\frac{11}{30} + \frac{1}{2}} = \frac{\frac{1}{2}}{\frac{26}{30}} = \frac{15}{26}$$

B	w	Class Label
x_5	2	x
x_8	5	0
x_{10}	4	0

$$p(X|B) = \frac{\frac{1}{3}}{\frac{1}{3} + \frac{11}{30}} = \frac{\frac{1}{3}}{\frac{21}{30}} = \frac{30}{3 \times 21} = \frac{10}{21}$$

$$p(0|B) = \frac{\frac{11}{30}}{\frac{1}{3} + \frac{11}{30}} = \frac{\frac{11}{30}}{\frac{21}{30}} = \frac{11}{21}$$

پرسیده کوئال  
 $\hookrightarrow p(C|x) = \frac{\text{vote}_c(x)}{\sum_{i \in C} \text{vote}_i(x)} \Rightarrow \text{vote}(x) \Rightarrow \text{unnormalized vote assigned to class } Y \text{ for sample } x$

3) 2021/2/10

Where

Date

Who

Why

## K-means clustering - Centroids

You given a dataset containing 9 points in 2 dimension ( $X_1, X_2$ )

A-I

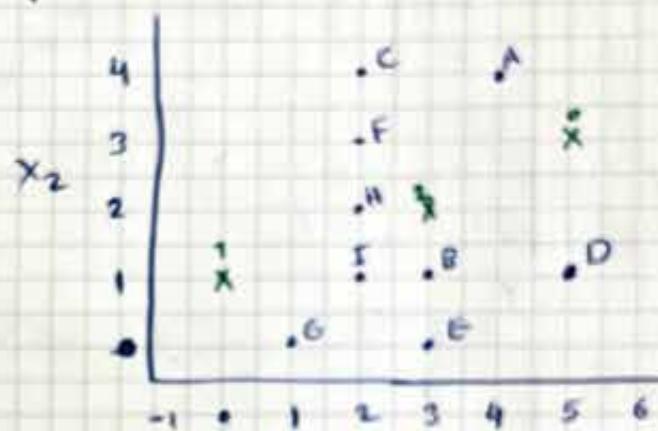
Apply K-means clustering with  $K=3$   $\Rightarrow$  (A-I)  
 $\begin{cases} \text{(1) Blue points} \\ \text{(2) red stars} \end{cases} \Rightarrow \text{Centroids (c1, c2)}$

what are the new centroids after 1 iteration of K-means algorithm?

Use Euclidean distance  $\Rightarrow \sqrt{\sum (x_i - y_i)^2} \Rightarrow \sqrt{(x_p - x_i)^2 + (y_p - y_i)^2}$

(1,3) (4,1) (3,1)

	$dist(x_i)$	$c_0$	$c_1$	$c_2$	
A	(4,4)	$\sqrt{3}$	5	$\sqrt{5}$	$c_0$
B	(3,1)	$\sqrt{8}$	3	1	$c_2$
C	(2,4)	$\sqrt{10}$	$\sqrt{13}$	$\sqrt{5}$	$c_2$
D	(5,1)	2	5	$\sqrt{5}$	$c_0$
E	(3,0)	$\sqrt{13}$	$\sqrt{10}$	2	$c_2$
F	(2,3)	3	$\sqrt{8}$	$\sqrt{2}$	$c_2$
G	(1,0)	5	$\sqrt{2}$	$\sqrt{8}$	$c_1$
H	(2,2)	$\sqrt{10}$	$\sqrt{5}$	1	$c_2$
I	(2,1)	$\sqrt{13}$	2	$\sqrt{2}$	$c_2$



$$\begin{aligned} & c_0 \Rightarrow \{(A, D)\} = \frac{1}{2} \left( \begin{pmatrix} 5, 1 \\ 4, 4 \end{pmatrix} \right) = \left( \frac{9}{2}, \frac{5}{2} \right) \\ & c_1 \Rightarrow \{(G\}\} = (1, 0) \end{aligned}$$

$$\begin{aligned} & c_2 \Rightarrow \{B, C, E, H, I\} = \frac{1}{5} \left( \begin{pmatrix} 3, 1 \\ 2, 4 \\ 2, 2 \\ 2, 1 \\ 2, 3 \end{pmatrix} \right) = \left( \frac{14}{5}, \frac{11}{5} \right) \end{aligned}$$

کاں کا مکان میرے ہار مرے طاہر بھر جائے گا ①  
 اپنے کمرے میں میرے ہار میرے دست قبایل رنجو ②

Q1 & Q2

Where

Date

Who

Why

A binary classifier is trained to separate between images of cats and dogs. The test set used to evaluate this model is balanced, with 10,000 images of cats and 10,000 images of dogs.

The classifier only predicts 50 images as being cats. All those predictions are correct.

What can said about this classifier?

$$\text{Cats} = 10$$

$$\text{dogs} = 10$$

$$\text{predict(cat)} = 5$$

Dog (10)	-5
Cat (10)	+5

out of 5 cat is 5

5 out of 10 dog

$$\text{Recall(cat)} = \frac{5}{10} = 0.5$$

$$\text{Recall(dog)} = \frac{10}{10} = 1$$

It has high recall for the class "dog"

There are 3 classes (Squares, Circle, triangle) After training. (Ans)

where the label assigned predicted class = ■

by C is given by the background = ●

color of the region in. = ▲



Suppose that the symbols in the image are test points whose shape represents the respective ground truth label.

$$P_{\square} = \frac{3}{4}$$

$$P_{\Delta} = \frac{5}{5} = 1$$

$$P_{\circ} = \frac{8}{11}$$

$$R_{\square} = \frac{3}{5}$$

$$R_{\Delta} = \frac{5}{6}$$

$$R_{\circ} = \frac{8}{9}$$

$$\text{Accuracy} \Rightarrow \frac{3 + 8 + 5}{20} = \frac{16}{20} = 0.8$$

4)

Where

Date

Who

Why

The Pearson Correlation between two variables:  $X = \{x_1, x_2, \dots, x_n\}$  having mean  $\bar{x}$  and  $Y = \{y_1, y_2, \dots, y_n\}$  having mean  $\bar{y}$  is

defined as:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \times \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Two sequences A and B have been sampled from a distribution. The following Seq are included below both in numerical and graphical:

$$B = 95A^2 + 5 \quad A \Rightarrow X=Y$$

What is the Pearson Correlation between A, B ( $r_{AB}$ )?

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 0 \Rightarrow r_{AB} = 0$$

$$\bar{x} = 0 \quad (\Rightarrow x \neq 0) = -37,5$$

$$\bar{y} = 10 \Rightarrow +1 - (-10) = -12$$

$$+(-2 \times -3) = +6$$

$$+(-1 \times -8) = +8$$

$$+(0 \times 5) = 0$$

$$+(1 \times 5,5) = +5,5$$

$$+(2 \times -3) = -6$$

$$+(3 \times -7,5) = -37,5$$

$$+(4 \times 3) = 12$$

$$+(5 \times 7,5) = 37,5$$

A	B
-5	17,5
-4	13
-3	9,5
-2	7
-1	5,5
0	5
1	5,5
2	7
3	9,5
4	13
5	17,5

5)

Where

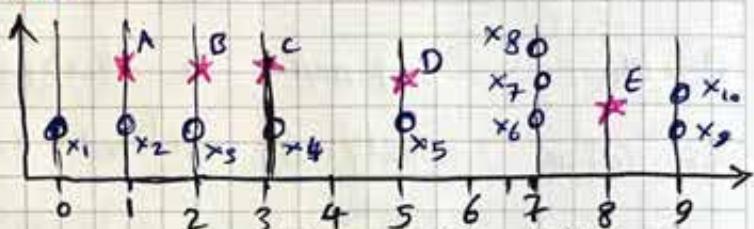
Date

Who

Why

## K-Means algorithm - centroids

$$K=5 \leftarrow A, B, C, D, E, F$$



$$A = \frac{1}{2} (1+0) = \frac{1}{2}$$

$$B = 2$$

$$C = 3$$

$$D = 5$$

$$E = \frac{1}{5} \left( \underbrace{7+7+7}_{27} + \underbrace{9+9}_{18} \right) = \frac{39}{5}$$

		(1)	2	3	5	8	A
0	x <sub>1</sub>	1	2	3	5	8	A
1	x <sub>2</sub>	2	1	4	7		A
2	x <sub>3</sub>	1	2	3	6		B
3	x <sub>4</sub>	2	1	4	5		C
5	x <sub>5</sub>	4	3	2	6	3	D
7	x <sub>6</sub>	6	5	4	2	1	E
7	x <sub>7</sub>	6	5	4	2	1	E
7	x <sub>8</sub>	6	5	4	2	1	E
9	x <sub>9</sub>	8	7	6	4	1	E
9	x <sub>10</sub>	8	7	6	4	1	E

## 7) Cosine Similarity - one hot encoding - Minkowski

Where

Date

Who

Why

The following definitions refer to any two vectors

$$a = (a_1, a_2, \dots, a_n) \text{ and } b = (b_1, b_2, \dots, b_n) \in \mathbb{R}^n$$

The cosine similarity is defined as:

$$\cos(a, b) = \frac{a \cdot b}{\|a\| \|b\|}$$

The Minkowski distance of order  $p$  between  $a$  and  $b$  is defined as:

$$D_p(a, b) = \left( \sum_{i=1}^n |a_i - b_i|^p \right)^{\frac{1}{p}}$$

You are given a categorical dataset with two independent features  $x_1$  and  $x_2$ . The two features can assume one of 13 and 7 distinct values respectively. more specifically.

$$x_1 = (n_{00}, n_{01}, n_{02}, \dots, n_{12})$$

$$x_2 = (m_{00}, m_{01}, \dots, m_{06})$$

From this dataset two records  $r_1 = (\underline{n_{04}}, m_{06})$  and  $r_2 = (\underline{n_{11}}, m_{01})$  are extracted.

① What is the cosine similarity between  $r_1, r_2 \Rightarrow \cos(r_1, r_2) = ?$

$$n_{04} = (0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11 \ 12)$$

$$\cos(r_1, r_2) = ?$$

$$n_{11} = (0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0)$$

$$m_{06} = (0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1)$$

$$m_{01} = (1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0)$$

② What is the Minkowski distance of order 1 between  $r_1, r_2$ ?

$$d_1(r_1, r_2) = \sum_{i=1}^2 |a_i - b_i| = 4 \quad |1-0| + |0-1| + |0-1| + |1-0| \quad D_1(r_1, r_2)$$

$$d_2(r_1, r_2) = \sqrt{2} \quad D_2(r_1, r_2)$$

$$d_2(r_1, r_2) = \sqrt{\sum_{i=1}^2 |a_i - b_i|^2} = \sqrt{4} = 2$$

8)

Where

Date

Who

Why

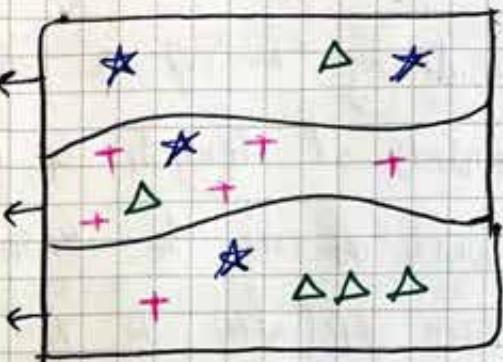
# Random forest classifier

Precision | Recall

predict(A)  $\star$

predict(B)  $+$

predict(C)  $\Delta$



$$P(A) = \frac{2}{3}$$

$$R(A) = \frac{2}{4} = \frac{1}{2}$$

$$P(B) = \frac{5}{7}$$

$$R(B) = \frac{5}{6}$$

$$P(C) = \frac{3}{5}$$

$$R(C) = \frac{3}{5}$$

9)

Where

Date

Who

Why

## Apriori Algorithm

considering  $\text{minSup} = 2 \Rightarrow$  only consider the extraction of itemsets of length up to 3

- 1) which ~~are~~ are the itemsets that can be pruned during the purify step (after the join step)?
- 2)  $= = = =$  that are pruned after computing their support count?

ADE	A	7	AB	5	ABC	2	ABCD	X
AB	B	7	AC	2	ABD	2	ABCE	X
BCD	C	4	AD	3	ABE	3	ABDE	P
TDE	D	6	AE	5				
ABCDE	E	6	BC	3				
			BD	4				
			BE	4				
			CD	3				
			CE	1				
			DE	4				

$\Rightarrow$   $\boxed{ACD}$   $\boxed{1} \Rightarrow < \text{minSup}$

$\Rightarrow$   $\boxed{ACE}$  X  $\Rightarrow$

$\Rightarrow$   $\boxed{ADE}$  3

$\Rightarrow$   $\boxed{BCD}$  2

$\Rightarrow$   $\boxed{BCE}$  X

$\Rightarrow$   $\boxed{BDE}$  3

ABE

ABCDE

ABDE

(1) {ACE, BCE, ABCD, ABCE}

(2) {CE, ACD}

ii)

Where

Date

Who

Why

## clustering, Euclidean distance - cluster Similarity

	x	y
C <sub>1</sub>	P <sub>1</sub>	2
C <sub>1</sub>	P <sub>2</sub>	4
C <sub>1</sub>	P <sub>3</sub>	3
C <sub>2</sub>	P <sub>4</sub>	11
C <sub>2</sub>	P <sub>5</sub>	13
C <sub>3</sub>	P <sub>6</sub>	8
	7	

$$C_1 = \frac{1}{3} \begin{pmatrix} (2, 2) \\ (4, 2) \\ (3, 5) \end{pmatrix} = \left( \frac{9}{3}, \frac{9}{3} \right) = (3, 3)$$

$$C_2 = \frac{1}{2} \begin{pmatrix} (11, 2) \\ (13, 4) \end{pmatrix} = \left( \frac{24}{2}, \frac{6}{2} \right) = (12, 3)$$

$$P_6 = (8, 7)$$

① distance(C<sub>1</sub>, C<sub>2</sub>) =  $\sqrt{9^2 + 9^2} = 9\sqrt{2}$

$$\text{distance}(C_1, P_6) = \sqrt{91} = 9.50$$

$$\text{distance}(C_2, P_6) = \sqrt{32} = 5.657$$

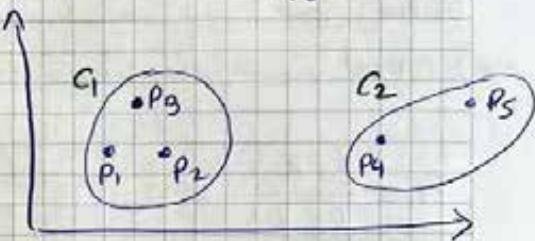
2) C<sub>1</sub> ⇒ (P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>)

C<sub>2</sub> ⇒ (P<sub>4</sub>, P<sub>5</sub>, P<sub>6</sub>)

① what is the proximity matrix at this step  
write the answer as distances between pairs  
of points/clusters?

② what are the clusters obtained after applying  
an additional step?

• P<sub>6</sub>



السؤال ١: مسافة بين C1 و C2

السؤال ٢: مسافة بين C1 و P6

السؤال ٣: مجموع مسافات كل نقطة في C1 من مركزها

2)

Where

Date

Who

Why

Classification:

Given the following confusion matrix

Q<sub>1</sub> ⇒ Compute the accuracy score.

Q<sub>2</sub> ⇒ Compute F-measure (F1) of class b :

$$\text{Accuracy} = \frac{10 + 4 + 10 + 6}{40} = \frac{30}{40} = \frac{3}{4}$$

(Accuracy جزوی است: تعداد موارد مطابق با کلاس اصلی بر تعداد مجموع موارد)

why: Incorrect classification happened

	a	b	c	d
a	10	0	2	0
b	0	4	0	4
c	0	4	10	0
d	0	2	0	6

← predicted

actual

$\sum_{i=1}^4 10 = 40$

$$F_1 = \frac{2PR}{P+R} \Rightarrow \frac{2 \times \frac{2}{5} \times \frac{1}{2}}{\frac{2}{5} + \frac{1}{2}} = \frac{\frac{2}{5}}{\frac{9}{10}} = \frac{20}{95} = \frac{4}{19}$$

$$P(B) = \text{precision}(B) = \frac{4}{10} = \frac{2}{5}$$

$$R(B) = \text{Recall}(B) = \frac{4}{8} = \frac{1}{2}$$

3) Regression: multi Regression Pipeline

Given the following dataset with 2 features ( $x_0, x_1$ ) and 3 data points:  $x = [[2, 4], [1, 2], [2, 0]]$

Apply to x the following multi regression pipeline

1) Feature extraction step  $\Rightarrow [x_0, x_1, x_0^2, x_1^2, x_0 x_1]$

2) Regression parameters (to be applied on the extracted features)

$$\beta = [0, 2, 0, 1, 1/2], \text{ Bias} = 1$$

Q<sub>1</sub> ⇒ What is the output vector with the predictions?  $y_{\text{pred}} = [?]$

Q<sub>2</sub> ⇒ Given the ground truth predictions  $\Rightarrow y_{\text{truth}} = [28, 9.5]$

Compute the mean Absolute error (MAE) of the obtained predictions ( $y_{\text{pred}}$ ).

