

EXERCISES - PART I

KNN

1. 3-NN $(1,0,0,5) \rightarrow P$

NN:	$(0,75; 0,25)$	$\rightarrow d = \sqrt{0,25^2 + 0,25^2} = 0,3536$	N
	$(0,75; 0,75)$	$\rightarrow d = \sqrt{0,25^2 + 0,25^2} = 0,3536$	N
	$(1,25; 0,25)$	$\rightarrow d = \sqrt{0,25^2 + 0,25^2} = 0,3536$	N
	$(1,25; 0,75)$	$\rightarrow d = \sqrt{0,25^2 + 0,25^2} = 0,3536$	P

(C) is not able to classify

2. 3-NN $(1,75; 2) \rightarrow$

$\rightarrow (2,2)$	$\rightarrow d = \sqrt{0,25^2} = 0,25$	P
$\rightarrow (1,5; 2)$	$\rightarrow d = \sqrt{0,25^2} = 0,25$	N
$(1,2)$	$d = \sqrt{0,75^2 + 0} = 0,75$	N
$(1,75; 0)$	$d = \sqrt{0 + 2^2} = 2$	
$(1,25; 3)$	$d = \sqrt{0,5^2 + 1^2}$	

4.

x_1	x_2	x_3	Class
0	0	1	0
0	1	1	1
1	0	1	1
1	1	1	0
0	0	0	0
0	1	0	0
1	0	0	0
1	1	0	0

a) $(1,1,1)$

$$d = \sqrt{1^2 + 0 + 0} = 1$$

$$d = \sqrt{0 + 1^2 + 0} = 1$$

3NN

(B)

$$d = \sqrt{0 + 0 + 1^2} = 1$$

b) leave-one-out

x_1	x_2	x_3	Class	
0	0	0	0	TN
0	0	1	0	TN
0	1	0	0	TN
0	1	1	0	FN
1	0	0	0	TN
1	0	1	0	FN
1	1	0	1	FP
1	1	1	1	FP

$$\text{Accuracy} = \frac{TP + TN}{All} = \frac{4}{8} = \frac{1}{2}$$

NAIVE BAYES

1.

a) instance (A, B)

$$P(C=P | (A, B)) = P(C=P) \times P(A|P) \times P(B|P) = \\ = \frac{3}{8} \times 1 \times \frac{1}{2} = \frac{1}{8}$$

(B) NEGATIVE

$$P(C=N | (A, B)) = P(C=N) \times P(A|N) \times P(B|N) = \\ = \frac{5}{8} \times \frac{2}{5} \times \frac{4}{5} = \frac{1}{2} >$$

2. instance (T, T, F) $P(P) = 2/8 = 1/4$ $P(N) = 3/4$

$$P(C=P | (T, T, F)) = P(C=P) \times P(A=T|P) \times P(B=T|P) \times P(C=F|P) = \\ = \frac{1}{4} \times \frac{2}{5} \times 0 = 0$$

(B) NEGATIVE

$$P(C=N | (T, T, F)) > 0$$

4. x_1 numerical e x_2 ordinal instance (0.75; 1)

$$P(C=P | (0.75; 1)) = P(C=P) \times P(x_1=0.75 | P) \times P(x_2=1 | P)$$

$$= \frac{7}{12} \times P(x_1=0.75) \times \frac{2}{7}$$

$$P(x_1=0.75) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x-\bar{x})^2}{2\sigma^2}}$$

DECISION TREES

1. A, B, C numerical

a) ENTROPY do atributo class

$$E(P, N) = -P(P) \cdot \log_2 P(P) - P(N) \cdot \log_2 P(N) = \\ = -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) \approx 0,97$$

b) ATTRIBUTE TESTED IN THE ROOT \rightarrow CH.5

(INFORMATION GAIN CRITERIA)

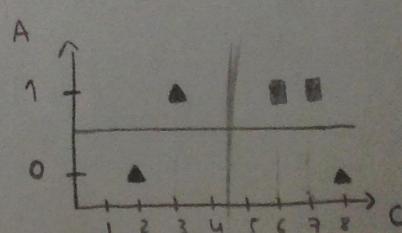
$$G(A) = E(P, N) - y_0 \times E(0+, 0-) - y_1 \times E(1+, 1-) \\ = 0,97 - \frac{2}{5} \times E(2, 0) - \frac{3}{5} \times E(1, 2) = 0,42$$

$$E(2, 0) = -\frac{2}{2} \log_2 \frac{2}{2} - 0 = -1 \times 0 = 0$$

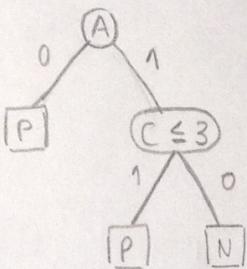
$$E(1, 2) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0,92$$

\downarrow Fazem os contas p/ os outros

$$G(A) = G(C)$$



c)



remover lâminas q A = Ø

$$E(\text{classe}) = E(1,2) \approx 0,92$$

$$G(B) = \dots \rightarrow 0 \text{ que dir } \rightarrow \text{val para o ramo}$$

$$G(C) = \dots$$

EX. 8

$$a) x_2 \leq 1,25 \quad x_1 \leq 2,75$$

$$\text{gini split } (x_1 \leq 2,75) = p(x_1 \leq 2,75) \times \text{gini } (x_1 \leq 2,75) + \\ p(x_1 > 2,75) \times \text{gini } (x_1 > 2,75) = \\ = \frac{8}{20} \times 0 + \frac{12}{20} \times 0,49 = 0,29$$

$$\text{Gini } (x_1 \leq 2,75) = 1 - \left(\frac{8^2}{82} \right) = 0 \rightarrow \text{quantos não P e não } > 2,75$$

$$\text{Gini } (x_1 > 2,75) = 1 - \left(\frac{(5^2 + 7^2)}{12^2} \right) = 0,49 \rightarrow \text{quantos não N e não } > 2,75$$

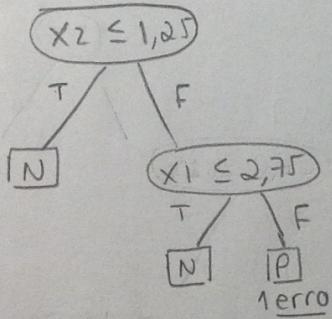
$$\text{gini split } (x_2 \leq 1,25) = p(x_2 \leq 1,25) \times \text{gini } (x_2 \leq 1,25) + \\ p(x_2 > 1,25) \times \text{gini } (x_2 > 1,25) = \\ = \frac{10}{20} \times 0 + \frac{10}{20} \times 0,5 = 0,25$$

$$\text{Gini } (x_2 \leq 1,25) = 1 - \left(\frac{10^2}{10^2} \right) = 0$$

$$\text{Gini } (x_2 > 1,25) = 1 - \left(\frac{5^2 + 5^2}{10^2} \right) = 0,5$$

\Rightarrow escolher gini < : 0,25 $\rightarrow x_2$

b)



gini split $\begin{cases} x_2 \leq 1,25 \\ x_1 \leq 2,75 \end{cases}$

$$c) \text{Accuracy} = \frac{TP+TN}{All} = \frac{19}{20}$$

$$\text{sensitivity} = \frac{TP}{TP+FN} = \frac{5}{5+0} = 1$$

$$\text{specificity} = \frac{TN}{TN+FP} = \frac{14}{15}$$

EX.12 A and B Boolean, C numerical — 0,25; 0,5; 0,75 (values)

a) (A)

b) (A)

c) ENTROPY

$$E(P,N) = -\frac{8}{12} \log_2 \frac{8}{12} - \frac{4}{12} \log_2 \frac{4}{12}$$

$$= 0,918 \quad (\text{C})$$

$$\text{d) ACCURACY} = \frac{TP+TN}{\text{All}} = \frac{6+4}{12} = \frac{10}{12} \approx 0,83$$

$$83\% \quad (\text{B})$$

$$\text{e) SENSIBILITY} = \frac{TP}{TP+FN} = \frac{6}{8} = 0,75$$

$$\text{SPECIFICITY} = \frac{TN}{TN+FP} = \frac{4}{4} = 1$$

(B)