Where

Date

$$X = \begin{bmatrix} 2 & 4 \\ 1 & 2 \\ 2 & 0 \end{bmatrix} \Rightarrow X_{\text{poly}} \Rightarrow \begin{bmatrix} x_0 & x_1 & x_0^2 & x_1^2 & x_0x_1 \end{bmatrix}$$

$\begin{bmatrix} 2 & 4 \\ 1 & 2 \\ 2 & 0 \end{bmatrix} \quad \begin{bmatrix} 16 & 8 \\ 4 & 2 \\ 0 & 0 \end{bmatrix} \quad \begin{bmatrix} * & * \\ 1 & 12 \end{bmatrix}$

$$\Rightarrow \hat{y}^* = \begin{bmatrix} 4x_2 + 16x_1 + \frac{8}{2} \\ 2x_2 + 4x_1 + \frac{2}{2} \\ 2x_0 + 1x_0 + 1x_0 \end{bmatrix}$$

$B \Rightarrow \begin{bmatrix} x & x \\ x & x \end{bmatrix}$

Q1)

$$\Rightarrow \hat{y}^* = \begin{bmatrix} 28 \\ 9 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 29 \\ 10 \\ 1 \end{bmatrix}$$

Bias

$$Q_2) \Rightarrow MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \Rightarrow \frac{1}{3} \left[ \begin{bmatrix} 29 \\ 10 \\ 1 \end{bmatrix} - \begin{bmatrix} 28 \\ 9 \\ 0 \end{bmatrix} \right] = \frac{1}{3} (1+1+4) = 2$$

$y_{\text{predicted}}$        $y_{\text{truth}}$

#### 4) computation of indices - clustering RI $\Rightarrow$ Rand Index

Given the labels predicted by a clustering algorithm and ground truth labels:

$$y_{\text{truth}} = [1, 1, 1, 2]$$

$$y_{\text{pred}} = [3, 3, 1, 1]$$

compute the Rand Index (RI)  $\Rightarrow \frac{TP + TN}{\binom{n}{2}} = \frac{n(n-1)}{2}$

	truth	predicted	TP	TN
[0,1]	1	1	1	0
[0,2]	1	0	0	0
[0,3]	0	0	0	1
[1,2]	1	0	0	0
[1,3]	0	0	0	1
[2,3]	0	1	0	0
			1	2

$$\Rightarrow \frac{1+2}{4(3)} = \frac{3}{6} = \boxed{\frac{1}{2}}$$

Jun 4 J 2020  
2021, 1, 28

5) Where

Date

Who

Why

## clustering - DBSCAN - silh(p<sub>i</sub>)

Given the following distance matrix (each cell describes the distance between two points)

Apply DBSCAN clustering - hyperparam  $\rightarrow \epsilon = 5$ , minpoints = 2

Q<sub>1</sub>  $\Rightarrow$  label each point with B(border), C(core), N(noise)

a	b	c	d	e	f	g
C	B	C	C	N	C	B

Q<sub>2</sub>  $\Rightarrow$  Assign a cluster id to each point:

a	b	c	d	e	f	g
0	1	0	1	-1	0	1

Q<sub>3</sub>  $\Rightarrow$  Compute the silhouette score of point g:

$$\text{silhouette}(g) = \frac{\text{inter}(g) - \text{intra}(g)}{\max(\text{inter}(g), \text{intra}(g))}$$

$$\text{intra}(g) = \frac{1}{2} (6+9) = \frac{15}{2} = 7.5$$

{d, b}

class(0, a, c, f)

$$\text{inter}(g) = \min\left(\frac{6+6+9}{3}, 8\right) = \min(7, 8)$$

class(-1, e)

$$\frac{21}{3} = 7$$

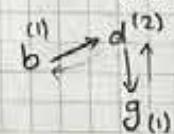
(min)

$$\text{silh}(g) = \frac{7-5}{7} = \frac{2}{7}$$

a	b	c	d	e	f	g
6	4	7	8	3	6	1
6	6	3	7	7	6	6
4	6	7	7	3	9	1
7	3	7	6	8	4	1
8	7	7	6	7	8	1
3	7	3	8	7	6	1
6	6	9	4	8	6	1



$\Rightarrow \text{class} = 0$



$\Rightarrow \text{class} = 1$

$$\text{Silh}(p_i) = \frac{\text{inter}(p_i) - \text{intra}(p_i)}{\max(\text{inter}(p_i), \text{intra}(p_i))}$$

$\Rightarrow$   $\text{intra}(p_i) \Rightarrow$  Any dissimilarity of  $p_i$  with all other point which  $p_i$  is member of

$\text{inter}(p_i) \Rightarrow$  be the lowest Any dissimilarity of  $p_i$  to any other cluster of which  $p_i$  is not a member or it

## Maximal Frequent itemset

## Where

Date

**Who**

why

An itemset is frequent maximal if none of its immediate supersets is frequent. close: A (count) is not greater than

Close: A (count) is not greater than its immediate support.

~~checklist~~

ABCD  
ABCD  
ABC  
BCD  
CD

	1
A	3
B	4
C	5
D	4

		(2)
A	B	3
A	C	3
A	D	2
B	C	4
B	D	3
C	D	4

$$A(3) \Rightarrow A\beta(3), A\zeta(3), AD^2$$

$\leftarrow \text{minsup}$  { A is not close  
A is not max }

minSup Count > 3

`{close: 5} 3,4,4`

$\max = A$ 's immediate superset, itemset

$$C(5) \Rightarrow A(3), B(4), C(4) \quad \text{not Max: } 3, 4, 4 < 3$$

minifare preset with min

$$D(9) \Rightarrow AD(2), BD(3), CD(4) \Rightarrow \begin{cases} \text{Not close} \Rightarrow 4 > 2, 3, 4 \\ \text{Not Max} \Rightarrow 2, 3, 4 < 3 \end{cases}$$

3

$ABC$	3
$BCD$	3

$$ABC(3) \Rightarrow ABC(3) \rightarrow \begin{cases} \text{not close} \rightarrow 3 > 3 \\ \text{not max} \rightarrow 3 \leq 3 \end{cases}$$

$$AC(3) \Rightarrow ABC(3) \Rightarrow \begin{cases} \text{not close} \rightarrow 3 > 3 \\ \text{not max} \end{cases}$$

$$BC(4) \Rightarrow ABC(3), BCD(3) \Rightarrow \begin{cases} \text{close : } 4 > 3, 3 \\ \text{not max: } 3, 3 < 3 \end{cases} \checkmark$$

$$BD(3) \rightarrow BCD(3) \Rightarrow \begin{cases} \text{not close: } 3 > 3 \\ \text{not max} \end{cases} \times$$

$$CD(4) \rightarrow BC(3) \Rightarrow \begin{cases} close \Rightarrow 4 > 3 \\ not max \Rightarrow 3 < ? \times \end{cases}$$

(A) ABCD 2

$$ABC(3) \Rightarrow ABCD(2) \Rightarrow \begin{cases} \text{close} \Rightarrow 3 > 2 \\ \text{max} \Rightarrow 2 \leq 3 \end{cases} \checkmark$$

$$BCD(3) \Rightarrow ABCD(2) \Rightarrow \begin{cases} \text{close} \Rightarrow 3 > 2 \\ \text{not max} \Rightarrow 2 \leq 3 \end{cases} \checkmark$$

Where

Date

Who

Why

## Frequent - Item Set :

is an itemset whose support is greater than or equal to a minsup threshold

## Support :

$$A \rightarrow B$$

is the fraction of transactions that contains an itemset  
exp:  $\{Beer, milk\} = 2/5 \rightarrow$

Transactions with itemset      Support Transaction

$$\text{Sup} = \frac{\# \{A, B\}}{T}$$

## Confidence:

$$A \rightarrow B$$

is the frequency of B in transactions containing A

$$\text{conf} = \frac{\text{Sup}(A, B)}{\text{Sup}(A)}$$

\* if  $A \subseteq B$  then  $\text{Sup}(A) \geq \text{Sup}(B)$

## Maximal Frequent Itemset

An itemset is frequent maximal if none of its immediate supersets is frequent.

## Closed Itemset

If none of its immediate superset has the same support as the itemset.



## maximal frequent itemset

Where	Date	Who	Why
ABC	A 3	AB 2	
ABCD	B 3	AC 3	
BCE	C 4	AD 2	
ACDE	D 3	AE 1	ABC 2
DE	E 3	BC 3	ABD 1
		BD 1	ACD 2
		BE 1	CDE 1
		CD 2	
		CE 2	
		DE 2	

$\text{minSup} = 2$

closed

- ① C(4)  $\Rightarrow$  AC(3), BC(3), CD(2), CE(2)  $\Rightarrow$  4 > 3, 3, 2, 2 ✓
- ② D(3)  $\Rightarrow$  AD(2), BD(1), CD(2), DE(2)  $\Rightarrow$  3 > 2, 1, 2, 2 ✓
- ③ E(3)  $\Rightarrow$  AE(1), BE(1), CE(2), DE(2)  $\Rightarrow$  3 > 1, 1, 2, 2 ✓
- ④ AC(3) ⑤ BC ⑥ CE ⑦ DE ⑧ ABC ⑨ ACD

maximal

- ① CE(2)  $\Rightarrow$  CDE(1)  $\Rightarrow$  1 < 2 ✓
- ② DE(2)  $\Rightarrow$  CDE(1)  $\Rightarrow$  1 < 2 ✓
- ③ ABC(2)  $\Rightarrow$  ABCD(1)  $\Rightarrow$  1 < 2 ✓
- ④ ACD(2)  $\Rightarrow$  ~ = ~ ✓

# Robustness to the noise

Where

Date

Who

Why

(classifications): ✓ 1) Decision Tree (hunt's algorithm)

✓ 2) Classification Rules

✓ 3) Association Rules

a) Neural network

$p(X|C)$

$p(x_1, x_2, \dots, x_n | C)$

← 5) Bayesian networks  $\Rightarrow$  fully incremental  
assign  $x$  to class with maximal  $p(C|x)$

majority vote ① → d on distance  
6) k-nearest neighbours (kNN) in one metall

weighted vote ② →  $d_{ij} = \sqrt{\sum (x_i - x_j)^2}$   $\rightarrow$  weight  $w_{ij} \propto \frac{1}{d_{ij}}$   $\rightarrow$   $\sum w_{ij} = 1$

✓ 7) support vector machines (SVM)

according to distances

✓ 8) Random forest

## Classification :

Given the following confusion matrix

Q1: Compute the accuracy score

Q2: Compute F-measure (F1) of class b

Correct

$$\Sigma = 10 + 4 + 10 + 6 = 30 \quad \text{Majority of predictions are correct}$$

$$\text{Accuracy} = \frac{\text{total correct prediction}}{\text{total prediction}} = \frac{30}{30+10} = \frac{3}{4}$$

(accuracy score)

$$Q2: F1 = \frac{2 P R}{P + R}$$

$$\begin{aligned} \text{For class } b \rightarrow & P = \frac{4}{10} \rightarrow \text{Correct predictions} \\ & R = \frac{4}{8} \rightarrow \text{Total predictions} \end{aligned}$$

$$F1 = \frac{2 \times \frac{4}{10} \times \frac{4}{8}}{\frac{4}{10} + \frac{4}{8}} = \frac{4}{9}$$

- precision : P

Recall : R

## Regression : (3 point)

Given the following dataset, with 2 features ( $x_0, x_1$ ) and 3 data points:

$$X = [[2, 4], [1, 2], [2, 0]]$$

Apply to X the following multinomial regression pipeline

1 - Feature extraction step

$$[x_0, x_1, x_0^2, x_1^2, x_0 x_1]$$

2 - Regression parameters (to be applied on the extracted features)

$$B = [0, 2, 0, 4, 1/2], \text{ Bias} = 1 \quad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ 1 \\ 16 \\ 8 \end{bmatrix} = \begin{bmatrix} 29 \\ 10 \\ 1 \end{bmatrix} = Y_{pred}$$

Q1 - What is the output vector with the predictions?

$$Y_{pred} = [?]$$

$$X = \begin{bmatrix} 2 & 4 \\ 1 & 2 \\ 2 & 0 \end{bmatrix} \Rightarrow X_{poly} = \begin{bmatrix} 1 & 2 & 1 & 1 & 2 \\ 2 & 4 & 1 & 16 & 8 \end{bmatrix} \quad * \begin{bmatrix} 4 \times 2 + 1 \times 1 + \frac{8}{2} \\ 2 \times 2 + 4 \times 1 + \frac{1}{2} \\ 1 \end{bmatrix} = \begin{bmatrix} 29 \\ 10 \\ 1 \end{bmatrix}$$

Q2: Given the ground truth predictions

$$-\hat{y}_{\text{truth}} = [28, 2, 5]$$

- compute the Mean Absolute Error (MAE) of the obtained predictions ( $\hat{y}_{\text{pred}}$ )

$$\hat{y} = \begin{bmatrix} 28 \\ 9 \\ 5 \end{bmatrix}$$

ground truth

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \Rightarrow \frac{1}{3} \left| \begin{bmatrix} 29 \\ 10 \\ 1 \end{bmatrix} - \begin{bmatrix} 28 \\ 9 \\ 5 \end{bmatrix} \right| = \frac{1}{3} |1+1+4| = 2$$

↑  
Prediction

$y_i$

↓

$\hat{y}_i$

↓

$y_i$

#### 4. Computation of indices

- Given the labels predicted by a clustering algorithm and ground truth labels:

$$\begin{array}{c} \cdot \quad 1 \quad 2 \quad 3 \\ - \hat{y}_{\text{true}} = [1, 1, 1, 2] \end{array}$$

$$\begin{array}{c} \cdot \quad 1 \quad 2 \quad 3 \\ - \hat{y}_{\text{pred}} = [3, 3, 1, 1] \end{array}$$

- Compute the Rand Index score (RI)

$$RI = \frac{TP + TN}{\binom{n}{2}} = \frac{TP + TN}{\frac{n(n-1)}{2}}$$

- Where  $TP$  = number of pairs of elements that are in the same set in  $\hat{y}_{\text{true}}$  and in the same set in  $\hat{y}_{\text{pred}}$
- $TN$  = number of pairs of elements that are in different sets in  $\hat{y}_{\text{true}}$  and different sets in  $\hat{y}_{\text{pred}}$
- $n$  = number of data points

$$n = 4$$

$\hat{y}_{\text{true}}$	$\hat{y}_{\text{pred}}$	TP	TN
[0, 1]	1	1	0
[0, 2]	1	0	0
[0, 3]	0	0	1
[1, 2]	1	0	0
[1, 3]	0	0	1
[2, 3]	0	1	0
		1	2

$$RI = \frac{1+2}{\binom{4}{2}} = \frac{3}{6} = \frac{3}{\frac{4!}{2!2!}} = \frac{3}{24/2!} = \frac{3}{12} = \frac{1}{4}$$

$$RI = \frac{3}{6} = \frac{1}{2}$$

## 6- Python-related questions

- Given two Numpy vectors  
 $X$  with shape  $(100, 50)$   
 $y$  with shape  $(50,) = (1, 50)$

$$\sqrt{\sum(x-y)^2}$$

✓

axis = 1  
ستعمل على  
العمودي  
العمودي

a)  $\text{np.sqrt}((X-y)**2).sum(\text{axis}=1)$

is the euclidean distance between rows of  $X$  and  $y$  and the result has shape  $(100, 1)$

b)  $\text{np.sqrt}(((X-y)**2).sum(\text{axis}=1))$   
has shape  $(100, )$

c)  $\text{np.sqrt}(((X-y)**2).sum(\text{axis}=0))$   
has shape  $(100, )$

d)  $\text{np.sqrt}(((X-y).sum(\text{axis}=1))**2)$   
has shape  $(100, 1)$

$$X = \begin{bmatrix} & & 1 \\ & & \\ 100 & & \\ & & \\ -50 & & \end{bmatrix} \quad Y = \begin{bmatrix} & & & \\ & & & \\ & & -50 & \\ & & & (1 \times 50) \\ & & & \\ & & & \end{bmatrix}$$

~~we can broadcast~~

$$\text{euclidean} = \sqrt{\sum(y_i - x_i)^2} \rightarrow \text{involves many additions}$$

one dimensional  $\Rightarrow$  row factor

$\Rightarrow \text{axis} = 1$

## T - Python-related questions



- Given a Dataframe with four columns (category, year, month, # subscriptions)
  - a) `df[['category', 'year']].pivot_table(['#subscriptions'], index='category', columns='year')`
  - b) `df.groupby('category').sum().unstack()`
  - c) `df.pivot_table(['#subscriptions'], index='category', columns='year', aggfunc='sum')`
  - d) `df.drop('month').groupby(['category', 'year']).sum().unstack()`
  - e) None of the previous answers is correct

## 1- Preprocessing

(2021-12-03)

- The following list represents training set values of a specific attribute  $[10, 0, 5, 3, 3, 0, 3, 4, 4, 7, 5, 7, 8, 4, 9] = 14.15$

- Use these values to train an equal-frequency based discretization with three bins (low, medium, high). Which statement is correct?

- The test vector  $[1, 7, 9]$  is discretized to [low, medium, high]
- The test vector  $[0, 7, 4]$  is discretized to [high, medium, medium]
- The test vector  $[3, 4, 7]$  is discretized to [low, medium, high]
- The test vector  $[5, 4, 2]$  is discretized to [high, medium, low]

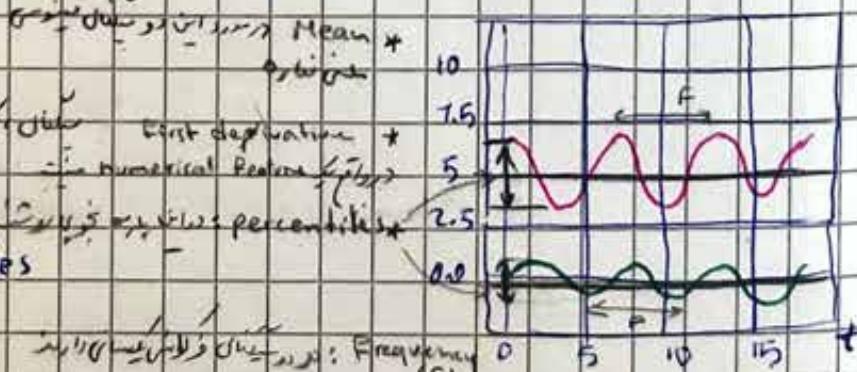
Width-Based		Number of bins: 15	Frequency
MIN	MAX	Bin width: 3 bins	Bin centers, intervals
0	2	0	0
2	5	1	2
5	8	2	3
8	11	3	4
11	14	2	5
14	17	2	7
17	20	1	8
20	23	1	9
23	26	1	10

Frequency distribution is equal to (low, medium, high) 3 bins method

## 2- Time Series

- Which is the most significant pair of features for distinguishing between the two periodic time series depicted in the figures below?

- Mean, First derivative



- Mean, Percentiles

- First derivative, Percentiles

- Percentiles, Frequency

- All of the pairs above are equivalent for distinguishing between the two series

### 3- Classification

The two dataset splits depicted in the figure represent an intermediate step of Hunt's algorithm.

a) Compute the Gini index of the two splits

$$- \text{Gini}(X), \text{Gini}(Y)?$$

b) Which of the two attribute splits will be selected by the algorithm?

$$\begin{array}{l} a) X \\ b) -Y \\ \text{partition : } p_a \end{array}$$

		Y		n = 120
		class a	class b	
partition : p <sub>a</sub>	1x	20	0	
	2y	60	60	

		X		n = 120
		class a	class b	
partition : p <sub>a</sub>	1x	60	10	
	2y	40	10	

$$X) \text{Gini}_{\text{split } j} = \sum_i \frac{n_{ij}}{n} \text{Gini}_{p_a}$$

$$\text{Gini}_{p_a} = 1 - \sum_k \left( \frac{m_k}{n_i} \right)^2 \Rightarrow 1 - \left( \frac{x_a}{X} \right)^2 - \left( \frac{x_b}{X} \right)^2$$

$$1X) 1 - \left( \frac{60}{100} \right)^2 - \left( \frac{40}{100} \right)^2 = \frac{12}{25} \leftarrow \text{Gini}_{p_a(1x)}$$

$$2X) 1 - \left( \frac{10}{20} \right)^2 - \left( \frac{10}{20} \right)^2 = \frac{1}{2} \leftarrow \text{Gini}_{p_a(2x)}$$

$$\text{Gini}_X = \frac{1x}{n} \times \text{Gini}_{p_a(1x)} + \frac{2x}{n} \times \text{Gini}_{p_a(2x)} = \frac{100}{120} \times \frac{12}{25} + \frac{20}{120} \times \frac{1}{2} = \frac{58}{120}$$

$$1Y, 1 - \left( \frac{20}{20} \right)^2 - \left( \frac{0}{20} \right)^2 = 0$$

$$2Y, 1 - \left( \frac{40}{100} \right)^2 - \left( \frac{60}{100} \right)^2 = \frac{1 - \frac{14}{100}}{100} = \frac{66}{100} = \frac{12}{25}$$

$$\text{Gini}_Y = \frac{20}{120} \times 0 + \frac{100}{120} \times \frac{12}{25} = \frac{12}{10} = \frac{4}{10}$$

درواب بایوسنٹری جیئنی اینڈ کلیسیفیکیشن  
جیئنی اینڈ کلیسیفیکیشن

$$\text{Gini}_Y = \frac{40}{100}$$

$$\text{Gini}_X = \frac{58}{120}$$

$$\text{Gini}_{p_a} = 0$$

$$\text{Gini}_{p_a \cdot 2x} = \frac{1}{2}$$

$$\text{میں جسے } \text{Gini}_{p_a} = 1 \text{ تو}$$

## 4 - Hierarchical Clustering

✓✓✓

- Given the following distance matrix, apply agglomerative hierarchical clustering with single-linkage (min).

Which statement is correct?

a) with  $k=3$  clusters, a and b are in the same cluster

b) with  $k=2$  clusters, c and d are in different clusters

c) with  $k=3$  clusters, b and c are in different clusters

d) with  $k=2$  clusters b and c are in the same cluster

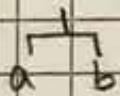
X	a	b	c	d
a	0	1	4	3
b	1	0	2	4
c	4	2	0	3
d	3	4	3	0

c) All of the previous answers are correct.

$$d(x,y) = d(y,x)$$

distance  $\rightarrow$

	A	B	C	D
A	0	2	3	
B	2	0	3	
C	3	3	0	



X1 (0,0)

	A	B	C	D
A	0	2	3	
B	2	0	3	
C	3	3	0	

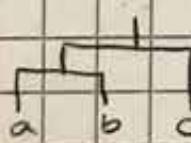
a, b = 1 ✓  $\min_{\{a,b\}} \{1, 2\}$  ① p6

a, c = 4 ✓  $\min_{\{a,c\}} \{4, 3\}$  ② p6

a, d = 3 ✓  $\min_{\{a,d\}} \{3, 4\}$

b, d = 4 ✓  $\min_{\{b,d\}} \{4, 3\}$

c, d = 3 ✓  $\min_{\{c,d\}} \{3, 0\}$

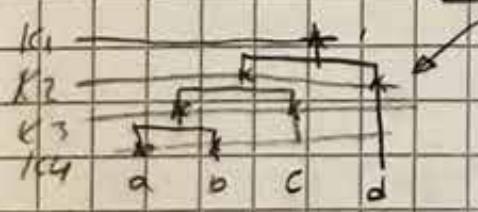


	A	B	C	D
A	0	2	3	
B	2	0	3	
C	3	3	0	

A, B, C = 2 =  $\min_{\{A,B,C\}} \{2, 0, 3\}$  X1 (0,0) ③ p6

A, B, D = 3  $\min_{\{A,B,D\}} \{3, 2, 4\}$  X1 (0,0) ④ p6

C, D = 3  $\min_{\{C,D\}} \{3, 0\}$



↳ min. link

k=1 ABCD

k=2 ABC D

k=3 ABD C D

k=4 A B C D

⑤ p6

1-(3 points)

2020/6/16 visual

Let  $X$  be a feature vector  $X \in \mathbb{R}^N$ , i.e.,  $X = (x_1, x_2, \dots, x_n)$ .  
Naive Bayes classifiers compute the output class based on the following definition:

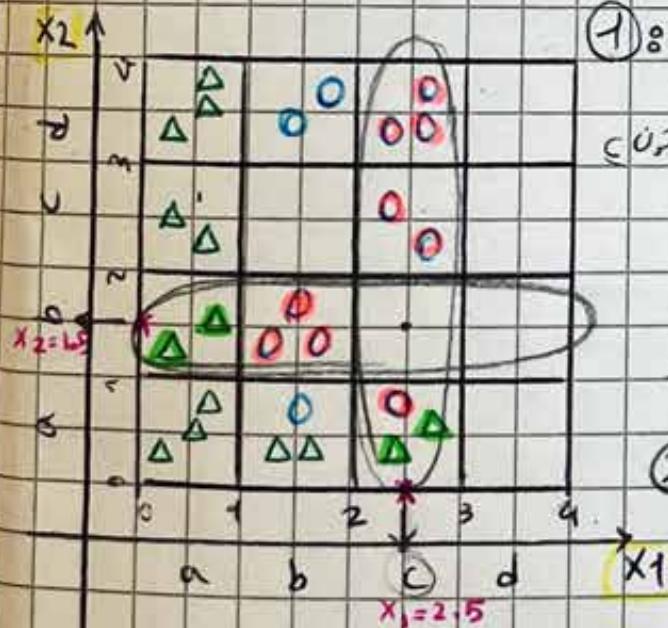
$$P(\text{class} | X) = P(x_1 | \text{class}) * P(x_2 | \text{class}) * P(x_3 | \text{class}) \dots * P(x_n | \text{class}) * P(\text{class})$$

It is given a dataset represented in Figure, where  $x_1$  and  $x_2$  are two continuous features, while triangle and circle are two class labels. The two features are discretized in four bins  $[a, b, c, d]$  depicted in the figure with a grid.

Suppose that a Naive Bayes classifier has been trained on the provided dataset, which contain 14 triangles and 12 circles.

classify the test data sample  $X$  with Feature  $x_1 = 2.5, x_2 = 1.5$ .

Write in the box the following values one answer per row.



$$\begin{aligned} ①: P(\Delta | X) &= P(x_1 | \Delta) * P(x_2 | \Delta) * P(\Delta) \\ &= P(c | \Delta) * P(b | \Delta) * P(\Delta) \\ &= \frac{2}{14} * \frac{2}{14} * \frac{14}{26} \rightarrow \frac{1}{13} \end{aligned}$$

$$\Rightarrow P(\Delta | X) = \frac{1}{13}$$

$$\begin{aligned} ②: P(O | X) &= P(x_1 | O) * P(x_2 | O) * P(O) \\ &= P(c | O) * P(b | O) * P(O) \\ &= \frac{3}{12} * \frac{3}{12} * \frac{12}{13} \\ &= \frac{3}{52} \end{aligned}$$

1 =  $P(\text{triangle} | X)$

2 =  $P(\text{circle} | X)$

3. The class assigned to  $X$  (circle or triangle)  $\rightarrow O \rightarrow$

2-(3 points) ✓✓

The Manhattan distance is defined as follows:

$$\text{dist}(P_1, P_2) = |x_1 - x_2| + |y_1 - y_2|$$

in the DBSCAN algorithm the datasets points can be labeled as:

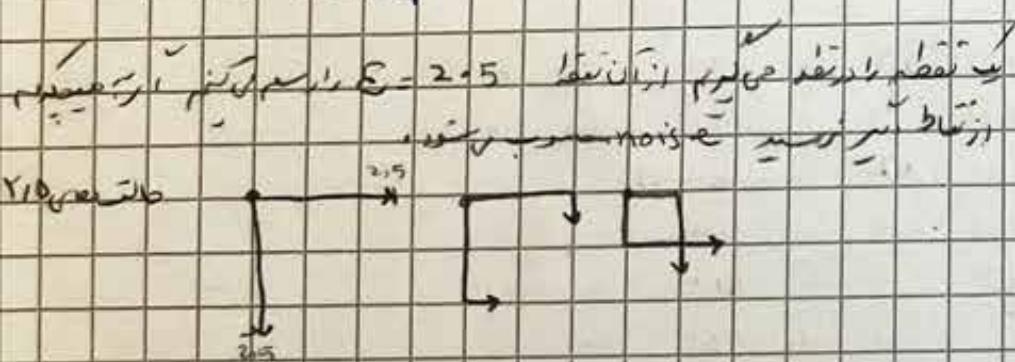
- **core** (points with at least minpoints neighbors within a distance  $\epsilon$ )
- **border** (at least a core point in the neighborhood)
- **noise** (neither a core point, nor a border point)

Using the Manhattan distance, apply the DBSCAN algorithm to the following points in the bidimensional space.

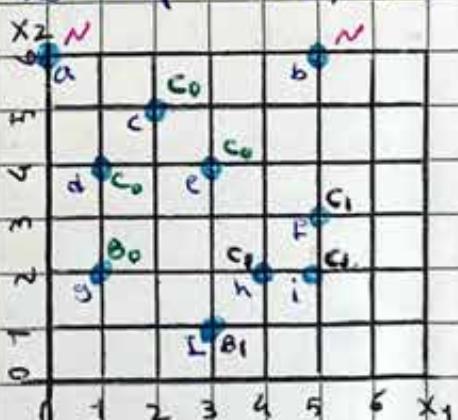
use the following hyperparameters:  $\epsilon = 2.5$ ,  $\text{minpoints} = 2$  (at least 2 points as neighbors). For each point write:

The assigned label ( $N$ =noise,  $B$ =border,  $C$ =core)

The assigned cluster id (order of cluster ids is not important, use -1 for noise points)



a	N	-1
b	N	-1
c	C	0
d	C	0
e	C	0
F	C	1
g	B	0
h	C	1
i	C	1
l	B	1



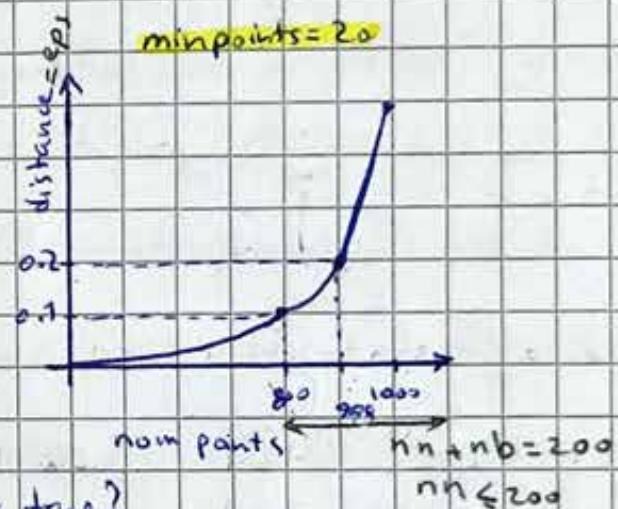
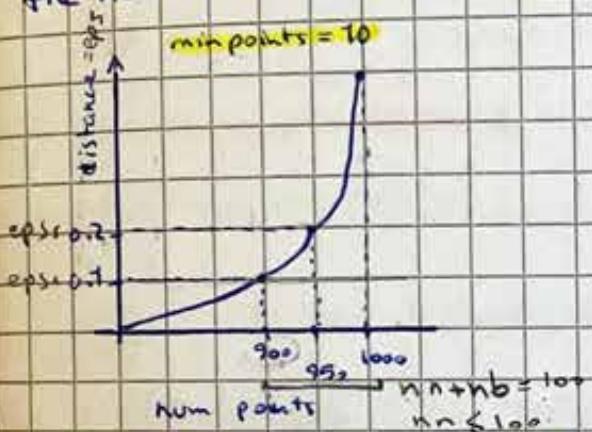
3-(1 points)

The chart shown in the figure below is typically used to set up epsilon (eps) and minpoints values in DBSCAN.

It shows on the y-axis, for each point in the dataset, the distance from the K-th nearest point, where K is equal to minpoints.

The x-axis represents all the dataset's points, ordered by increasing distance from the k-th nearest point.

Let nc the number of core points, nb the number of border points, nn the number of noise points.



which of the following statements is true?

- (a) with  $\text{minpoints} = 10 \Rightarrow \text{eps} = 0.1, \text{nc} \leq 900, \text{nn} = 100 \times \text{nc} = 900$
- (b)  $= 20 \Rightarrow \text{eps} = 0.2, \text{nc} < 900, \text{nb} \leq 100 \times \text{nc} = 900$
- (c)  $= 10 \Rightarrow \text{eps} \leq 0.2, \text{nc} + \text{nb} = 950, \text{nn} = 50 \times \text{nc} = 950 \text{ if nb} = 0 \times \text{nc} = 950$
- (d)  $= 20 \Rightarrow \text{eps} \leq 0.1, \text{nc} = 800, \text{nn} \leq 200 \times \text{nc} = 800$
- (e)  $- 10 \Rightarrow \text{eps} < 0.1, \text{nc} = 900, \text{nb} < 50 \times \text{nc} = 900$
- (f)  $- 20, \text{eps} \leq 0.2, \text{nc} = 900, \text{nn} = 100 \times \text{nc} = 1000 \times \text{nc} = 1000$

$$\text{nc} = 1000$$

Core (1)  $\times$  (2)  $\times$  (3)  $\times$  (4)  $\times$  (5)  $\times$  (6)

$$\left\{ \begin{array}{l} \text{nn} + \text{nb} = \text{Core}(1) \cdot \text{Core}(2) \dots \text{Core}(6) \\ \text{nn} \leq \\ \text{nb} \leq \end{array} \right.$$

4-(1 point) ✓

Broadcasting rules for executing an operation between two Numpy vectors are defined as follows:

- The shape of the array with fewer dimensions is padded with leading ones
  - if the shape along the dimension is 1 for one of the arrays and  $>1$  for the other, the array with shape=1 in that dimension is stretched to match the other array
  - if there is a dimension where both arrays have shape $>1$  then broadcasting cannot be performed

Apply the broadcasting rules to the following numpy vectors.

```
X = np.array ([[1,2,3],  
              [0,0,0],  
              [1,1,1]])
```

```
y=np.array([1,2,3])    z=np.array([[1],  
[0],  
[-1]])
```

$$x + y = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} =$$

row vector

$$\Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 2 & 4 & 6 \\ 1 & 2 & 3 \\ 2 & 3 & 4 \end{bmatrix}$$

shape: (3, 3) (1, 3)

$$Y+2 = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 3 & 4 \end{bmatrix}$$

$\swarrow$   
row vector       $\downarrow$   
Column vector

(b)  $x+y = \{ [2, 4, 6], [1, 2, 3], [2, 3, 4] \}$  and  $y+z = \{ [2, 3, 4], [1, 2, 3], [0, 1, 2] \}$

5 (1 point, -15% penalty for a wrong answer)

An itemset is defined to be frequent maximal if none of its immediate supersets (i.e. including one more item) is frequent.