A	B	C	class
0	0	0,25	N
0	1	0,25	P
1	0	0,25	P
1	1	0,25	N
0	0	0,5	N
0	1	0,5	P
1	0	0,5	P
1	1	0,5	N
0	0	0,75	N
0	1	0,75	P
1	0	0,75	P
1	1	0,75	N

EX.13. A, B, C numeric

a) (D)

b) (B)

$$\text{c) ENTROPY} = -\frac{9}{27} \log_2 \frac{9}{27} - \frac{18}{27} \log_2 \frac{18}{27} = 0,764 \quad (\text{B})$$

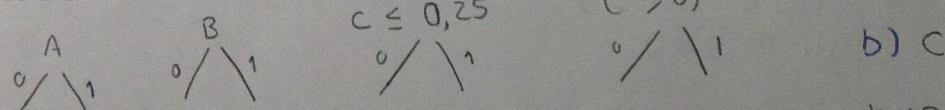
$$\text{d) ACCURACY} = \frac{TP+TN}{\text{All}} = \frac{9+15}{27} = 0,89 \quad (\text{B})$$

$$\text{e) SENSIBILITY} = \frac{TP}{TP+FN} = \frac{6}{9} = 0,67 \quad (\text{C})$$

$$\text{SPECIFICITY} = \frac{TN}{TN+FP} = \frac{18}{18} = 1$$

ENSEMBLE

EX.1. weak classifiers : all decision tree with only one binary node



$$c \leq 0,25$$

$$(c \geq 0,75)$$

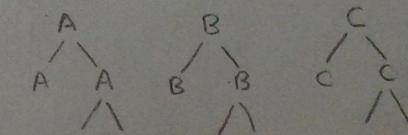
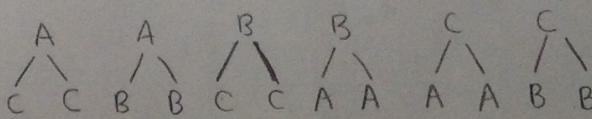
a) B

b) C

EX.2. weak classifiers : regular trees with just two levels of nodes

A, B e C 3 values

arvore 2 níveis



ASSOCIATION RULES

Ex.2. ABCDF, ABEG, AD, BCDE, CD min support = 40%

C1:	L1:	C2:	L2:	C3:
A 3/5	A	AB 2/5	AB	ABD
B 3/5	B	AC 1/5	AD	BCD (A)
C 3/5	C	AD 2/5	BC	BCE
D 4/5	D	AE 1/5	BD	BDE
E 2/5	E	BC 2/5	BE	
F 1/5		BD 2/5	CD	
		BE 2/5		
		CD 3/5		
		CE 1/5		
		DE 1/5		

b) CONFIDENCE ($A \rightarrow B$) = $P(B|A)$ SUPPORT($A \rightarrow B$) = $P(A \cap B)$

LIFT ($A \rightarrow B$) = $\frac{P(B|A)}{P(B)}$

$$\text{CONF}(A \rightarrow D) = P(D|A) = \frac{P(D \cap A)}{P(A)} = \frac{\frac{2}{5}}{\frac{3}{5}} = \frac{2}{3}$$

$$\text{CONF}(D \rightarrow A) = P(A|D) = \frac{P(A \cap D)}{P(D)} = \frac{\frac{2}{5}}{\frac{4}{5}} = \frac{2}{4} = \frac{1}{2}$$

$$\text{LIFT}(A \rightarrow D) = \frac{P(D|A)}{P(D)} = \frac{P(D \cap A)}{P(D) \times P(A)} = \frac{\frac{2}{5}}{\frac{4}{5} \times \frac{3}{5}} = \frac{2}{\frac{12}{25}} = \frac{50}{60} = \frac{5}{6}$$

$$\text{LIFT}(D \rightarrow A) = \frac{\text{CONF}(D \rightarrow A)}{P(A)} = \frac{\frac{1}{2}}{\frac{2}{5}} = \frac{5}{6} \quad (\text{A})$$

Ex.1.

C1:	L1:	min support = 40%			C4: \emptyset
A 2/4	A	AB 1/4	AC	BCG 2/4	BCG
B 3/4	B	AC 2/4	BC		
C 3/4	C	AC 1/4	BE		
D 1/4	E	BC 2/4	CE		
E 3/4		BE 3/4			
		CE 2/4			

C: candidates
L: frequent item sets

$$\text{CONF}(X \rightarrow Y) = P(Y|X) = \frac{P(Y \cap X)}{P(X)}$$

b) ASSOC. RULES WITH CONF > 50%.

X	Y	CONF
{BCE}	X	$2/4 / 3/4 = 2/3$
B \rightarrow BCE		$2/4 / 3/4 = 2/3$
C \rightarrow BCE		$2/4 / 3/4 = 2/3$
E \rightarrow BCE		$2/4 / 2/4 = 1$
{BCY} \rightarrow E		$2/4 / 2/4 = 1$
{CEY} \rightarrow B		$2/4 / 2/4 = 1$
{BEY} \rightarrow C		$2/4 / 3/4 = 2/3$

A \rightarrow B
A \rightarrow D
B \rightarrow C
B \rightarrow D
B \rightarrow E
C \rightarrow D
B \rightarrow A
D \rightarrow A
C \rightarrow B
D \rightarrow B
E \rightarrow B
D \rightarrow C

EX.3. 12346, 1257, 14, 2345, 34 min supp: 40%

C1:	L1:	C2
1 315	1	12 215
2 315	2	13 115
3 315	3	14 215
4 415	4	15 115
5 215	5	23 215
6 115		24 215
7 115		25 215
		34 315
		35 115
		45 115

235, 245

23, 25, 38

24, 25, 45

(C)

EX.7. {A, B, C, D, E} ABCD ACDE

a) L4
 $\begin{array}{l} \text{ABCD} - \text{ACOB} \\ \text{ACDE} - \text{ACDE} \end{array} \rightarrow \begin{array}{l} \text{ACDBE} \\ \downarrow \\ \text{ACDB} \\ \text{ACDE} \\ \text{BCDE} \end{array}$ (4)
 $\text{ACDE} \neq \text{L4}$ porque não scan
 4 vezes
 (L1, L2, L3, L4)

b) L4 = {ABCD, ACDE}
 $L_3 = \{\underline{ABC}, \underline{ABD}, \underline{BCD}, \underline{ACD}, \underline{ADE}, \underline{CDE}, \underline{ACE}\}$
 $L_2 = \{AB, AC, AD, BC, BD, AE, CD, CE\}$
 $L_1 = \{A, B, C, D, E\}$
 $C_3 = \{\underline{ABC}, \underline{ABD}, ABE, \underline{ACD}, \underline{ACE}, \underline{ADE}, \underline{BCD}, \underline{CDE}\}, \dots\}$
 $C_3 \neq L_3 \rightarrow$ only for all the subsets (with 3 items)
 of each pattern.

EX.11. L4 = {ABCD, ABCE, ABDE, ACDE}

EX.12 {ABCD, ABCE}
 \downarrow
ABCDE

EX.13

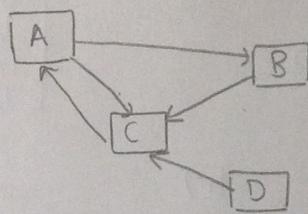
SOCIAL NETWORK ANALYSIS

EX.1.

a) cintânia DB \rightarrow 3

b)

	IN-DEGREE	OUT-DEGREE	DEGREE
A	1	2	3
B	1	1	2
C	3	1	4
D	0	1	1



c) DEGREE CENTRALITY

$$DC(A) = \frac{2}{3} \quad DC(C) = \frac{1}{3}$$

$$DC(B) = \frac{1}{3} \quad DC(D) = \frac{1}{3}$$

CLOSENESS CENTRALITY

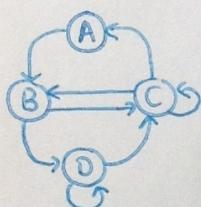
$$CC(A) = \frac{3}{1+1+\infty} = 0$$

$$CC(B) = 0$$

$$CC(C) = 0$$

$$CC(D) = \frac{3}{1+2+3} = \frac{3}{6} = \frac{1}{2}$$

EX.2.



a) INITIALIZATION FOR PAGERANK SCORES

$$\begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \end{bmatrix}$$

b) TRANSITION PROBABILITY MATRIX

$$\begin{array}{l} \begin{array}{cccc} & A & B & C & D \end{array} \\ \begin{array}{c} A \\ B \\ C \\ D \end{array} \left[\begin{array}{cccc} 0 & 1 & 0 & 0 \\ 0 & 0 & 1/2 & 1/2 \\ 1/3 & 1/3 & 1/3 & 0 \\ 0 & 0 & 1/2 & 1/2 \end{array} \right] \end{array}$$

c) PAGERANKS

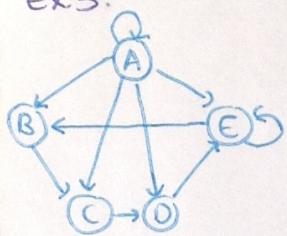
$$PR(A) = 1 - 0,8 + 0,8 \times \left(\frac{1/4}{3} \right) =$$

$$PR(B) = 1 - 0,8 + 0,8 \times \left(\frac{1/4}{1} + \frac{1/4}{2} \right)$$

$$PR(C) = 1 - 0,8 + 0,8 \times \left(\frac{1/4}{3} + \frac{1/4}{2} + \frac{1/4}{2} \right)$$

$$PR(D) = 1 - 0,8 + 0,8 \times \left(\frac{1/4}{2} + \frac{1/4}{2} \right)$$

Ex.3.



a)

$$\begin{bmatrix} 1/5 \\ 1/5 \\ 1/5 \\ 1/5 \\ 1/5 \end{bmatrix}$$

b)

$$\begin{array}{ccccc} & A & B & C & D & E \\ A & \left[\begin{array}{ccccc} 1/5 & 1/5 & 1/5 & 1/5 & 1/5 \end{array} \right] \\ B & \left[\begin{array}{ccccc} 0 & 0 & 1 & 0 & 0 \end{array} \right] \\ C & \left[\begin{array}{ccccc} 0 & 0 & 0 & 1 & 0 \end{array} \right] \\ D & \left[\begin{array}{ccccc} 0 & 0 & 0 & 0 & 1 \end{array} \right] \\ E & \left[\begin{array}{ccccc} 0 & 1/2 & 0 & 0 & 1/2 \end{array} \right] \end{array}$$

c)

$$PR(A) = 1 - 0,7 + 0,7 \times \left(\frac{1/5}{5} \right)$$

$$PR(B) = 1 - 0,7 + 0,7 \times \left(\frac{1/5}{5} + \frac{1/5}{2} \right)$$

$$PR(C) = 1 - 0,7 + 0,7 \times \left(\frac{1/5}{1} + \frac{1/5}{5} \right)$$

$$PR(D) = 1 - 0,7 + 0,7 \times \left(\frac{1/5}{1} + \frac{1/5}{5} \right)$$

$$PR(E) = 1 - 0,7 + 0,7 \times \left(\frac{1/5}{1} + \frac{1/5}{2} + \frac{1/5}{5} \right)$$

Exercises - Part II

Feature selection

1.1 Spearman correlation $\rightarrow 1 - \frac{6 \cdot \sum D^2}{n(n^2-1)} = 1 - \frac{6 \cdot (4+4+0,25+0,25)}{4(16-1)}$
 (between y_1 and y_4)

y_1	y_4	Rank y_1	Rank y_4	D	D^2	
-0,4	2	1	3	-2	4	$= 1 - \frac{6 \cdot 8,5}{60} =$
0,7	1	3,5	1,5	2	4	$= 1 - 0,85 = 0,15$
0,7	3	3,5	4	-0,5	0,25	
-0,3	1	2	1,5	0,5	0,25	

$$\frac{3+4}{2} = 3,5$$

$$\frac{1+2}{2} = 1,5$$

$$\bar{y}_1 = 0,175$$

$$\bar{y}_3 = -0,325$$

pearson correlation $= \frac{\text{cov}(y_1, y_3)}{\sqrt{\text{var}(y_1) \cdot \text{var}(y_3)}} = \frac{0,2317}{\sqrt{0,608 \cdot 0,478}} = 0,43$
 (between y_1 and y_3)

$$\text{cov}(y_1, y_3) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} =$$

$$= (-0,4 - 0,175) \cdot (0,1 + 0,325) + (0,7 - 0,175) \cdot (-0,3 + 0,325) + \\ (0,7 - 0,175) \cdot (-0,1 + 0,325) + (-0,3 - 0,175) \cdot (-1 + 0,325)$$

$$= \frac{0,244 + 0,013 + 0,178 + 0,32}{3} = 0,2317$$

$$\text{var}(y_1) = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{0,33 + 0,276 + 0,276 + 0,226}{3}} = 0,608$$

$$\text{var}(y_3) = \sqrt{\frac{0,18 + 0,000625 + 0,05 + 0,456}{3}} = 0,478$$

1.3. Discretization

EQUAL-DEPTH / Frequency

$$(y_3) \rightarrow <0,1, -0,3, -0,1, -1>$$

$$y_3: <1, 0, 1, 0>$$

$$[-1; -0,3] \quad [-0,1; 0,1]$$

$$y_2: <1, 1, 0, 1>$$

EQUAL-WIDTH / Range

$$(y_3) \rightarrow <0,1, -0,3, -0,1, -1>$$

$$<1, 1, 1, 0>$$

$$w = \frac{\max - \min}{2} = \frac{0,1 + 1}{2} = 0,55$$

$$[-1, -0,45] \quad [-0,45, 0,1]$$

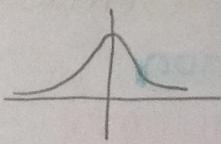
$$y_2: <1, 1, 0, 1>$$

GAUSSIAN CUT-OFF POINTS

$$\bar{y}_3 = -0,325 \quad \sigma_{y_3} = 0,4145$$

Δ já estão normalizados

$$\begin{array}{l} < 0 \rightarrow 0 \\ > 0 \rightarrow 1 \end{array} \quad \langle 1, 0, 0, 0 \rangle$$