After the analysis of a transactional dataset, the complete set of maximum frequent itemsets with  $\text{minsup} = 900$  are. ( $\text{abc} = 930$ ,  $\text{def} = 980$ )

- a)  $\sup(ab)$  can be in range  $[900, 929]$ ,  $\sup(abcd)$  cannot be in range  $[0, 930]$

b)  $\sup(ab)$  can be a value in range  $[900, 929]$ ,  $\sup(abcd)$  can be a value in range  $[700, 930]$

c)  $\sup(ab)$  can be a value in range  $[930, 1000]$ ,  $\sup(cd)$  cannot be in range  $[900, 930]$

d)  $\sup(ab)$  can be a value in range  $[930, 1000]$ ,  $\sup(abc)$  can be a value in range  $[200, 930]$

e) None of the other answer is correct X

f)  $\sup(ab)$  can be a value in range  $[900, 929]$ ,  $\sup(cd)$  cannot be in range  $[900, 930]$

$\Rightarrow \sup(ab) > 930$  ✓

b-(1 point / -15% penalty)

- A neural network is trained on a training set with 2 features  $x_1, x_2$  ✓  
•  $x_1$  is a continuous value  
•  $x_2$  is a nominal value among the following: "high", "med", "low"

Before training, the features are preprocessed in the following way:

- $x_1$  is normalized with z-score:  $x_1 \text{ norm} = (x_1 - \text{np.mean}(x_1)) / \text{np.std}(x_1)$  (training)
- $x_2$  is encoded with a one hot vector ('high': 001, 'med': 010, 'low': 100)

when applying the trained classifier on unlabeled data (denoted as test set)  
which of the following statements are true (multiple answers maybe true)

- a) the  $x_1$  feature of the test record must be normalized with the mean and standard deviation previously computed for the training set.
- b) the  $x_1$  feature of the test record can be provided to the classifier without normalization.
- c) the  $x_1$  feature of the test record must be normalized with the mean and standard deviation computed for the test set.
- d) it is suggested to process a record with a value of  $x_2$  not included in 'high', 'med', 'low' by choosing a random value for the one hot features.
- e) it is suggested to process a record with a value of  $x_2$  not included in 'high', 'med', 'low' by adding an input neuron to the neural network for new values
- 1) a record with a value of  $x_2$  not included in 'high', 'med', 'low' cannot be processed.

$x_1$   
 $x_{\text{Hot}}$   
 $x_{\text{Med}}$   
 $x_{\text{Low}}$

7-(1 point/-15) penalty)

A model's capacity defines its ability of fitting more complex data

Regularization is defined as a set of techniques to avoid overfitting.

Referring to a polynomial regression problem, which of the following statements is true?

- (A) All the statements are correct.
- (B) Adding higher power degrees (e.g.  $x^4x^5$ ) is always useful to obtain better predictions, since it increases model capacity.
- (C) Regularization decreases model capacity and it is only suitable for output variables that are linearly correlated with features
- (D) Regularization is always necessary when the output variable presents complex non linear relationships with the features
- (E) Regularization helps to decrease model capacity and fit more complex data.
- (F) Adding higher power degrees (e.g.  $x^4x^5$ ) to the input features increases model capacity and the likelihood of having overfitting

8-(1 point) -15% penalty)  
To access the quality of a classification result the following measures were defined.

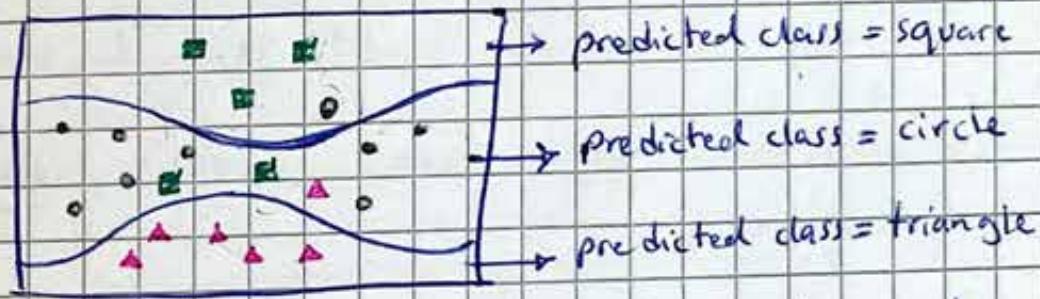
accuracy : The percentage of correctly classified items.

precision ( $c$ ) : the percentage of correctly classified items among those predicted with class  $C$

recall ( $c$ ) the percentage of correctly classified items among those with actual class  $C$

Let  $C$  be a supervised classifier and a two-dimensional input space.

The points drawn from / belong to one of the three classes • circle/square/ triangle. After the training stage,  $C$  learns the following decision boundaries:



where the label assigned by  $C$  is given by the background color of the region in  $I$ .

Suppose that the symbols in the image are test points whose shape represents the respective ground truth label (either a circle, a square, a triangle)

which of them is True?

④ None of the answers is correct

⑤ The recall for square is  $\frac{3}{4} \neq \frac{3}{5}$

⑥ The precision for Triangle is 1

⑦ The precision for square is  $\frac{3}{5} \neq \frac{3}{4}$

⑧ The accuracy is 1 → we have some points that are not predicted correctly.  $\frac{3+8+5}{20} = \frac{16}{20} \neq 1$

$$P_{\square} = \frac{3}{4} \quad R_{\square} = \frac{3}{5}$$

$$P_{\Delta} = \frac{5}{5} = 1 \quad R_{\Delta} = \frac{5}{6}$$

$$P_{\circ} = \frac{8}{11} \quad R_{\circ} = \frac{8}{9}$$

1) (2 points)

2022-01-13

(تم 2021-01-25)

A KNN classifier is trained on the training set in the table below which only contains categorical features.

With categorical features, the distance between two samples a and b can be computed as the number of features with different values.

$$\text{dist}(a, b) = \sum_i \delta(a_i, b_i)$$

where:

$$\delta(m, n) = \begin{cases} 0 & \text{if } m=n \\ 1 & \text{otherwise} \end{cases}$$

For each neighbor  $X_i$  of  $x$  the vote is weighted:

$$\text{weight}(x, X_i) = \frac{1}{\text{dist}(x, X_i) + 1}$$

width	weight	speed	class	dist(t1)
$X_1$ big	light	fast	A	1
$X_2$ big	heavy	fast	B	0
$X_3$ small	heavy	slow	B	2
$X_4$ big	heavy	slow	A	1
$X_5$ small	light	slow	A	3

Given the test point below, write in the answer box

1- the list of the 3 neighbors of  $t1 \rightarrow$   $\text{dist}(t1, X_i)$

2- the class assigned to  $t1$  with  $k=3$

width	weight	speed	class
$t1$ big	heavy	fast	B

use the following notation

neighbors = {list of neighbors for  $t1$ }

class = class assigned

	$\text{dist}(t1)$	$\text{weight}(x, X_i)$	class
$X_1$	1	$\frac{1}{1+0} = \frac{1}{2}$	(A)
$X_2$	0	$\frac{1}{1+2} = \frac{1}{3}$	B
$X_3$	1	$\frac{1}{1+3} = \frac{1}{4}$	(A)

$$\text{similar A} = \frac{2}{3} + \frac{2}{3} = \frac{4}{3}$$

$$\text{similar B} = \frac{1}{3}$$

جواب سؤال

①

B و معاذ الله A و معاذ الله  
لهم انت ذي الصلة نستعين بغيرك  
لهم انت ذي الصلة نستعين بغيرك

• k-means, k-centroids, k-medians  
• k-means, k-centroids, k-medians

2) (1 point)

The z-score normalization is applied to a feature  $x$  by removing the mean  $\mu$  and dividing by its standard deviation  $\sigma$ . The resulting feature  $x'$  is 
$$x' = \frac{x - \mu}{\sigma}$$
  $\mu = 0$   $t_1 - \text{Mean}_{\text{train}}$   $t_1 - \sigma_{\text{train}}$

Consider the regression task described by the training set and the test set below. In which the features have not been normalized. The dataset contains three input features (Freq, rpm, power) and an output variable ( $y$ ).

The input features have been normalized with z-score (not shown in the tables before training) a regression model to predict the value of  $y$ .

You are required to perform the normalization on the test point  $t_1$  (see test set table). Write in the answer box the normalized row corresponding to the test point  $t_1$ , after applying the proper z-score normalization on the features.

training set  $\rightarrow$

	Freq	rpm	Power	$y$
$x_1$	10	20	0	0
$x_2$	5	0	2	10
$x_3$	5	0	0	10
$x_4$	10	20	2	40
mean	7.5	10	1	15
std	2.5	10	1	15

$t1 = [freq, rpm, power, y]$

✓  $freq = \frac{5.5 - 7.5}{2.5} = \frac{-2}{5}$

$rpm = \frac{1 - 10}{10} = \frac{9}{10}$

power  $s = \frac{4 - 1}{1} = 3$

$y = 3$

freq =  $\frac{t1(\text{prop}) - \text{Mean}_{\text{train}}}{\text{Std}_{\text{train}}}$

	Freq	rpm	power	Y	test
t1	5.5	1	4	3	
t2	0.5	6	8	6	
mean	3	3.5	6	4.5	
std	2.5	2.5	2	1.5	

Apply the input normalization to obtain output

$y_{\text{pred}} = y_{\text{true}} (\text{class}(A))$

تصحیف داده A در ورودی نتیجه تم ← TP

تصحیف داده A در ورودی نتیجه تم A بود ← TN

تصحیف داده A در ورودی نتیجه تم A نبود ← FP

تصحیف داده A در ورودی نتیجه تم A نبود ← FN

$y_{\text{true}} = [A \ A \ B \ C]$

$y_{\text{pred}} = \left[ \begin{array}{cccc} A & B & A & B \\ \hline \text{TP} & \text{TN} & \text{FP} & \text{FN} \end{array} \right]$

3) (2 points)

The Gini index of the node is computed as follows:

$$gini(\text{node}) = 1 - \sum_j P(j|t)^2$$

where  $P(j|t)$  is the relative frequency of class  $j$  at node  $t$

The Gini index of a split with parent  $P$  and children  $c_i$  is computed as

follows:  $gini(\text{split}) = \sum_i \frac{n_i}{n} gini(c_i)$

where  $n_i$  is the number of records at child  $c_i$  and  $n$  is the number of records in  $P$ .

In the Figure it is shown a split  $x$ , with three children (a, b, c) for each child you are given the number of elements belonging to each of three classes. ( $c_1, c_2, c_3$ )

$x$			$gini(a)$ -value	$gini(x)$ -value
a	b	c		
$c_1: 80$	$c_1: 10$	$c_1: 0$		
$c_2: 0$	$c_2: 10$	$c_2: 50$		
$c_3: 20$	$c_3: 30$	$c_3: 0$		
$n_1 = 100$	$n_2 = 50$	$n_3 = 50$		

$$Gini(a) = 1 - \left( \frac{80}{100} \right)^2 - \left( \frac{20}{100} \right)^2 = 1 - \frac{16}{25} - \frac{1}{25} = \frac{8}{25}$$

$$Gini(b) = 1 - \left( \frac{10}{50} \right)^2 - \left( \frac{10}{50} \right)^2 - \left( \frac{30}{50} \right)^2 = 1 - \frac{1}{25} - \frac{1}{25} - \frac{9}{25} = \frac{14}{25}$$

$$Gini(c) = 1 - \left( \frac{50}{50} \right)^2 = 0$$

$$Gini(x) = \frac{100}{200} \times \frac{8}{25} + \frac{50}{200} \times \frac{14}{25} = \frac{4}{25} + \frac{7}{50} = \frac{15}{50} = \frac{3}{10} = 0.3$$

(1) (2 points)

Given the similarity matrix in the figure, apply agglomerative hierarchical clustering with single (MIN) linkage. The similarity metric ranges from 0 (more distant) to 1 (less distant).

Let K be the number of obtained clusters. Write in the answer box:

1. The clusters obtained with  $K=3$

2. The clusters obtained with  $K=2$

Use the following notation:

$K=3 \{ \text{points of the first cluster} \}$

$\{ \text{points of the second cluster} \}$

$\{ \text{points of the third cluster} \}$

$K=2 \{ \text{points of the first cluster} \}$

$\{ \text{points of the second cluster} \}$

A	B	CE	D	
1	0.4	0.5	0.8	$D, E \rightarrow 1/4$
B	0.4	1	0.9	0.3
CE	0.5	0.4	1	0.5
D	0.8	0.3	0.5	1

	a	b	c	d	e
a	1	0.4	0.4	0.8	0.5
b	0.4	1	0.3	0.3	0.4
c	0.4	0.3	1	0.4	0.9
d	0.8	0.3	0.4	1	0.5
e	0.5	0.4	0.9	0.5	1

AD	B	CE	
1	1/4	1/5	$B, A \rightarrow 1/4$
B	1/4	1	0.9
CE	1/5	0.4	1

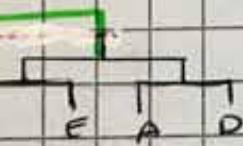
(.19) . . . . . (1/6)  
(C, E)  $\rightarrow$



(.8) . . . . . (1/6)  
(A, D)  $\rightarrow$



(.15) . . . . . (1/6)  
(AD, CE)  $\rightarrow$



ADCE	B	
1	1/4	$AD, B \rightarrow 1/4$
B	0.4	$CE, B \rightarrow 1/4$



q1)  $K=3$  clusters =  $(B, CE, AD)$

q2)  $K=2$  clusters =  $(B, ACDE)$



5) (1 point - 15) penalty)

Consider the following numpy array:

$x = \text{np.array}([[[2, 4, 6, 8], [3, 6, 9, 12], [4, 8, 12, 16]]])$

The following indexing techniques are applied to obtain a and b

$a = x[1:, ::-1]$  start from the end

$$a = [12, 9, 6, 3]$$

$b = x[[\text{False}, \text{True}, \text{True}], :0:-1]$

$$b = [12, 9, 6]$$

which of the following assertion is correct?

- a) The script will provide an error because b cannot be computed.
- b) The value of a is  $[[3, 6, 9, 12], [4, 8, 12, 16]]$ . The variables a & b contain the same values.
- c) The value of a is  $[[12, 9, 6, 3], [16, 12, 8, 4]]$ . The variables a & b contain different values.
- d) The values of a is  $[[3, 6, 9, 12], [4, 8, 12, 16]]$ . The variables a & b contain different values.
- e) The value of a is  $[[12, 9, 6, 3], [16, 12, 8, 4]]$ . The variables a & b contain the same values.
- f) None of the other answers is correct.

6) (2.5 points)

Silhouette

- range between -1 and +1
- \*typically between 0 and 1
- \*the closer to 1 the better

for each data point  $P_i$

- let  $\text{intra}(P_i)$  be the average dissimilarity of  $P_i$  with all other points within the same cluster
- let  $\text{inter}(P_i)$  be the lowest average dissimilarity of  $P_i$  to any other cluster of which  $P_i$  is not a member.