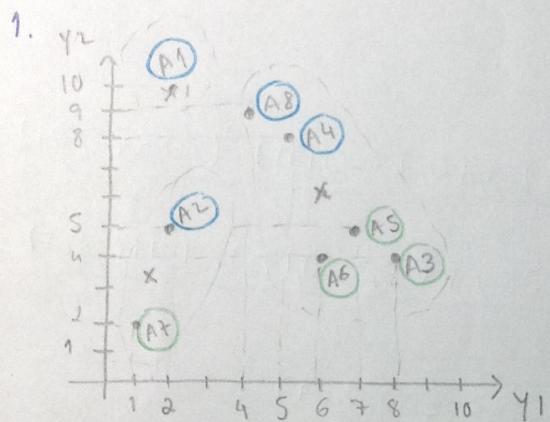
Reduce overfitting in decision trees

Pruning

enforce a minimum nr of samples in leaf nodes

enforce a maximum depth for the tree

CLUSTERING



1.1. $A_1, A_4, A_7 \rightarrow$ needs
 $\downarrow \quad \downarrow \quad \downarrow$
 K_1, K_2, K_3

- distância de cada um a A_1, A_4 e A_7

- \downarrow
 $(K_1) A_1$
 $(K_2) A_3, A_5, A_6, A_8, A_4$
 $(K_3) A_2, A_7$

K-means

$$\begin{array}{l} K_1 = A_1 \\ 1^{\text{st}} \text{ epoch} \quad K_2 = (6, 6) \\ \quad \quad \quad K_3 = (1.5, 3.5) \end{array}$$

$$\begin{array}{l} (K_1) : A_1, A_8 \\ \rightarrow (K_2) : A_4, A_5, A_6, A_3 \\ \quad \quad \quad (K_3) : A_2, A_7 \end{array} \quad \begin{array}{l} K_1 = (3, 9.5) \\ \rightarrow K_2 = (6.5, 5.25) \end{array}$$

K-median

$$K_1 = (2, 10)$$

$$K_2 = (6, 5) \checkmark$$

$$K_3 = (1, 5, 3, 5)$$

K-medoids

↓

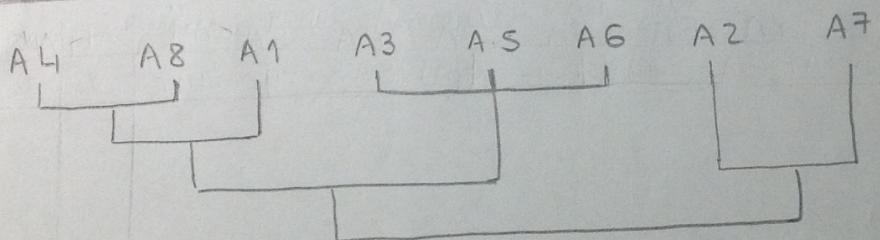
ver o ponto que está + centro

1.2.

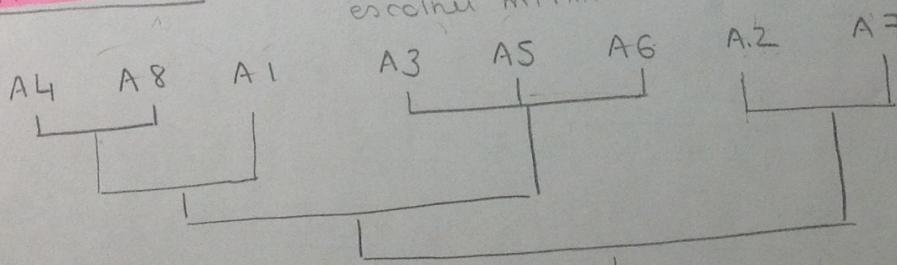
	A1	A2	A3	A4	A5	A6	A7	A8
A1	0	5	6	$\sqrt{13}$	$\sqrt{50}$	$\sqrt{52}$	$\sqrt{5}$	
A2		0	$\sqrt{37}$	$\sqrt{18}$	5	$\sqrt{7}$	$\sqrt{10}$	$\sqrt{20}$
A3			0	5	$\sqrt{2}$	$\sqrt{2}$	$\sqrt{53}$	$\sqrt{41}$
A4				0	$\sqrt{3}$	$\sqrt{4}$	$\sqrt{52}$	$\sqrt{2}$
A5					0	$\sqrt{2}$	$\sqrt{45}$	5
A6						0	$\sqrt{29}$	$\sqrt{29}$
A7							0	$\sqrt{58}$
A8								0

AGGLOMERATIVE: initially every point is a cluster of its own and we merge cluster until we end-up with one unique cluster containing all points

SINGU UNIC: ver as distâncias entre clusters e escolher a mínima



COMPUTE UNIC: escolher máx de distância entre clusters e depois escolher mínima das máximas



AVERAGE UNIC:

UNIFORM UNIC / K-mean

↳ calcular centro (ponto médio) de cada cluster, para a distância entre cada um desses pontos e seja o menor

1.3.

$$(1) \text{ PURITY}(C3) = \frac{1}{2} \cdot 1 = \frac{1}{2} \checkmark$$

$$(4) \text{ PURITY}(\text{clustering}) = \frac{1}{8} (3+3+1) = \frac{7}{8} \checkmark$$

$$(3) \text{ RI} = \frac{\text{TP} + \text{TN}}{\text{All}} = \frac{15+6}{15+6+1+6}$$

$\text{TN} = \{A1, A3\}, \{A1, A5\}, \{A1, A6\} \cup \{A4, A5\} \cup \{A4, A6\} \cup \{A7, A8\} \cup \{A8, A9\} \cup \{A8, A10\} \cup \{A9, A10\}$
 $\{A2, A3\} \cup \{A2, A5\} \cup \{A2, A6\} \cup \{A7, A10\} \cup \{A7, A9\} \rightarrow 15$

≠ classe em ≠ cluster

$$\text{TP} = \{A1, A4\} \cup \{A1, A8\} \cup \{A2, A5\} \cup \{A3, A6\} \cup \{A5, A6\} \rightarrow 6$$

= classe no mesmo cluster

$$\text{FP} = \{A2, A7\} \rightarrow 1$$

estão no mesmo cluster e não deviam

$$\text{FN} = \{A2, A1\} \cup \{A2, A4, A2, A8\} \cup \{A7, A3\} \cup \{A7, A5, A7, A6\} \rightarrow 6$$

estão separados e deviam estar juntos

2.

	y_1	y_2	cluster
x_1	2	2	C1
x_2	2	4	C1
x_3	4	3	C2
x_4	4	5	C2

(1) SILHOUETTE INDEX x_2

exclusão $a = d(x_2, u_1) = \sqrt{0+2^2} = 2$

$$d(x_2, u_3) = \sqrt{2^2+1^2} = \sqrt{5} \quad > \quad b = \frac{\sqrt{5}+2}{2} = \sqrt{5} > 2$$

$$d(x_2, u_4) = \sqrt{2^2+1^2} = \sqrt{5}$$

$$b > a \quad \text{s.i.} = \frac{b-a}{b} = \frac{\sqrt{5}-2}{\sqrt{5}} = 0,1$$

$a = d(u_2, u_1) = \max(12-21, 14-21) = 2$

$$d(u_2, u_3) = \max(14-21, 13-21) = 2 \quad > \quad b = \frac{2+2}{2} = 2$$

$$d(u_2, u_4) = \max(14-21, 15-21) = 2$$

$a = b \rightarrow \text{silhouette} = 0$ (equally distant to C1 and C2)

CHBYHEV

(3) AVERAGE SILHOUETTE OF CLUSTER C₂

SILHOUETTE X₃

$$a = d(x_3, x_4) = |4-4| + |5-3| = 0+2=2$$

$$d(x_3, x_1) = |4-2| + |5-2| = 2+1=3 \quad \rightarrow b = \frac{3+3}{2} = 3$$

$$d(x_3, x_2) = |4-2| + |3-4| = 2+1=3$$

$$b > a$$

$$\Delta.i(x_3) = \frac{b-a}{b} = \frac{3-2}{3} = \frac{1}{3}$$

SILHOUETTE X₄

$$a = d(x_4, x_1) = 2$$

$$d(x_4, x_1) = |4-2| + |5-2| = 2+3=5 \quad \rightarrow b = \frac{5+3}{2} = 4$$

$$d(x_4, x_2) = |4-2| + |5-4| = 2+1=3$$

$$b > a$$

$$s.i(x_4) = \frac{4-2}{4} = \frac{2}{4} = \frac{1}{2}$$

$$\text{average silhouette}(C_2) = \frac{\frac{1}{3} + \frac{1}{2}}{2} = 0,417 //$$

2.2. 4 Euclidean

COHESION ERROR \downarrow

$$\begin{aligned} SSE &= (2-2)^2 + (2-2)^2 + \\ &\quad (2-3)^2 + (4-3)^2 + \\ &\quad (4-4)^2 + (4-4)^2 + \\ &\quad (3-4)^2 + (5-4)^2 \\ &= 1+1+1+1=4 \end{aligned}$$

$$\left| \begin{array}{l} \bar{c}_1 = (2, 3) \\ \bar{c}_2 = (4, 4) \end{array} \right.$$

SEPARATION ERROR \downarrow

$$= 2 \cdot 1,25 + 2 \cdot 1,25 = 5$$

PONTO MÍDIO ENTRE CLUSTERS

$$\downarrow$$

$$\frac{\bar{c}_1 + \bar{c}_2}{2} = (3, 3,5)$$

$$d(c_1, p_n) = \sqrt{1^2 + (0,5)^2} = \sqrt{1,25}$$

$$d(c_2, p_n) = \sqrt{1^2 + (0,5)^2} = \sqrt{1,25}$$

TOTAL ERROR = SSE + BSS
 $= 4+5=9$

-D

PCA

3. $20 + 10 + 5 + 4 + 3 + 2 + 1 = 45$

$$\frac{45}{45+x} = 0,9 \Leftrightarrow \frac{45}{0,9} = 45 + x \Leftrightarrow x = 5$$

$$7 + 5 = \underline{\underline{12}} \text{ attributes}$$

4.

$$S = \begin{bmatrix} 91.43 & 171.92 & 297.99 \\ 171.92 & 373.92 & 545.21 \\ 297.99 & 545.21 & 1297.26 \end{bmatrix}$$

$$S \vec{u}_1 = \lambda_1 \vec{u}_1 \Leftrightarrow (S - \lambda_1 I) \cdot \vec{u}_1 = 0$$

$$\begin{bmatrix} 91.43 & 171.92 & 297.99 \\ 171.92 & 373.92 & 545.21 \\ 297.99 & 545.21 & 1297.26 \end{bmatrix} \begin{bmatrix} 0.22 \\ 0.41 \\ 0.88 \end{bmatrix} = \begin{bmatrix} 352,833 \\ 670,91 \\ 1430,68 \end{bmatrix}$$

$$\begin{bmatrix} 352,833 \\ 670,91 \\ 1430,68 \end{bmatrix} = \begin{bmatrix} \lambda_1 \cdot 0,22 \\ \lambda_1 \cdot 0,41 \\ \lambda_1 \cdot 0,88 \end{bmatrix} \Leftrightarrow \begin{array}{l} \lambda_1 = 1603,7 \\ \lambda_1 = 1636,37 \\ \lambda_1 = 1625,77 \end{array} \quad \begin{array}{l} \lambda_2 = 127,65 \\ \lambda_2 = 129,45 \\ \lambda_2 = 127,85 \end{array}$$

$$\lambda_2 = 127$$

$$\lambda_3 = 7,543$$

$$\frac{\lambda_1 + \lambda_2}{\lambda_1 + \lambda_2 + \lambda_3} = \frac{1603,7 + 127,65}{1603,7 + 127,65 + 7,543} = 0,99$$

5.

(1) TRANSFORM THE INPUT

$$\begin{aligned} \text{Transformed Data} &= \text{Feature vector}^T \times \text{Data Adjusted} \\ &= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 10 & -1 \\ -1 & 0 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 2 & -1 & -1 \\ 0 & 1 & -1 \end{bmatrix} \end{aligned}$$

remover media

(2) RECOVER THE ORIGINAL DATA using the most informative component

$$y_1 \downarrow \bar{p}q \quad \lambda_1 > \lambda_2$$

$$\frac{1}{\sqrt{2}} \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{2}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{0}{\sqrt{2}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$$

REGRESSION

1.1.

$$b) \text{MSE} = \frac{\sum_{i=1}^n \text{residuos}}{n} = \frac{4,813^2 + (-1,231)^2 + (-0,081)^2 + (-4,955)^2 + 0,387^2 + 0,16^2}{6} = 8,24$$

$$\text{RMSE} = \sqrt{8,24} =$$

$$\text{MAE} = \frac{1}{n} \sum |z_i - \hat{z}_i|$$

$$\text{MAPE} = \frac{1}{n} \sum \left| \frac{z_i - \hat{z}_i}{z_i} \right|$$

$$1.3. d(x_i, x_j) = |10(a_{i1} - a_{j1}) + (a_{i2} + a_{j2})| / 10$$

similardades

DISTANCE WEIGHTS

$$(3) s_{11} = 1 - \frac{d}{d_{\max}} = 1 - 0 = 1$$

$$w_j \\ s_{12} = 1 - \frac{d(x_1, x_2)}{100} = 1 - \frac{74}{100} = 0,26 \dots$$

$$s = s_{11} + s_{12} + s_{13} + s_{14} + s_{15} + s_{16}$$

(4) DISTANCE WEIGHTS FOR x_1

$$w_{12} = \frac{s_{12}}{s} \quad w_{13} = \frac{s_{13}}{s}$$

$$(5) k=2 \\ x_1 \times 5 \quad x_1 \times 6$$

$$0,39 \times 13,75 + 0,61 \times 18,11$$

$$w_{15} = 0,39 \quad w_{16} = 0,61$$

BICLUSTERING

4.

SEQUENTIAL PATTERN MINING

1.