The Silhouette score of point  $P_i$  is defined as:

$$\text{Silh}(P_i) = \frac{\text{inter}(P_i) - \text{intra}(P_i)}{\max(\text{inter}(P_i), \text{intra}(P_i))}$$

For two  $n$ -dimensional points  $a = (a_1, a_2, \dots, a_n)$  and  $b = (b_1, b_2, \dots, b_n)$

the Manhattan distance is defined as follows

$$d(a, b) = \sum_{i=1}^n |a_i - b_i|$$

in the figure below, three different clusters are represented with the following shapes: circle, star, square. Use the Manhattan distance as the dissimilarity metric.

consider point  $i$  belonging to the square cluster

$$\text{intra}(i) = \frac{1}{5} (1+1+2+3+3) = 2$$

$$\text{inter}(i) = \min(d(i, a), d(i, b), d(i, c), d(i, d), d(i, e), d(i, f), d(i, g))$$

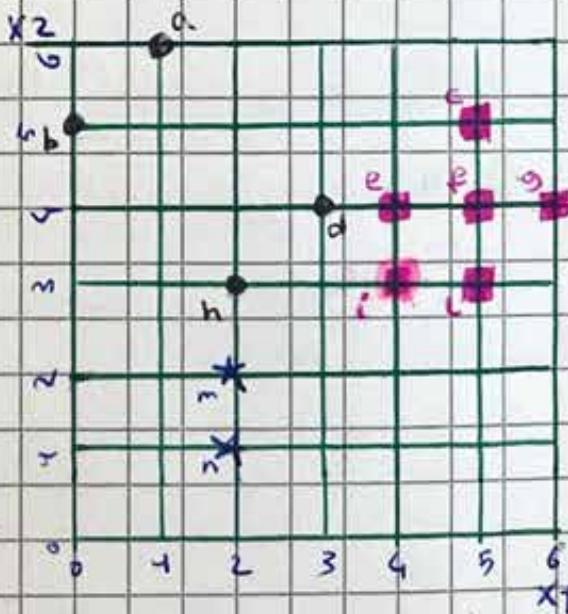
$$\text{silh}(i) = \text{value}$$

notation:

$\text{intra}(i) = \text{value}$

$\text{inter}(i) = \text{value}$

$\text{silh}(i) = \text{value}$



$$\text{inter}(i) = \min(d(i, a), d(i, b), d(i, c), d(i, d), d(i, e), d(i, f), d(i, g)) = \frac{1}{7} \sum_{j \in \text{cluster}} d(i, j)$$

$$d(i, a) = \frac{1}{5} (1+1+2+3+3) = 2$$

$$d(i, b) = \frac{1}{5} (2+2+6+6) = 4$$

$$d(i, c) = \frac{1}{5} (3+4) = 3.5 \rightarrow \text{inter}(i) = 3.5$$

$$\text{silh}(i) = \frac{3.5 - 2}{3.5} = \frac{3}{7}$$

F7 (1.5 points - 15% penalty) ✓

An itemset  $I$  is closed if none of its immediate supersets has the same support as  $I$ .

An itemset  $I$  is represented by a collection of literals (e.g. abc) and a number representing its support count (e.g. abc: 12).

After the analysis of a transactional dataset, the list of frequent itemsets found (with support counts) is the following:

length = 1, a: 150 b: 160 c: 180 d: 150 e: 150

length = 2, ab: 140 ac: 120 ad: 130 ae: 140

length = 3 abc: 120 acd: 110 abe: 130

length = 4, abce: 120

which of the following statement is true?

- a) The only non-closed itemsets are ac, ad, abc, abce
  - b) The only non-closed itemsets are ac, ad
  - c) The only non-closed itemsets are ac, ad
  - d) None of the other answers is correct
  - e) The only closed itemsets are ac, abc
- ② The only non-closed itemsets are ac, abc

		<u>closed</u>		<u>non-closed</u>
4	120	ABCE		3 120 ABC
3	110	ACD		2 120 AC
3	130	ABE		
2	140	AB		
3	130	AD		
2	140	AE		
1	150	A		
1	160	B		
1	180	C		
1	150	D		
1	190	E		

↓↓↓↓↓      sup      up  
                  count

8) (1 point - 15% penalty)



- **Precision (c)** is the fraction of correct predictions among the samples predicted with class c
- **Recall (c)** is the fraction of correct predictions among the samples with actual class c

after training two models ( $M_1, M_2$ ) we obtain the **confusion matrices** shown below.

consider **class b** only. Which of the following statements is true?

	$M_1$	a	b	c	d
true class	a	15	10	1	0
	b	0	50	5	5
	c	5	20	70	0
	d	5	10	0	60

Predicted class

	$M_2$	a	b	c	d
true class	a	14	10	0	2
	b	5	40	5	10
	c	5	5	80	5
	d	10	5	10	50

- (a)  $M_1$  has higher recall than  $M_2$  and lower precision
- (b)  $M_1$  has the same recall as  $M_2$  and lower precision
- (c) None of the other answers is correct
- (d)  $M_1$  has lower recall than  $M_2$  and higher precision
- (e)  $M_1$  has higher recall than  $M_2$  and higher precision
- (f)  $M_1$  has lower recall than  $M_2$  and lower precision

$$\begin{array}{cc} P & R \\ M_1 & \frac{50}{90} = \frac{5}{9} & \frac{50}{60} = \frac{5}{6} \end{array}$$

Λ

✓

$$\begin{array}{cc} P & R \\ M_2 & \frac{40}{60} = \frac{2}{3} & \frac{40}{60} = \frac{2}{3} \end{array}$$

$$\frac{2}{3} > \frac{5}{9}$$

$$\frac{2}{3} > \frac{5}{6}$$

9) (1.5 points, -15% penalty)

✓ ✓ ✓

consider the following Python operations:

$L_1 = [1, 2, 3]$

$d = \{ 'mon': [1, 2, 3], 'tue': [2, 4, 6] \}$

$L_2 = [2, 4, 6]$

$[1, 2, 3, 4]$

$[2, 4, 6, 8]$

$d = \{ 'mon': L_1, 'tue': L_2 \}$

$res = d.value()$

$L_1.append(4)$

$d['tue'].append(8)$

$print(res)$

●  $\text{dict\_value}([L_1, L_2], [L_3, L_4])$  ✓

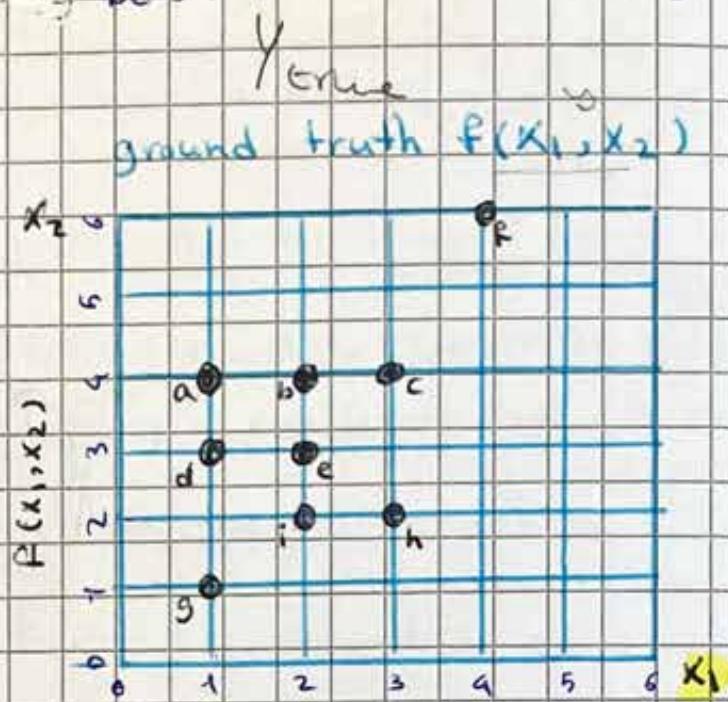
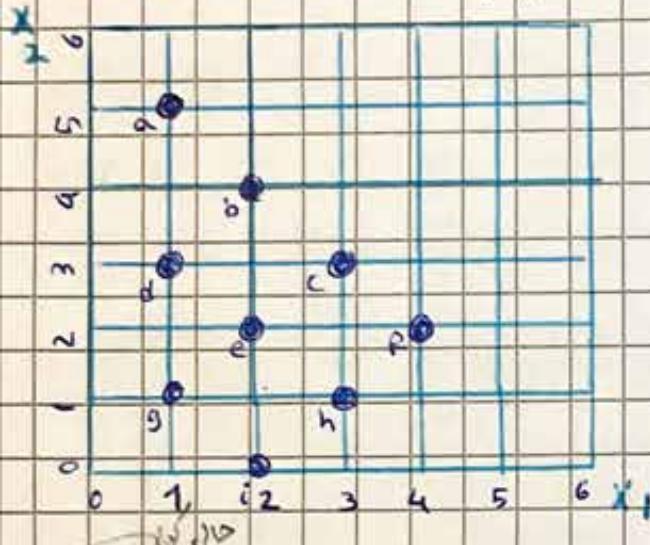
40) (1.5 points - 15% penalty)



The Mean Squared Error between two vectors  $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^n$  is defined as

$$MSE(\mathbf{z}, \mathbf{z}') = \frac{1}{n} \sum_i (z_i - z'_i)^2$$

Consider the points in the  $\mathbb{R}^2$  below with features  $x_1, x_2$  and ground truth values  $f(x_1, x_2)$  input space prediction  $y_{\text{pred}}$



Apply the following regression model to all the data points:

$$f'(x_1, x_2) = 0.5x_1 + 0.5x_2 + 1$$

Write in the answer box the value of the MSE between the prediction  $f'(x_1, x_2)$  and the ground truth  $f(x_1, x_2)$

(with 2 decimal places)

$$X_1 \quad X_2 \quad f(x_1, x_2) \quad f'(x_1, x_2) = 0.5x_1 + 0.5x_2 + 1$$

$$f(x_1, x_2) - f'(x_1, x_2)$$

	$x_1$	$x_2$	$f(x_1, x_2)$	$f'(x_1, x_2)$	$\frac{1}{n} \sum (y - \hat{y})^2$
a	1	5	4	$\frac{5}{2}$	4.5
b	2	4	4	4	0
c	3	3	4	4	0
d	1	3	3	3	0
e	2	2	3	3	0
f	4	2	6	6	0
g	1	1	1	1	0
h	3	1	2	2	0
i	2	0	2	2	0
	$\bar{x}$	$\bar{y}$			

Q2 (2 points)  
 An itemset is closed if none of its immediate supersets has the same support as the itemset. Given the transactional dataset shown in the figure below, apply the apriori algorithm to extract all frequent itemsets. The value of minsup is 2. Support  $\geq$  minsup  
 an itemset is considered to be frequent if its support count is equal to or higher than the minsup.

Q1) list all frequent itemsets having length 2, along with their support count.

Q2) list all itemsets of length 3 that have been generated by Apriori after the join and prune steps, before counting their support in the database.

Q3) list all frequent itemsets that are not closed, along with their support count.

in A1) { list of itemsets w/ support count }

		L1	L2	L3			
transactions	itemsets	sup	itemsets	sup	itemset	sup	
0	BC	close + A	4	AB	2	ABC	x
1	AE	close - B	7	AC	1	ABD	1
2	ACD	- C	6	AD	2	ABE	x
3	AB	- D	9	AE	1	ACE	x
4	BE	- E	4	BC	5	ACD	x
5	DE		not close	BD	2	ADE	x
6	BCE			BG	2	BCD	2
7	BC			CD	3	BCE	x
8	BCD		close	CE	1	BDE	x
9	ABCD			DE	1	CDE	x

Support 2,2

A1) { AB: 2, AD: 2, BC: 5, BD: 2, BE: 2, CD: 3 }

A2) { BCD: 2 } ABD, before counting sup? 1, 2, 3, 4, 5, 6, 7, 8, 9

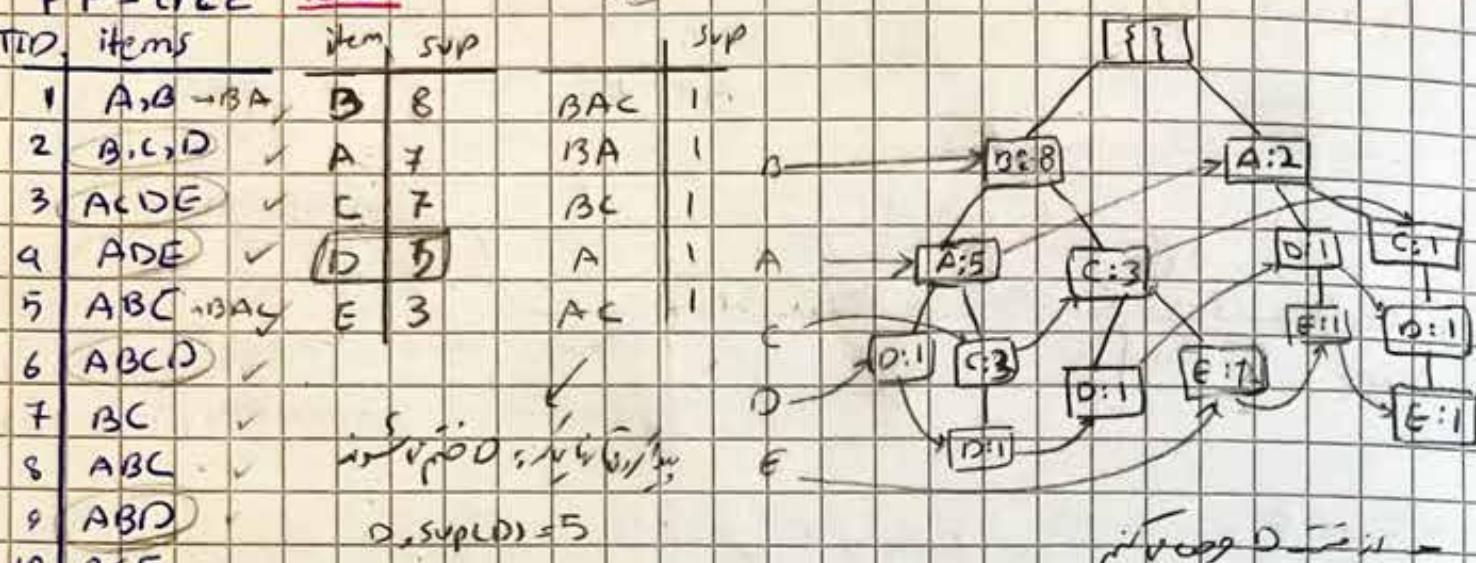
A3) { AC: 1, AE: 1, CG: 1, DE: 1, ABD: 1 }

# FP-tree

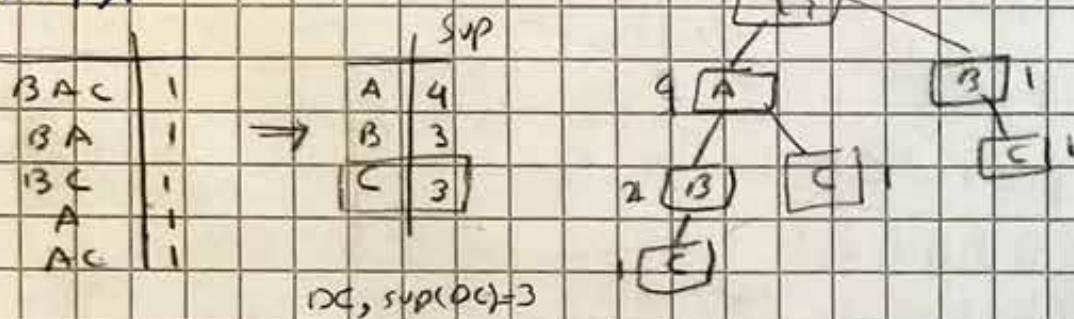
DB

TID	items	item, sup
1	A,B → BA,	B 8
2	B,C,D	A 7
3	A,C,D,E	C 7
4	A,D,E	D 5
5	A,B,C → B,A,Y	E 3
6	A,B,C,D	
7	B,C	
8	A,B,C	
9	A,B,D	
10	B,C,E	

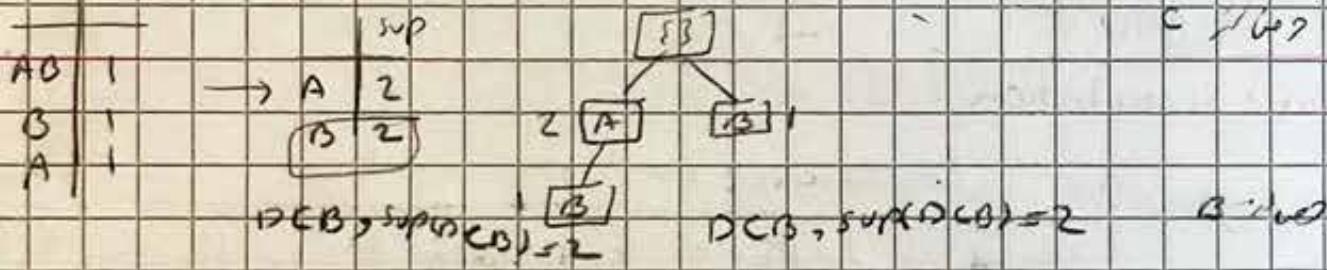
sup



minsup > 1



DC, sup(DC)=3



DCB, sup(DCB)=2

DCA, sup(DCA)=2

مقدمة بحث متقدم - د. جعفر عباس - ٢٠١٧  
العنوان: فضاءات FP - tree  
العنوان: فضاءات FP - tree

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

Euclidean distance مسافة أوروبية

$$P | \begin{matrix} x \\ y \end{matrix}$$

$$d(x, y) = \left( \sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

Minkowski dist

Cubic metric if r=3, r=2 if

Similarity Between Binary Vectors:

$$SMC = \frac{M_{11} + M_{00}}{M_{00} + M_{01} + M_{10} + M_{11}} = \frac{\text{number of matches}}{\text{number of attributes}}$$

Binary  
vectors

$$M_{00} = P_0 Q_0 / M_{01} = P_0 Q_1 / M_{10} = P_1 Q_0 / M_{11} = P_1 Q_1$$

$$(SMC \text{ Versus Jaccard}) \quad J = \frac{M_{11}}{M_{01} + M_{10} + M_{11}} \rightarrow \text{no non-zero entries}$$