(1) {1,2,3,4} → 4 items ✓

MAXIMAL SEQ. PAT : not a subsequence of a seq. pattern

CLOSED SEQ. PAT : $a \sqsubset \dots \sqsubset b \sqsubset c \sqsubset d$ with
name support

2. PREFIXSPAN

min abs supp = 3

$$\begin{aligned} a &= 6 \\ \underline{b} &= 2 \\ \underline{c} &= 2 \\ d &= 3 \\ \underline{e} &= 2 \\ p &= 5 \end{aligned}$$

DP
a(cac)cad(c)
(ba)(pb)a
(ab)bpb(ae)
a(cop)d
d(pac)
(oop)(ae)

no def has co-occurrences

DA
a(cad)
fa
fa
(ap)d
d(-p)
(-dp)a

$$\begin{aligned} a &= 5 \\ \underline{d} &= 2 \\ p &= 4 \end{aligned}$$

DD

DAG	Daf
(-d)	

Time series FORECASTING

1. (1) residuals increasing \rightarrow multiplicative model

(2) Irregular component:

- baseline
- / seasonal index

$$\langle -1,42; 0; 2,67; \dots \rangle$$

2.

$$\langle 1,6,10,13,16,19,21,23 \rangle$$

$$(1) \hat{x}(10) = 23$$

(2) one DIFFERENZUNG

$$\begin{aligned} &\begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \end{matrix} \\ &\langle 1,6,10,13,16,19,21,23 \rangle \\ &\langle 5,4,3,3,3,2,2 \rangle \end{aligned}$$

$$\hat{x}(9) = 23 + 2 = 25 \quad \hat{x}(10) = 27$$

$$(3) \hat{x}(t) = 2x(t-2) - 1,2x(t-1)$$

$$\hat{x}(9) = 2x(7) - 1,2x(8) = 2 \times 2 - 1,2 \times 2 = 4 - 2,4 = 1,6$$

$$\hat{x}(10) = 2x(8) - 1,2x(9) = 2 \times 2 - 1,2 \times 1,6 = 2,08$$

$$\hat{x}(10) = 23 + 1,6 + 2,08 = 26,68 \approx 26,7$$

$$\hat{x}(10) = 23 + 1,6 + 2,08 = 26,68 \approx 26,7$$

$$(5) \hat{x}(t) = \underbrace{2x(t-2) - 1,2x(t-1)}_P + \frac{\varepsilon(t-1)}{9} \quad d = \# \text{ differings}$$

$$3. \text{ mean Forecasting Error} = \frac{\sum x - \hat{x}}{n} = \frac{2+3+1+2+3+3+3+1+2+4+1+2}{12} = \frac{27}{12} = 0,583$$

$$\text{mean absolute deviation} = \frac{\sum |x - \hat{x}|}{n} = \frac{2+3+1+2+3+3+3+1+2+4+1+2}{7} = \frac{27}{12} = 2,25$$

→ D

Time series REPRESENTATIONS

1.

1.1.(2) PAA of x_1

$$\langle 1, 4, 2, 2, 4, 1 \rangle$$

↓ normalizar (-media / σ)

$$\langle -1,07; 1,34; -0,26; -0,26; 1,34; -1,07 \rangle$$

$$\frac{-1,07 \times 1 + 1,34 \times 0,2}{1,2} = -0,67$$

$$\frac{1,34 \times 0,8 + (-0,26) \times 0,4}{1,2} = 0,80$$

$$\frac{(-0,26) \times 0,6 + (-0,26) \times 0,6}{1,2} = -0,27$$

$$\frac{(-0,26) \times 0,4 + 1,34 \times 0,8}{1,2} = 0,80$$

$$\frac{1,34 \times 0,2 + (-1,07) \times 0,8}{1,2} = -0,67$$

6 segmentos

↓
5 segmentos

$$\frac{6}{5} = 1,2$$

$$\rightarrow \langle -0,67; 0,8; -0,27; 0,8; 0,67 \rangle$$

(4)

$$\langle 0, 1, 2, 3 \rangle \cdot \langle 1, 3, 4, 2 \rangle$$

CHEBYSIEV: $\max(|0-1|, |1-3|, |2-4|, |3-2|) = 2$

(5) x_2 Hamming

$$\langle 1, 3, 4, 2 \rangle$$

1.2.

(2)

$$\langle \underline{abab'a} \rangle$$

→ motif with support = 2

(3)

aa	ab
ba	bb

0	2/4
2/4	0

(4)

aaa	0
aab	0
aba	2
abb	0
bac	0
bab	1
bbc	0
bbb	0

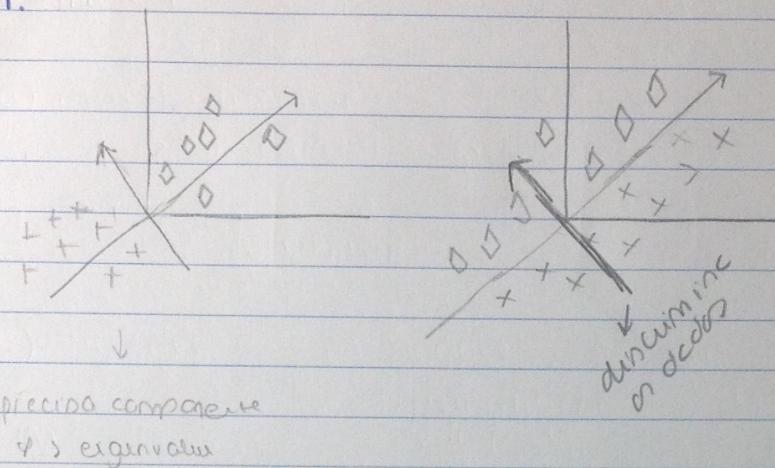
$$\langle abab'a \rangle$$

aba → 2/3 bab → 1/3

PRINCIPAL COMPONENT ANALYSIS

> eigenvector \rightarrow variable

1.



(ii) one of the principal components can accurately discriminate the class.

2. (3)

3. TOP-7 eigenvalues explain 90%.

$$\lambda_1 = 20$$

$$\lambda_2 = 10$$

$$\lambda_3 = 5$$

$$\lambda_4 = 4$$

$$\lambda_5 = 3$$

$$\lambda_6 = 2$$

$$\lambda_7 = 1$$

$$20^2 + 10^2 + 5^2 + 4^2 + 3^2 + 2^2 + 1^2 = \\ = 400 + 100 + 25 + 16 + 9 + 4 + 1 = \\ = 555 \times$$

$$\frac{45}{45+n} = 0,9 \quad \Leftrightarrow 45+n = \frac{45}{0,9} \quad \text{E1}$$

$$\text{E1 } 45+n = 50 \quad \text{E1 } n=5$$

↓

7+5 attributes
↓
12 attributes

4.

covariance matrix $S = \begin{bmatrix} 91,43 & 171,92 & 297,99 \\ 373,92 & 545,21 & 1297,26 \end{bmatrix}$