Cosine Similarity

$$\cos(d_1, d_2) = \frac{(d_1 \cdot d_2)}{\|d_1\| \|d_2\|}$$

$$\|d_1\| = (\bar{d}_{11} \times \bar{d}_{11} + \bar{d}_{12} \times \bar{d}_{12} + \dots)^{1/2}$$

$$\|d_1\| = (\bar{d}_1 \times \bar{d}_1)^{1/2}$$

$$-1 \leq \cos \leq 1$$

Pearson's correlation

$$\text{corr}(x, y) = \frac{\text{covariance}(x, y)}{\text{standard deviation}(x) * \text{standard deviation}(y)} = \frac{S_{xy}}{S_x S_y}$$

$$\text{covariance}(x, y) = S_{xy} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$$

$$\text{standard deviation}(x) = S_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2}$$

$$\text{standard deviation}(y) = S_y = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$\begin{aligned} d_1 &= 000000 \\ d_2 &= 000000 \end{aligned}$$

$$\cos(\theta) = \frac{x_A x_B + y_A y_B}{\sqrt{x_A^2 + y_A^2} \times \sqrt{x_B^2 + y_B^2}}$$

$$0 \leq \cos \leq 1$$

$$\begin{array}{c} 0110 \\ 0110 \\ 0110 \\ 0110 \end{array}$$

لما زادت ميله من  $\bar{x}$   
لما زادت ميله من  $\bar{y}$

## Accuracy Score :

جودة التصنيف

accuracy

: classification

تصنيف صحيح

Correct

تصنيف خاطئ

incorrect

تصنيف خاطئ

measure

$2P \cdot R$

$P + R$

مقدار الدقة

مقدار الدقة

= (Percision) P

Recall =  $\frac{\text{No of objects correctly assigned to } C}{\text{No of obj belonging to } C}$

(Recall) R

Precision =  $\frac{\text{No of obj correctly assigned to } C}{\text{No of obj assigned to } C}$

مجموع العناصر المطلوبة

مجموع العناصر المطلوبة

## Evaluation Regression

$$MAE \text{ (Mean Absolute Error)} = \frac{1}{n} \sum_i^n |y_i - \hat{y}_i|$$

$$MSE \text{ (Mean Squared Error)} = \frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2$$

(tend to penalize less errors close to 0)

$$R^2 = 1 - \frac{MSE}{SST^2} \quad R^2 = \begin{cases} \approx 1 & \rightarrow \text{good prediction} \\ < 0 & \rightarrow \text{wrong prediction} \end{cases}$$

## Minkowski Distance

$$d(x, y) = \left( \sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

-  $r=1$  Manhattan distance ( $L_1$  norm)

-  $r=2$  Euclidean distance

-  $r \rightarrow \infty$  supremum ( $L_{\max}$  norm,  $L_\infty$  norm)

$\frac{p_1}{p_1} \times \frac{p_2}{p_2} \dots \frac{p_n}{p_n} \Rightarrow \frac{p_1}{p_1}$

$$\text{Euclidean : } \sqrt{\sum (x_i - y_i)^2} = d(x, y)$$

$\text{shape}(X) = (100, 50)$

$\text{shape}(y) = (1, 50)$

$$\text{np.sqrt}((x - y) * * 2).sum (axis = 1)$$

$(1, 100) = \text{euclidean dist}$

## Gini index (classification)

$$n = n_1 + n_2 + n_3$$

$$\text{gini}_a = 1 - \left( \frac{n_1}{n_1} \right)^2 + \left( \frac{n_2}{n_2} \right)^2 \dots$$

$$\text{gini}_b = 1 - \left( \frac{n_1}{n_2} \right)^2 + \left( \frac{n_3}{n_3} \right)^2 \dots$$

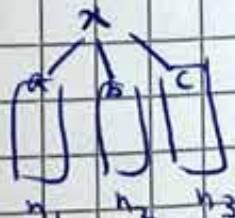
$$\text{gini}_c = 1 - \left( \frac{n_1}{n_3} \right)^2 + \left( \frac{n_2}{n_3} \right)^2 \dots$$

$$\Rightarrow \text{Gini} = \frac{n_1}{n} \times \text{gini}_a + \frac{n_2}{n} \times \text{gini}_b + \frac{n_3}{n} \times \text{gini}_c$$

نسبة العناصر في كل فرع

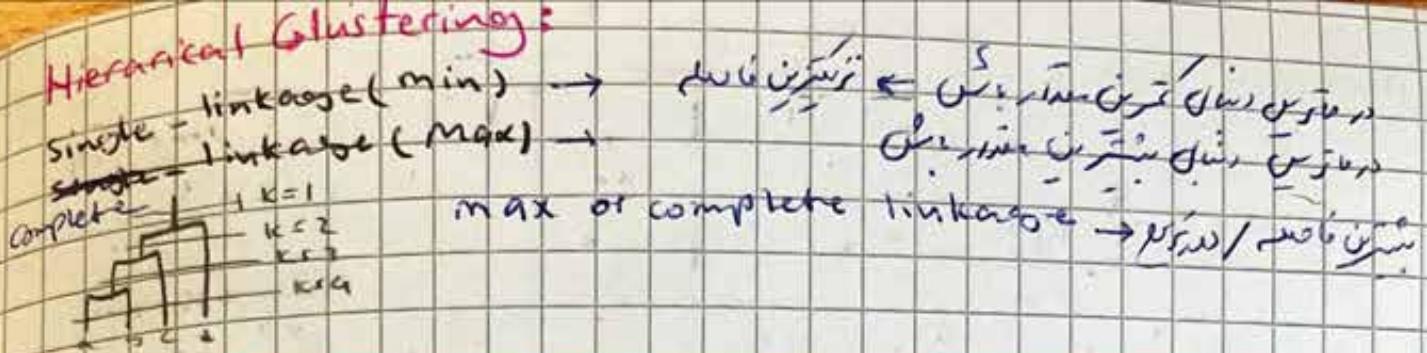
الكتل التي تدخل في كل فرع

y



✓

## Hierarchical Clustering:



## Naive Bayes classifier

$$x \in \mathcal{X}^N \quad x = (x_1, x_2, \dots, x_n)$$

$$p(\text{class} | x) = p(x_1 | \text{class}) * p(x_2 | \text{class}) * \dots * p(x_n | \text{class}) * p(\text{class})$$

↓  
خواص  
جوابی  
جوابی

Naive Bayes  
Classification

$$\text{Manhattan distance} \quad \text{dist}(P_1, P_2) = |x_1 - x_2| + |y_1 - y_2|$$

DBSCAN - core : Points with at least minpoints neighbors within a distance  $\epsilon$   
- border : at least a core point in the neighborhood  
- noise : neither core nor border

support

$$A \Rightarrow B$$

confidence

$$\text{sup} = \frac{\text{itemsets}}{\text{total items}} = \frac{\# \{A, B\}}{|T|} \quad \text{conf} = \frac{\text{itemset}}{\text{total items}} = \frac{\sup(A, B)}{\sup(A)}$$

$\rightarrow \text{sup} \geq \text{minsup}$

$\text{conf} \geq \text{minconf}$

(Frequent) Itemset  $\emptyset$  is an itemset whose support is greater than or equal

Rand Index  $RI_{th}$  to a minsup threshold

$$\text{Rand Index} = \frac{P_{00} + P_{11}}{P_{00} + P_{01} + P_{10} + P_{11}}$$

$$= \frac{TN + TP}{TN + FP + TP + FN}$$

pred [ 0, 1 ]

True [ 0, 1 ]

$$Y_p \quad TP / TN$$

مقدار این مقدار

[ 0, 1 ]

[ 0, 1 ]

[ 0, 1 ]

[ 0, 1 ]

[ 0, 1 ]

[ 0, 1 ]

$$TP = 1 \leftarrow Y_{true} = 1, Y_p = 1$$

$$TN = 1 \leftarrow Y_{true} = 0, Y_p = 0$$

$$FP = 1 \leftarrow Y_{true} = 0, Y_p = 1$$

$$FN = 1 \leftarrow Y_{true} = 1, Y_p = 0$$

تصویر		تصویر	
تصویر	تصویر	تصویر	تصویر
FP	TP	بررسی	ا
P <sub>01</sub>	P <sub>11</sub>		

تصویر		تصویر	
تصویر	تصویر	تصویر	تصویر
TN	FN	بررسی	ا
P <sub>00</sub>	P <sub>10</sub>		

$$0.5N = TP \leftarrow Y_{true} = 1, Y_p = 1$$

$$0.5N = FN \leftarrow Y_{true} = 1, Y_p = 0$$

$$0.5N = FP \leftarrow Y_{true} = 0, Y_p = 1$$

$$0.5N = TN \leftarrow Y_{true} = 0, Y_p = 0$$

سیستم hyper parameters (متغیر RI,)

## Apriori Algorithm

	Example DB	$L_1$	$L_2$	$L_3$	$L_3$
1	{A, B}	1 A	1 AB	1 A, B	1 A, B, C
2	{B, C, D}	2 B	2 BC	2 A, C	2 A, B, D
3	{A, C, D, E}	3 C	3 CD	3 A, D	3 A, B, E
4	{A, D, E}	4 D	4 DE	4 A, E	4 A, C, D
5	{A, B, C}	5 E	5 BE	5 BC	5 A, C, E
6	{A, B, C, D} $\Rightarrow$			6 BD	6 A, D, E
7	{B, C}			7 BE	7 B, C, D
8	{A, B, C}			8 CD	8 B, C, E
9	{A, B, D}			9 CE	9 B, D, E
10	{B, C, E}			10 DF	10 C, D, E

$\text{minSup} > 1$

\* An Itemset is closed if none of its immediate supersets has the same support as the itemset.

$L_4$

items	sup
A, B, C, D	1
A, B, D, E	
A, C, D, E	

$\Rightarrow L_4 = \emptyset$

$y_{true} [ \quad , \quad , \quad , \quad , \quad ]$

$y_{pred} [ \quad , \quad , \quad , \quad , \quad ]$

accuracy  $\frac{\text{true}}{\text{true} + \text{false}}$

Precision  $(A) = \frac{\text{true}}{\text{true} + \text{false}}$

recall  $s$

$A \rightarrow \text{true}$

$y_{true} \rightarrow A$

		Predicted Class	
		+	-
Actual Class	+	$P_+$ (TP)	$P_-$ (FN)
	-	$P_-$ (FP)	$P_{-+}$ (TN)

← confusion Matrix  
(binary classifier)

$$\text{accuracy} = \frac{TP + TN}{TP + FN + FP + TN}$$

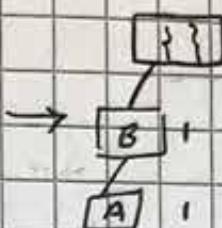
for positive class :  $\text{Precision}(P) = \frac{a}{a+c} = \frac{TP}{TP+FP}$

$$\text{Recall}(R) = \frac{a}{a+b} = \frac{TP}{TP+FN}$$

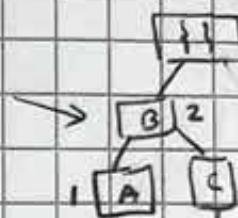
$$\text{F-measure}(F) = \frac{2RP}{R+P} = \frac{2a}{2a+b+c}$$

## FP-tree:

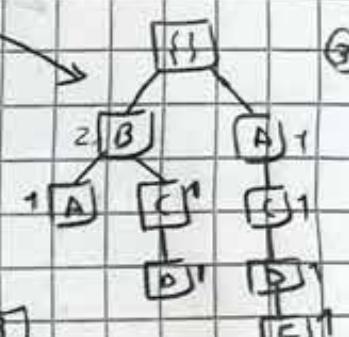
	items	sup	
1	B	8	$\Rightarrow \sqcup A, B$
2	A	7	
3	C	7	
4	D	5	
5	E	3	
6			$\downarrow$
7	B, C		$\sqcup B, C, D$
8	A, B, C		$\downarrow$
9	A, B, D		$\downarrow$
10	B, C, E		



①

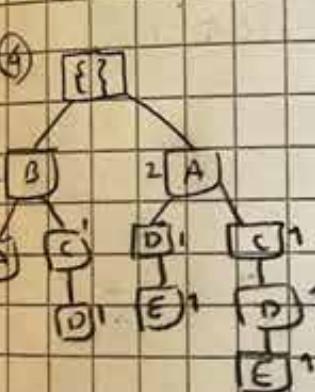


②

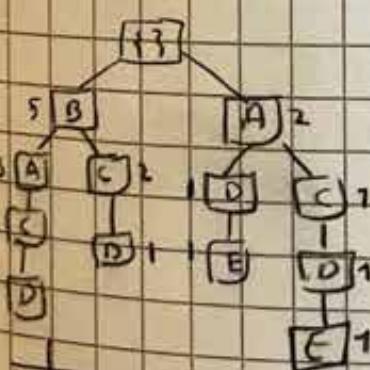


③

4 | A, D, E

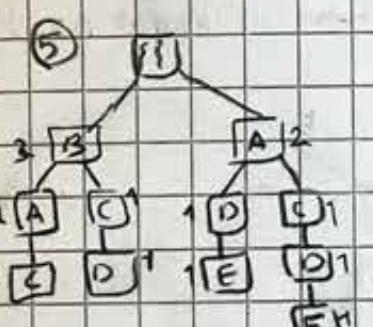


7 | B, C

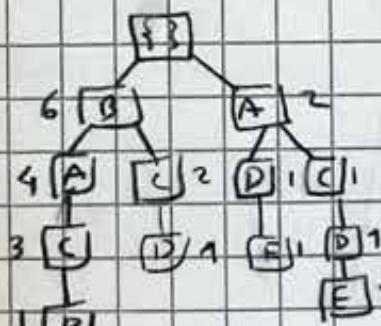


10 | B, C, E

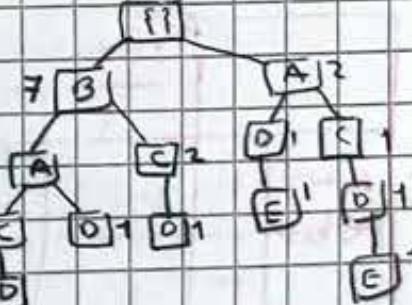
5 | B, A, C



8 | B, A, C

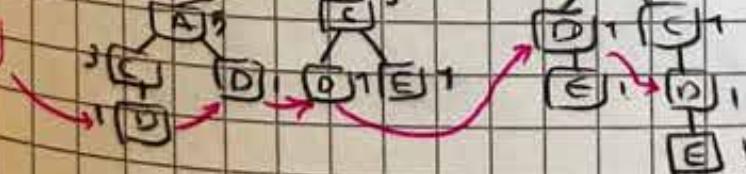


9 | B, A, D



Conditional pattern Base:  $\sqcup D$

$D, \text{sup}(D) = 5$



\* فرم D<sub>i</sub>/items sup

	items	sup	A العلامة
BAC	1	4	
BA	1	3	ازول
BC	1	3	حروف ممكنا
A	1		
AC	1		این کلمه ایست

$DC, sup(DC) = 3$  - این کلمه ایست

	items	sup	A العلامة
1 A B	1		
2 A	1		
3 B	1		

$DCB, sup(DCB) = 2$  - این کلمه ایست

empty = DCB - CPB ، باید از مجموع داشت شود ، ازینجا تعقیب نموده است - ایم A باید از مجموع داشت شود ، ایم A و B در پایان

$DCA, sup(DCA) = 2$

	items	sup	A العلامة
1 A	1		
2 [A] B	1		

$DCA - CPB = empty$

حال برسی کنم عصب  $\rightarrow$  (A, B, C) برداشته شده  $\rightarrow$  (A, B) برداشته شده  $\rightarrow$  (A)

	items	sup	A العلامة
BAC	1		
BA	1		
BC	1		
A	1		
AC	1		

$DB, sup(DB) = 3$

must  
be  
norm  
alized

✓ Correlation < Pearson →  $r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$  / linear corr if  $r \approx 1$   
 Spearman <  $r_s = \frac{1}{n(n-1)} \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{(r_{ij} - \bar{r})^2}{\sigma_r^2}$  / 1: 2 var correlated by an increasing monotonic func  
 -1: Function is decreasing monotonic  
 0: there isn't a monotonic func correlation

## Support Vector Regression

$$MAE = \frac{1}{n} \sum |y_i - \hat{y}_i|$$

$$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$RSE = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$
number of sample:  $n$

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{MSE}{\sigma^2 (std)^2}$$

$$RSS = \sum_i (y_i - \hat{y}_i)^2$$

$$TSS = \sum_i (y_i - \bar{y}_i)^2$$

$$\bar{y} = \frac{1}{n} \sum_i y_i$$

predicted value:  $\hat{y}$  / actual value:  $y$

In supervised learning, it's better to use MAE, MSE

$R^2 = 1$  perfect linear relationship between  $x$  &  $y$

$R^2 = \begin{cases} 1 & \text{perfect linear relationship between } x \text{ & } y \\ \text{root of the } y \text{ variation is explained by variation in } x & \end{cases}$

$R^2 = r^2 \text{ close to } 1$  very good linear relationship between  $x$  &  $y$

good prediction

$\cdot (R^2)$  weaker linear relationship  
 A portion of the variation in  $y$  is not explained by variation in  $x$

$\cdot R^2 \approx 0$  no linear relationship  
 The value of  $y$  doesn't depend on the value of  $x$

It's upper bound + 1

It has no lower bound (-∞)

It is used as a regression evaluation metric

Micro average scores:

$$\text{micro-P} = \frac{\text{total\_TP}}{\text{total\_TP} + \text{total\_FP}}$$

$$\text{micro-R} = \frac{\text{total\_TP}}{\text{total\_TP} + \text{total\_FN}}$$

$$\text{micro-F1} = \text{micro-P} = \text{micro-R}$$

Correlation or lift:

$$\text{correlation} = \frac{P(A, B)}{P(A)P(B)} = \frac{\text{Conf}(r)}{\text{Sup}(B)}$$

$r: A \rightarrow B$

$$\text{correlation} \begin{cases} > 1 & \rightarrow \text{positive correlation} \\ = 1 & \rightarrow \text{statistical independence} \\ < 1 & \rightarrow \text{negative correlation} \end{cases}$$

$$\text{lift}(A, B) = \frac{C(A \rightarrow B)}{\text{Sup}(B)} \leq \frac{S(A \cup B)}{S(A) \times S(B)}$$

ROC

$$\hookrightarrow \text{TPR} \text{ (True Positive Rate)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

(on the y-axis)

$$\hookrightarrow \text{FPR} \text{ (False Positive Rate)} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

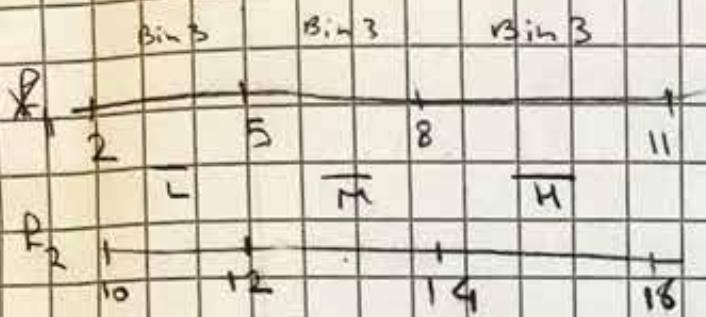
(on the x-axis)

The output of TF-IDF is:

A sparse matrix, since most terms will not occur in most documents t.

$$w_R = \frac{11-2}{3} = \frac{9}{3} = 3$$

$$w_{F_2} = \frac{16-10}{3} = \frac{6}{3} = 2$$



(P1, P2)

	P <sub>1</sub>	P <sub>2</sub>
x <sub>1</sub>	6	M
x <sub>2</sub>	2	L
x <sub>3</sub>	2	L
x <sub>4</sub>	11	M
	16	H

$$x_1 = [M, M]$$

$$x_2 = [L, L]$$

(V) ✓

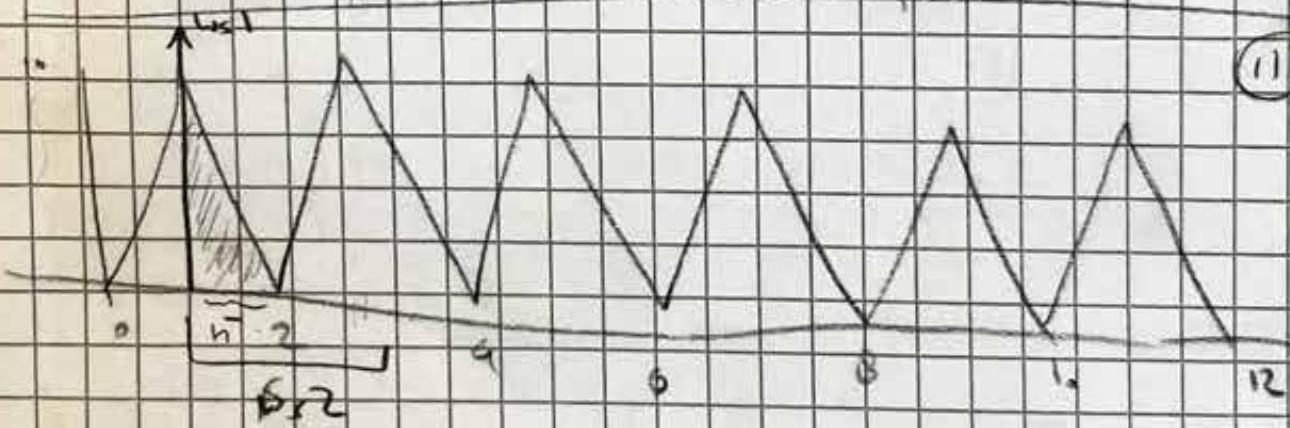
age | Gender | occupation

0	23	M	Student
1	27	M	Consult
2	34	F	manager
3	49	F	manager
4	54	M	manager
5	29	M	Freelancer
6	21	F	Student

Student	M	1	21
Consult	M	1	-
Freelancer	M	1	-
manager	F	2	34
	M	1	49
		54	→ means

$$\frac{34+49+54}{3} = 49$$

(11) ✓



Step & distance you move at each iteration

w = window size should be the piece you observe

w=7 s=1 df. arr. mean() is 5

1.11/1.10

		alpha area	days	C <sub>1</sub>	C <sub>1</sub>	C <sub>2</sub>	C <sub>2</sub>	det <sub>11</sub>	det <sub>12</sub>	ans.
x <sub>1</sub>	high	0.7	5	wed	80	0	5	15	$\frac{80}{85} + 0$	$\frac{5}{85} + \frac{15}{85}$
x <sub>2</sub>	low	0.8	100	mon	30				$\frac{2}{32} + \frac{2}{32}$	$\frac{2}{32} + \frac{2}{32}$
x <sub>3</sub>	low	0.4	50	Fri	30	5	2	70	$\frac{30+5}{32+5}$	$\frac{70}{32+5}$
x <sub>4</sub>	high	0.6	20	wed	30	0	5	15	$\frac{30}{32+5}$	$\frac{15}{32+5}$
					30	0	5	15	$\frac{2}{10} + \frac{2}{15}$	$\frac{2}{10} + \frac{2}{15}$
					30	0	5	15	$= 1$	$= 1$
					30	0	5	15	$\frac{2}{2} + \frac{2}{2}$	$\frac{2}{2} + \frac{2}{2}$
					30	0	5	15	$\frac{2}{2}$	$\frac{2}{2}$

[C<sub>2</sub>, C<sub>1</sub>, C<sub>2</sub>, C<sub>1</sub>]

$$\begin{array}{l}
 \begin{matrix} x_1 & x_2 \\ \alpha & 1 & 2 \\ x_1 & 2 & 1 \\ x_2 & 1 & 2 \\ c & 4 & 1 \\ d & 2 & 3 \end{matrix} \times \begin{bmatrix} x_1 & x_2 & x_1x_2 & x_1^2 & x_2^2 \\ 1 & 2 & 2 & 1 & 4 \\ 2 & 1 & 2 & 4 & 1 \\ 4 & 1 & 4 & 16 & 1 \\ 2 & 3 & 6 & 4 & 9 \end{bmatrix} = \begin{bmatrix} 1-2+1 \\ 2-1+1 \\ 4-1+2 \\ 2-3+3 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \\ 5 \\ 2 \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ 4 \\ 1 \end{bmatrix} = Y_{\text{pred}}
 \end{array}$$

$$\begin{array}{l}
 Y_{\text{true}} = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 \\ 1 & 3 & 2 & 3 & 2 \end{bmatrix} \\
 Y_{\text{pred}} = \begin{bmatrix} 1 & 1 & 2 & 3 & 2 \\ - & - & - & - & - \end{bmatrix}
 \end{array}$$

	Y <sub>true</sub>	Y <sub>pred</sub>	TP	TN	FP	FN
(0, 1)	0	1	0	0	1	1
(-, 2)	0	0	0	1	0	1
(-, 3)	0	0	0	1	0	1
(+, 4)	0	0	0	1	0	1
(1, 2)	0	0	0	1	0	1
(1, 3)	1	0	0	1	0	1
(1, 4)	0	0	0	1	0	1
2, 3	1	0	0	1	0	1
2, 4	1	1	1	1	0	0
3, 4	0	1	0	1	1	0
			1	3	1	1

$$FMs = \frac{1}{(1+1)(1+1)} = \frac{1}{4}$$

$$\begin{cases} Y_{\text{true}} = 1 : FN \\ Y_{\text{pred}} = 1 : FP \end{cases} \star$$

$$X_5 = \begin{bmatrix} 2 & 4 & 6 & 8 \\ 3 & 6 & 9 & 12 \\ 4 & 8 & 12 & 16 \end{bmatrix}$$

$$Z = \begin{bmatrix} 1 & 2 \\ 4 & 8 \end{bmatrix} \text{ out}$$

(9 JLR)

$$Y_5 = [1 \ 2]$$

$$\begin{bmatrix} 2 & 4 & 6 & 8 \\ 4 & 8 & 12 & 16 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 4 & 8 \end{bmatrix} \times \begin{bmatrix} 3 & 4 \\ 7 & 10 \end{bmatrix}$$

$$y = [1 \ 6 \ 8]$$

$n=3$

2020-9-2

10:10

$$\hat{y}_{\text{pred}1} = [1 \ 5 \ 6]$$

$$\hat{y}_{\text{pred}2} = [1 \ 6 \ 6]$$

$$\theta = [1/3 \ -2/3 \ 1 \ -1/3]$$

$$Q_2 = [4/3 \ -6/3 \ 2/3 \ -3/3]$$

MSE ( $y, \hat{y}_{\text{pred}}$ )

$$\frac{1}{3} \left[ (1-1)^2 + (6-5)^2 + (8-6)^2 \right] + \left( \frac{1}{3} + \frac{-2}{3} + \frac{3}{3} + \frac{-1}{3} \right)$$

$$\frac{5}{3} \leftarrow \frac{7}{3}, \quad \textcircled{81} = (Y - \hat{Y}_{\text{pred}})^2$$

$$(Y - \hat{Y}_{\text{pred}2}) / (\theta_2) =$$

$$\text{inter}(b) = \frac{1}{2} (1+3) = \frac{4}{2}$$

(20')

$$\text{inter}(e) = \frac{1}{1} (1) = 1$$

$$\text{extra}(b), (C_2 \text{ 以左}) = \frac{1}{2} (4+8) = 6 \quad (C_3 \text{ 以左}) = \frac{1}{2} (6+8) = 7$$

$$\text{sill}(b) = \frac{4+7}{\max(4, 7)} = \frac{4}{6} = \frac{2}{3}$$

$$\text{extra}(e), (C_2 \text{ 以左}) = \frac{1}{2} (7+9) = 8 \quad (C_3 \text{ 以左}) = \frac{1}{3} (6+6+6) = 6$$

$$\text{sill}(e) = \frac{8-1}{\max(6, 1)} = \frac{7}{6}$$

(2022, 1-م دیم ✓)

$x_0$	$x_1$	$x_2$	$x_0-a$	$x_0-c$	$x_1-a$	$x_1-b$	$x_1-c$	$x_2-x_1$	$x_2-y$	$x_2-z$
c	b	x	$x_1$	0	1	0	1	0	1	0
a	a	z	$x_2$	1	0	1	0	0	0	1
c	c	y	$x_3$	0	1	0	0	1	0	0
a	a	y	$x_4$	1	0	1	0	0	0	0

$$j(x_1, x_2) = \frac{M_{11}(x_1, x_2)}{M_{11}(x_1, x_2) + M_{01}(x_1, x_2) + M_{10}(x_1, x_2)}$$

$$j(x_1, x_2) = \frac{0}{0+3+3} = 0 \quad [M_{11}=0 \Rightarrow M_{00}=1, M_{10}=3, M_{01}=3]$$

$$j(x_3, x_4) = \frac{1}{2+2+1} = \frac{1}{5} \quad [M_{11}=1, M_{00}=3, M_{10}=2, M_{01}=2]$$

	SUP	SUP
A CD		
CD	A 7	C 8
AC	B 5	D 8
C	C 8	A 7
BD	D 8	B 5
ABCD	E 2	E 2

~~(CDA)~~ →

B - CPB

B\_CHT ✓

CDA	1	ANC	A 4
CA	1	D 3	
DA	1	C 2	
A	1		

(A2)

(A1)

CDE

ABC

AC

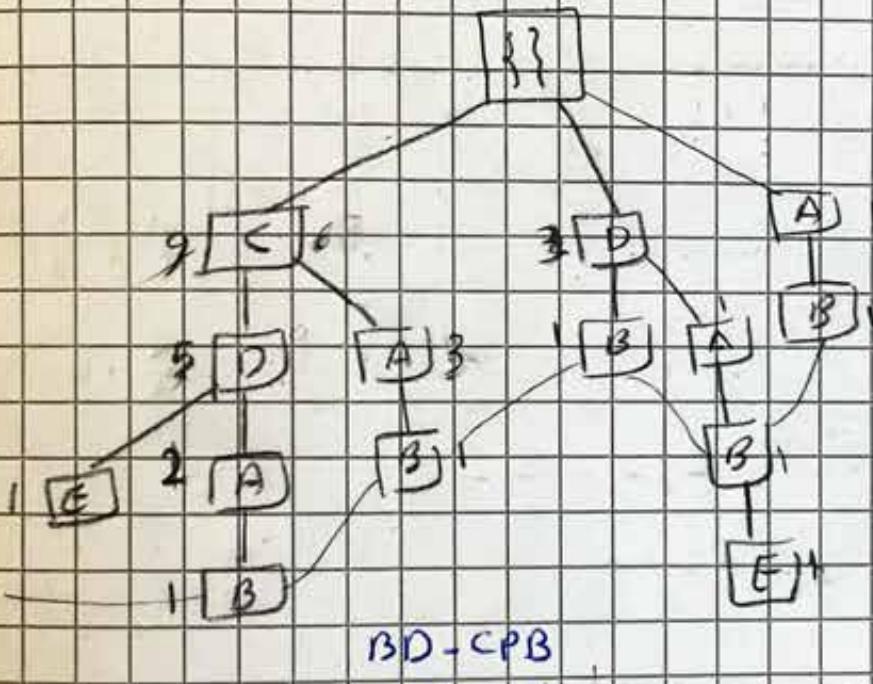
D

CD

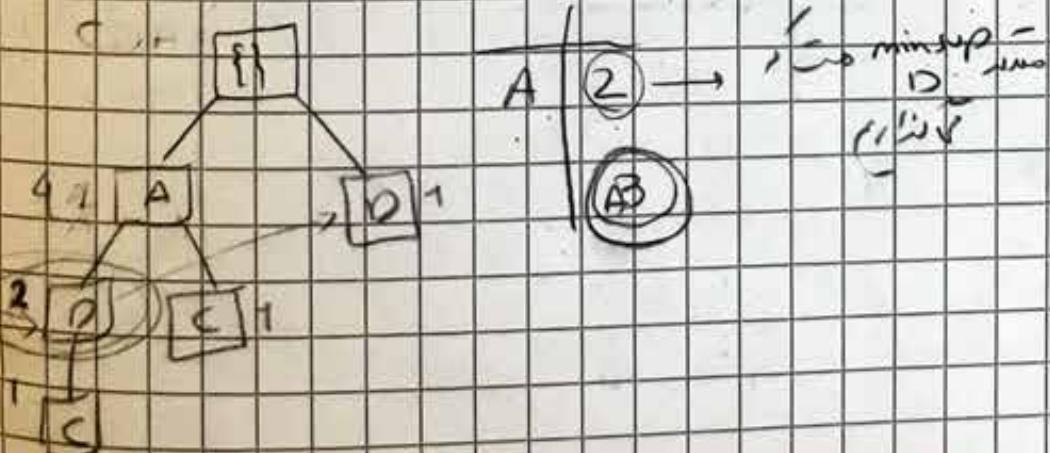
AB

ABCDE

⇒



BD - CPB



x	y	z	ground truth
2	1	3	Not Accept
0	2	-1	Accept
-3	1	1	Accept L
1	3	-2	reject
0	0	0	Accept H

$$R^2 = (x-3)^2 + (y-4)^2 + (z-2)^2$$

$1+9+4=14 \neq 15 \times$   
 $9+9+9=27 \checkmark$   
 $36+9+0=45 \checkmark$   
 $9+1+9=19 \times$   
 $9+16+15=26 \checkmark$