EIGENVECTORS:

$$\vec{u}_1 = \begin{bmatrix} 0,22 \\ 0,41 \\ 0,88 \end{bmatrix} \quad \vec{u}_2 = \begin{bmatrix} 0,25 \\ 0,85 \\ -0,46 \end{bmatrix} \quad \vec{u}_3 = \begin{bmatrix} 0,94 \\ -0,32 \\ -0,08 \end{bmatrix}$$

$$S\vec{u}_1 = \lambda_1 \vec{u}_1$$

$$(S - \lambda_1 I) \vec{u}_1 = 0 \quad \Leftrightarrow$$

$$\begin{bmatrix} 91,43 & 171,92 & 297,99 \\ 373,92 & 545,21 & 1297,26 \end{bmatrix} - \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_1 & 0 \\ 0 & 0 & \lambda_1 \end{bmatrix} \vec{u}_1 = 0$$

$$\begin{bmatrix} 91,43 - \lambda_1 & 171,92 & 297,99 \\ 0 & 373,92 - \lambda_1 & 545,21 \\ 0 & 0 & 1297,26 - \lambda_1 \end{bmatrix} \begin{bmatrix} 0,22 \\ 0,41 \\ 0,88 \end{bmatrix} = 0$$

$$\begin{bmatrix} (91,43 - \lambda_1) \times 0,22 + 171,92 \times 0,41 + 297,99 \times 0,88 \\ 0 + (373,92 - \lambda_1) \times 0,41 + 545,21 \times 0,88 \\ 0 + 0 + (1297,26 - \lambda_1) \times 0,88 \end{bmatrix} = 0$$

$$= \begin{bmatrix} 20,1146 - 0,22\lambda_1 + 70,4872 + 260,2312 \\ 153,3072 - 0,41\lambda_1 + 418,1848 \\ 1141,5888 - 0,88\lambda_1 \end{bmatrix} = 0$$

$$= \begin{bmatrix} 352,833 - 0,22\lambda_1 \\ 571,492 - 0,41\lambda_1 \\ 1141,5888 - 0,88\lambda_1 \end{bmatrix} = 0$$

$$-0,22\lambda_1 = -352,833 \quad \Leftrightarrow \quad \lambda_1 = 1603,78636$$

$$-0,41\lambda_1 = -571,492 \quad \Leftrightarrow \quad \lambda_1 = 1393,88$$

$$-0,88\lambda_1 = -1141,5888 \quad \Leftrightarrow \quad \lambda_1 = 1297,26$$

5.	y_1	y_2	$\lambda_1 = 3$	$\lambda_2 = 1$
x_1	1	-1		
x_2	0	1		
x_3	-1	0	$\bar{u}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$	$\bar{u}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

(1) Transform the input data using PCA

$$\text{Data Adjusted} = \begin{bmatrix} 1 & 0 & -1 \\ -1 & 1 & 0 \end{bmatrix}$$

$$\text{feature vector} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

$$\begin{aligned} \text{Transformed data} &= \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}^T \begin{bmatrix} 1 & 0 & -1 \\ -1 & 1 & 0 \end{bmatrix} = \\ &= \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}^T \begin{bmatrix} 1 & 0 & -1 \\ -1 & 1 & 0 \end{bmatrix} = \\ &= \begin{bmatrix} 2/\sqrt{2} & -1/\sqrt{2} & -1/\sqrt{2} \\ 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \end{aligned}$$

(2) Recover the original data

$$\text{Data Recovered} = (\text{Feature Vector} \times \text{Transformed Data}) + \text{Original Mean}$$

$$\text{Data recovered} = \text{Data Adjusted} + \text{Original mean}$$

Data recovered

REGRESSION

1.

$$\hat{z} = \hat{\beta}_0 + \hat{\beta}_1 \cdot y_1 + \hat{\beta}_2 \cdot y_2 \rightarrow \text{MULTI}$$

a) Residuals?

Errors

$$\hat{z} = 2,341 + 1,616 \cdot y_1 + 0,014 \cdot y_2$$

$$\hat{z} - z$$

$$\hat{z}_1 = 2,341 + 1,616 \cdot 7 + 0,014 \cdot 560 = 21,493$$

$$\hat{z}_2 = 2,341 + 1,616 \cdot 3 + 0,014 \cdot 220 =$$

⋮

$$21,493 - 16,68 = 4,813$$

$$b) \text{MSE}(\hat{z}, z) = \frac{1}{n} \sum_{i=1}^n (z_i - \hat{z}_i)^2 =$$

$$= \frac{1}{6} \cdot [4,813^2 + (-1,231)^2 + (-0,081)^2 + (-4,955)^2 + 0,387^2 + 0,163^2] = 8,24$$

$$\text{RMSE}(\hat{z}, z) = \sqrt{\frac{1}{n} \sum_{i=1}^n (z_i - \hat{z}_i)^2} = \sqrt{8,24} = 2,8698$$

$$\text{MAE}(\hat{z}, z) = \frac{1}{n} \sum |z_i - \hat{z}_i| = \frac{1}{6} (4,813 + 1,231 + 0,081 + 4,955 + 0,387 + 0,163) \\ = 1,9383$$

$$\text{MAPE}(\hat{z}, z) = \frac{1}{n} \sum \left| \frac{z_i - \hat{z}_i}{z_i} \right| = \frac{1}{6} \left(\frac{4,813}{16,68} + \frac{1,231}{11,5} + \right.$$

$$\left. \frac{0,081}{12,03} + \frac{4,955}{14,88} + \frac{0,387}{13,75} + \frac{0,163}{1,8,11} \right) \approx 0,1287 \\ (3)$$

1.2. $y = mx + b \rightarrow \text{UNIVARIATE}$

$$y_1 = 16,68 + 1,616 \cdot 7 \quad (1)$$

$$1.3. d(x_i, n_j) = 110(a_{i1} - a_{j1}) + (a_{iz} - a_{jz}) / 101$$

$$d(x_1, x_2) = 74$$

$$d(x_1, x_3) = 62$$

$$d(x_1, x_4) = 78$$

$$d(x_1, x_5) = 51$$

$$d(x_1, x_6) = 23$$

$$d(x_1, n_4) = 4$$

$$d(n_1, n_3) = 12$$

(1) VER distância menor em cada x e esse z é o \hat{z}

(2) ver quais pontos mais próximos seu valores de z , pote médio e é o \hat{z}

(3) \leftarrow similaridade

$$s = 1 - \frac{d}{d_{\max}} \quad s_{11} = 1 - 0 = 1$$

$$s_{12} = 1 - \frac{74}{100} = 1 - 0,74 = 0,26$$

...

(4)

distância weight \rightarrow é soma das similaridades da vizinha considerada (da sua parte)

$$0,26 + 0,38 + 0,22 + 0,49 + 0,77$$

$$\frac{\sum \text{similaridade}}{\text{soma total}}$$

$$w_{12} = \frac{s_{12}}{\text{soma similaridades}}$$

\checkmark (5)

$$\hat{z}_1 = 13,75 \times 0,39 + 18,11 \times 0,61 = 16,4 \\ k=2$$

$$s_{15} + s_{16} = 0,49 + 0,77 = 1,26$$

$$w_{15} = \frac{0,49}{1,26} = 0,39$$

$$w_{16} = \frac{0,77}{1,26} = 0,61$$

(6)

2.

2.1.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$x_2 \approx 0,8 \quad +0$

↓

(3)

without regard to y_2

↓
varia!!

2.2.

Least square error. - (3)

3.