$\therefore r^2 > 15 \rightarrow$   
 Accept  
 $\therefore r^2 < 15 \rightarrow$   
 Reject

$$P(\text{Accept}) = \frac{3}{3}$$

$$R(\text{Accept}) = \frac{3}{4}$$

$$P(\text{reject}) = \frac{1}{2}$$

$$R(\text{reject}) = \frac{1}{2} \leftarrow \text{reject}$$

$$\frac{1}{4} \leftarrow \text{reject}$$

$$RSS = \frac{1}{n} \sum_i (y'_i - y_i)^2$$

$x \mid y$   
 $2 \mid 3$   
 $3 \mid -4$   
 $2 \mid -3$   
 $1 \mid -2$

$$RSS = \frac{1}{4} [(2m+3)^2 + (3m+4)^2 + (2m+3)^2 + (m+2)^2]$$

$$= 4m^2 + 12m + 9 + 9m^2 + 24m + 16 + 4m^2 + 12m + 9 + m^2 + 4m + 4$$

$$= 18m^2 + 52m + 38 = 0$$

$$\cancel{36m+52} \rightarrow 36m+52 = 0 \quad m = \frac{-52}{36} = -\frac{13}{9}$$

$X_1$	$X_2$
-6.7	-11.4
-3.3	-9.8
9.7	4.1
-4.9	-11.4
-4.1	-7.6
11.2	3.8
-4.2	-9.8
-6.6	-9.8
8.5	3.5
-4.7	-9.1

KNN 3

$$3 \times 10 = 30$$

$$20 \times 10 \times 3 = 600$$

$$5 \times 20 \times 10 \times 3 = 3000$$

(per iteration)  $\downarrow \downarrow \downarrow$

(8)

high width weight

4	4	2
1	2	4
4	3	1
3	4	4
3	3	2
3	2	1
3	4	3
2	2	1

type: DataFrame

level=0

groupby

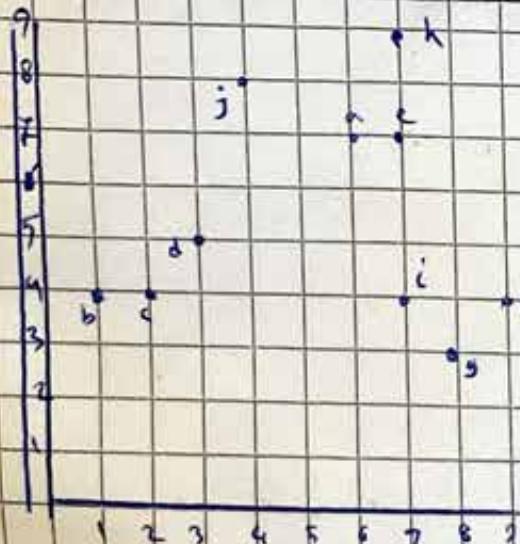
high	count
4	1
3	1
2	1
1	1
3	1
2	1
3	1
4	1

height

2 3 4 ← width

unstack →

index



a	B	0
e	C	0
h	B	0
i	C	1
g	C	1
f	C	1
b	B	2
c	C	2
d	B	2
j	N	-1



	x	y	z	L	
a	5	9	2	B	o
b	4	3	2	C	o
c	1	2	4	C	1
d	2	2	4	B	1
e	3	2	2	C	0
f	1	1	5	C	1
g	2	1	2	B	0
h	3	1	5	B	1
i	5	1	4	N	-1

8 ٢٠٢١٢٠٢١

① O/P ✓

- (a,b) (1,1,-1) = 2 ✓
- (a,c) (0,3,2) = 5 ✗
- (a,e) (2,2,0) = 4 ✗
- (b,e) (1,1,0) = 2 ✓
- (b,g) (2,2,0) = 4 ✗
- (e,g) (1,1,0) = 2 ✓
- (e,d) (1,0,-2) = 3 ✗
- (e,h) (0,1,3) = 4 ✗
- (d,g) (-1,1,2) = 3 ✗
- (d,c) (1,-1,-1) = 1 ✓
- (d,f) (1,1,1) = 3 ✗
- (d,h) (1,1,1) = 3 ✗
- (c,f) (0,1,1) = 2 ✓
- (f,h) (2,0,0) = 2 ✓
- (g,f) (1,0,3) = 4 ✗
- (g,h) (1,0,3) = 4 ✗

Ytrue : { A B C B A B C C }

② O/P ✓

Ypred : { A B C A A B B A }

accuracy :  $\frac{4}{8}$

precision(A) =  $\frac{1}{4}$

recall(A) =  $\frac{1}{1}$

(1,4) (9,2) (5,0)

C<sub>1</sub> C<sub>2</sub> C<sub>3</sub>

6									
5									
4	*	a		a(3,5)	3	4	7	C <sub>1</sub>	
3		b		b(5,4)	4	3	4	C <sub>2</sub>	
2	c	d	e	f	4(2,3)	2	3	6	C <sub>1</sub> {C <sub>1</sub> , C <sub>2</sub> , C <sub>3</sub> }
1					d(3,3)	3	2	5	C <sub>2</sub>
0					e(4,3)	9	1	9	C <sub>2</sub>
1					f(6,3)	6	3	9	C <sub>2</sub> {A, C, D, E}
2					g(5,2)	6	0	2	C <sub>2</sub> {B, D, E, F, G}
3					h(4,1)	4	5	6	C <sub>1</sub>
4					i(2,0)	5	4	3	C <sub>3</sub> {I}

C<sub>1</sub>(x,y) =  $\frac{1}{3} \left( \begin{pmatrix} (3,5) \\ (2,3) \\ (1,1) \end{pmatrix} \right) = \frac{1}{3} (5,9) = \left( \frac{5}{3}, \frac{9}{3} \right)$

C<sub>2</sub>(x,y) =  $\frac{1}{5} \left( \begin{pmatrix} (5,4) \\ (3,3) \\ (6,3) \\ (4,2) \\ (9,2) \end{pmatrix} \right) = \left( \frac{23}{5}, \frac{15}{5} \right)$

# Minkowski Distance:

~~r=1~~ Manhattan distance

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

r=2 Euclidean distance

$$d(i,j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

$r \rightarrow \infty$  Supremum distance

$$d(i,j) = \lim_{r \rightarrow \infty} \left( \sum_{k=1}^n |x_{ik} - x_{jk}|^r \right)^{1/r} = \max_k |x_{ik} - x_{jk}|$$

$$\hookrightarrow \max \{ |x_{i1} - x_{j1}|, |x_{i2} - x_{j2}| \}$$

↳ manhattan

Point	x	y	r=1	L1	P1	P2	P3	P4
P1	0	2		P1	0	4	4	6
P2	2	0		P2	4	0	2	4
P3	3	1		P3	4	2	0	2
P4	5	1		P4	6	4	2	0

$$d(P_1, P_2) = |x_2 - x_1| + |y_2 - y_1| = |2-0| + |0-2| = 4$$

$$d(P_1, P_3) = |x_3 - x_1| + |y_3 - y_1| = 1$$

$$d(P_1, P_4) = |x_4 - x_1| + |y_4 - y_1| = 1$$

$$d(P_2, P_3) = |x_3 - x_2| + |y_3 - y_2| = 1$$

$$d(P_2, P_4) = |x_4 - x_2| + |y_4 - y_2| = 1$$

$$d(P_3, P_4) = |x_4 - x_3| + |y_4 - y_3| = 1$$

$L_2$	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$	0			
$P_2$		0		
$P_3$			0	
$P_4$				0

$$d(P_1, P_2) = \sqrt{|x_2 - x_1|^2 + |y_2 - y_1|^2} = \sqrt{1^2 + 1^2} = \sqrt{2}$$

Euclidean

$L_\infty$	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$	0	2	3	
$P_2$	2	*		
$P_3$	3		*	
$P_4$				*

$$d(P_1, P_2) = \max \{ |x_2 - x_1|, |y_2 - y_1| \} = \{ 2, 2 \} = 2$$

$$d(P_1, P_3) = \max \{ |x_3 - x_1|, |y_3 - y_1| \} = \{ 3, 1 \} = 3$$

## ① SMC : (versus Jaccard)

$$\begin{matrix} M_{01} & p=0, q=1 \\ M_{10} & p=1, q=0 \end{matrix}$$

for binary

$$\begin{matrix} M_{00} & p=0, q=0 \\ M_{11} & p=1, q=1 \end{matrix}$$

Similarity between  
binary vectors

$$\begin{matrix} p = 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ q = 0 & 1 & 1 & 0 & 1 & 0 & 1 \end{matrix}$$

## ② Cosine Similarity

$$\cos(d_1, d_2) = (d_1 \cdot d_2) / \|d_1\| \|d_2\|$$

$$\frac{\downarrow}{\sqrt{(d_1 \cdot d_1)} \times \sqrt{d_2 \cdot d_2}}$$

$$d_1 = 3 \ 2 \ 0 \ 2 \ 1 \ 1$$

$$d_2 = 1 \ 1 \ 5 \ 4 \ 0 \ 1$$

$$\begin{aligned} \cos(d_1, d_2) &= \frac{3+2+0+8+0+1}{\sqrt{3^2 + 2^2 + 0^2 + 2^2 + 1^2 + 1^2} \times \sqrt{1^2 + 1^2 + 5^2 + 4^2 + 0^2 + 1^2}} \\ &= \frac{14}{\sqrt{19} \times \sqrt{44}} \end{aligned}$$

## Classification techniques:

- 1- Decision trees
- 2- Classification rules
- 3- Association Rules
- 4- Neural Networks
- 5- Naive Bayes and Bayesian Networks
- 6- K-Nearest Neighbours (K-NN) → Euclidean distance  
 $d(x_1, x_2) = \sqrt{\sum (x_i - x_j)^2}$
- 7- Support Vector Machines (SVM)

## Clustering Algorithms:

- 1- K-means and its variants → SSE
- 2- Hierarchical clustering
- 3- Density-based clustering

### Evaluating K-means:

SSE (Sum of squared Error)

$$SSE = \sum_{i=1}^k \sum_{x \in C_i} \text{dist}^2(m_i, x)$$

→ absolute error ( $\sum_{i=1}^k \sum_{x \in C_i} \text{dist}(m_i, x)$ )

→ error ( $\sum_{i=1}^k \sum_{x \in C_i} \text{dist}^2(m_i, x)$ )

One way to reduce SSE is to increase K, the number of clusters  
(A good clustering with smaller K can have a lower SSE than a poor clustering with higher K)

Running K-means with K=10 on the same dataset will typically produce lower mean SSE

	$\frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}}$
TN	number of pairs of objects having a different class and a different cluster
FP	different class
FN	same class
TP	different cluster
	same cluster

numbers = [1, -8, 5, -2, 5]

Filter: Filter the element of a list based on a condition

مثلاً اعطيك لست بـ 5 عناصر، وتحتاج إلى إيجاد العناصر التي تحقق شرط معين

neg = list(filter(lambda x: x < 0, numbers))

neg → [-8, -2]

map: map each element of a list with a new value

sqr = list(map(lambda x: x \* x, numbers))

sqr = [64, 4]

map(function, list)

ex: R = list(map(lambda x: 'big' if x > 10 else 'small', numbers))

filter(function, list)

\* فیلتر عبار خالقی، دست دهنده بروزگار True یا False است، باعدها در برخی از دستورات رسمیت،

## partitional clustering

## k-means

$$c_x = \frac{\sum x_i}{n}$$

$$c_y = \frac{\sum y_i}{n}$$

$c_1, c_2, c_3$

الهدف من k-means هو تعيين مركبات كل نقطة في مجموعة البيانات إلى إحدى k مجموعات

$$\sum_{i=1}^N \min ||x^i - c_j||^2$$

مترافق مع ترتيب المركبات

Rand Index بالمقدار

$$\bar{X} = \frac{\sum x_i}{n} = m$$

x	y
x <sub>1</sub>	x <sub>2</sub>
a	d
b	e
c	f

$$Var(x) = \frac{1}{n-1} [(a-m)^2 + (b-m)^2 + (c-m)^2]$$

$$Cov(x, y) = \frac{1}{n-1} [(a-m)(d-m) + (b-m)(e-m) + (c-m)(f-m)]$$

$$Cov(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

مترافق مع [x, y, z]

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

axis

np.sort(a1, axis = 1)

np.sort(a1, axis = 0)

np.sort(a1, axis = None)

a.sum(axis = 0) mean(axis=0)  
مقدار مجموع عناصر في كل سطر

a.min(axis = 1)

np.sort(a1; axis = -1)

## Accessing Series Elements

1) by index

s1=pd.Series({'a1':1,'a2':2,'a3':3})

s1=pd.Series([1,2,3],index=['a1','a2','a3'])

- Explicit → print(s1.loc['a1'])

print

- Implicit → print(s1.iloc[1])

print(s1.loc['a1']) = 1

print(s1.iloc[1]) = 2

print(s1.loc['a1':'a2'])

explicit

Stop element included

print(s1.iloc[1:3])

Implicit

Stop element excluded

print(s1.loc[a1,a2])

a1 1  
a2 2

Q: A decision tree trained on a dataset with categorical features only:

A: Each internal node of a decision tree can have at most C children, where C is equal to the number of distinct values of the attribute associated with the internal node.

$$RI(y_{true} \rightarrow y_{pred}) = \frac{a+b}{\binom{n}{2}} \quad [ ]_{n \times m}$$

unlabeled dataset X (n rows, m columns) and run a cluster analysis using DBSCAN

\* The number of clusters is not a hyper parameter that can be fine-tuned (based on the knee of the curve of the sum of squared errors (SSE))

\* The number of clusters extracted is dependent of the hyperparameter  $\epsilon$  (DBSCAN)

\* DBSCAN may extract at most n cluster

\* You can use the RI to assess the quality of the resulting classification (not clusters)

Feature  $\rightarrow x_0, x_1, x_2, \dots$

$\begin{bmatrix} 0 & 0 \\ ? & ? \end{bmatrix}$   $n \times m$

$n \rightarrow$  sample  $n=?$

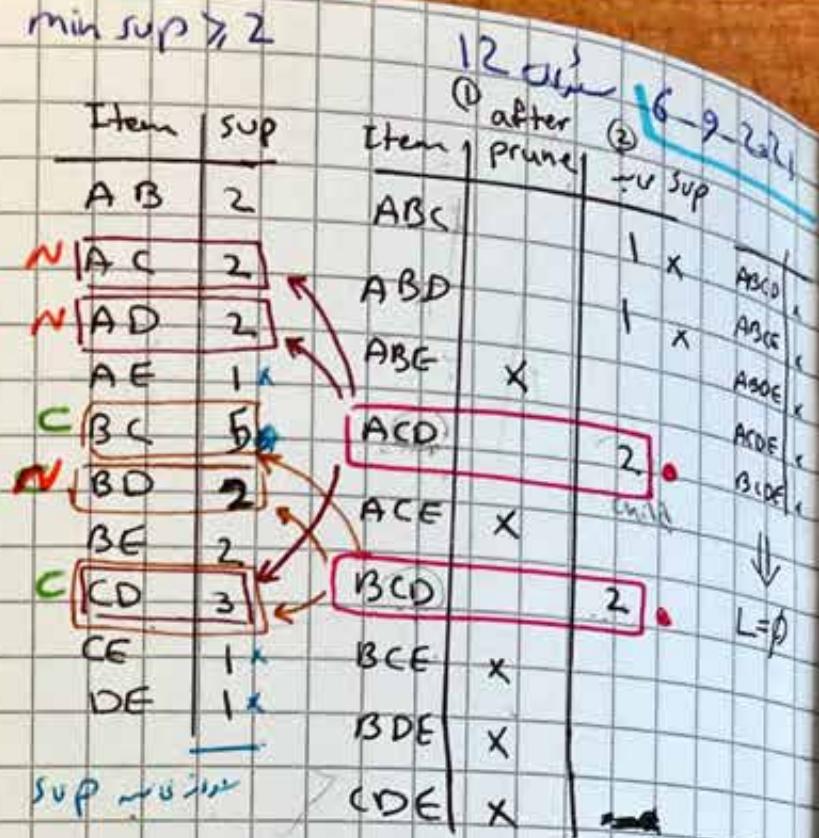
$m \rightarrow$  Feature

nb/  
sample

## Apriori Algorithm

$$\min \text{sup} \geq 2$$

	transaction	Item	sup
0	BC		
1	AE	C	4
2	ACD	C	7
3	AB	C	6
4	BF	C	4
5	DE	C	4
6	BCE	E	4
7	BC		
8	BCD		
9	ABCDE		



Q7 = list all frequent Itemsets having length 2 along with the Support Count

$$AT = \{AB:2, AC:2, AD:2, BD:2, BE:2, CD:3\}$$

Q2 = list all itemsets of length 3 that have been generated by Apriori after the join and prune steps. (before counting their support in the database)

$$AL = \{ ABC, ABD, ACD, BCD \}$$

نک: هم قائم می باشد از سه عبارت  $A \cup B$ ,  $A \cup C$ ,  $B \cup C$  برای  $(A \cup B) \cup C = A \cup (B \cup C)$  برای  $(A \cup B) \cup C = A \cup (B \cup C)$

Q3 = list all frequent itemsets that are not closed. along with their support counts.

$$A3 = \{ AC:2, AD:2, BD:2 \}$$

برای مجموعه  $A = \{C, D, A, C\}$  درست بین عبارت  $\text{sup}(A)$  و  $\text{sup}(\{C, D\})$  را بحث کنید.

أُنْتَهِيَّةُ تَمَكُّنِي مَعَ سَعْيِي

[ سَعْيٌ ، تَمَكُّنٌ ]