Tendencia

var

Ideal: erro



Tendencia: MAU



\bar{P} queles valores c. osnmeder ten mais resultados

↑
subirme sempre o
valor real

(4)

$$\begin{array}{c} 215 \\ 2 + \text{proj} \\ 13 - 75 + 18 - 11 = \end{array}$$

BICLUSTERING

encontra
relações entre variáveis

instâncias

variáveis

1. I: instâncias

J: features

constant
additive
multiplicative
order-preserving

(1) ORDER-PRESERVING:

	y_2	y_3	y_5
x_1	2	5	4
x_2	1	4	2
x_3	2	5	4

ver como variam as features $y_2 \leq y_5 \leq y_4$

(2) CONSTANT

	y_2	y_3	y_5
x_1	2	5	4
x_3	2	5	4

✓

(3) ADDITIVE: p todas as colunas o nome é o mesmo

$$\varphi_B = 2 \quad 1 \quad 4 \quad 2$$

	y_1	y_2	y_3	y_4
x_1	3	2	5	3
x_2	2	1	4	2

$$\varphi_B + 1$$

$$\leftarrow \varphi_B + 0$$

(4) MULTIPLICATIVE

	y_2	y_3	y_5
x_1	2	5	4
x_2	1	4	2
x_3	2	5	4

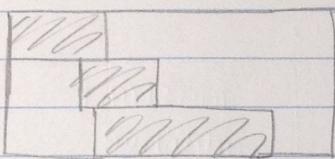
$$y \quad \times$$

$$\leftarrow \varphi_B \times 1$$

$$\leftarrow \varphi_B \times 2 \quad \times$$

(4)

2.



(3) (4)
(5)
(7)

3. (3) False negative !!

4. $\bar{A} \leftarrow$ amplitude of the range of values in a matrix A

Given A, coherence strength is a range of $\in (0,1)$ such that

$$a_{ij} = c_j + \gamma_i + \eta_{ij} \text{ (or other)} \text{ and}$$

$$\eta_{ij} \in [-\delta/2, \delta/2]$$

(1) constante = 2,1 $\bar{\rho} = \eta = \emptyset$ NOT ERROR \rightarrow noise

X (2) $3,2+0,1 = 3,3 < 3,4$ $\eta = (-0,1, 0,1)$

(3)

$$1,2 \ x_2$$

$$1,2 \ x_4$$

$$1,2 \ x_7$$

$$0 \ x_8$$

$$1,2 \ x_9$$

$$\eta = \pm 0,2$$

✓

(4)

	y_3	y_6	y_8	y_9	
x_8	0	1	3,1	0,5	$\leftarrow y_8 = \emptyset$
x_2	1,2	2,2	3,3	1,7	$y_2 = 1,2$
x_4					$y_4 = 1,2$
x_7					$y_7 = 1,2$
x_9					$y_9 = 1,2$

$$\eta \pm 0,2$$

V (4) Quality: amount of noisy elements

$$\frac{\text{noisy rows}}{\text{total number of columns}} = \frac{16-2}{16} = 0,875$$

+ type of noise $\rightarrow [y_{ij} - \delta/2, y_{ij} + \delta/2]$

SEQUENTIAL PATTERN MINING

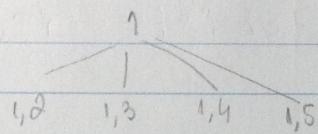
1. SID Sequence

- 1 $\langle \{1,5\} \{2\} \{3\} \{4\} \rangle$
- 2 $\langle \{1\} \{3\} \{4\} \{2,5\} \rangle$
- 3 $\langle \{1\} \{2\} \{3\} \{4\} \rangle$
- 4 $\langle \{1\} \{3,5\} \{5\} \rangle$
- 5 $\langle \{4\} \{5\} \rangle$

6^{SP}
6//

minsup = 0,4

3 items per pattern



(sup) $\frac{2}{5}, \frac{4}{5}, \frac{3}{5}, \cancel{\frac{1}{5}}$

(2) *



PATTERNS ↓

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

$\{1, (3,5)\} \text{ queries}$

sup $\frac{2}{5}, \frac{2}{5}, \frac{3}{5} = 96\% \quad (4) \checkmark$

$\{1,2,3,4\} \rightarrow (1) \checkmark$
 $\text{sup} = 2/5$

$\hookrightarrow \{1,3,5\} \text{ not sup!!} \quad (3) \checkmark$

2.

SID Sequence

- 1 $\langle (a)(a)(c)(a)(c)(d) \rangle$
- 2 $\langle (b)(a)(f)(b)(a) \rangle$
- 3 $\langle (a)(b)(b)(f)(b)(a)(e) \rangle$
- 4 $\langle a(f)(a)(p)(a) \rangle$
- 5 $\langle d(f)(p)(a)c \rangle$
- 6 $\langle (a)(d)(f)(a)(e) \rangle$

minimum absolute support $\rightarrow 3$

PREFIXSPAN

Dd	
1	a(cac)(adc)
2	(ba)(pb)a
3	(ab)bfbae)
4	a(cap)d
5	a(fac)
6	(cadp)(ae)

$$\begin{aligned}
 a &= 6 \\
 b &= 2 \\
 c &= 2 \\
 d &= 3 \\
 e &= 2 \\
 f &= 5
 \end{aligned}$$

co-occurring é só
uma posição!

Da	Dd	Df
a (a)(ad) 1	a (a)(ad) 1	a 2
b (b)(a) 2	(b)(ba) 5	a (a) 3
c f(a) 3	(c) (f)(a) 6	d 4
(c) (ap)d 4		a (a) 6
(c) (ap)(a) 6		

Da	Dd
(a) 1	

Da
(a)(ad)
(f)a
(c) f(a)
(a,f)d
(-f)
(-df)(a)

TIME SERIES FORECASTING

1. $\langle 2, 3, 7, 4, 4, 6, 3, 4, 8 \rangle$ time series
↳ trend

trend component was subtracted \rightarrow ADDITIVE

ADDITIVE MODEL

$$\text{DATA} = \text{seasonal effect} + \text{Trend} + \text{Cyclical} + \text{Residual}$$

MULTIPLICATIVE MODEL

$$\text{DATA} = (\text{seasonal effect}) \times \text{Trend} \times \text{Cyclical} \times \text{Residual}$$

seasonal indexes : $\langle 0.7, 0.8, 1.5 \rangle$

RESIDUALS

- (1) Residuals increase along time suggesting that the trend should be divided instead of subtracted
- (2) subtract baseline and remove seasonal component by dividing values by the seasonal indexes
- (3) Residuals increase along time

$$\hat{x}(t) = -t + t^2$$

$$2 + (-1) + 1^2 = 2$$

$$3 + (-2) + 2^2 = 5$$

$$7 + (-3) + 3^2 = 13 \quad \dots \quad (4) \checkmark$$

(2)

$$\frac{2-3}{0,7} =$$

$$\frac{3-3}{0,8} =$$

$$\frac{7-3}{1,5} =$$

$$\frac{4-3}{0,7} =$$

✓

2. Time series: $\langle 1, 6, 10, 13, 16, 19, 21, 23 \rangle$

naive forecast \rightarrow usa último valor da trend e
faz linhareta
(n'estima noda)

(1) ✓

DIFFERENZUNG \rightarrow diferença entre rada o ponto
 $(6-1, 10-6, 13-10, \dots)$

two time points ahead: $10-1, 13-6, 16-10, \dots$

$$23-19=4$$

$$\hat{x}(10)=23+4=27$$

(2) ✓ perco 1 punto para 1 differenzung
 $\langle 1, 6, 10, 13, 16, 19, 21, 23 \rangle$

$$\hat{x}(9)=23+2=25$$

$$25-23=2$$

$$\hat{x}(10)=25+2=27$$

naive forecast escolhe último

(3)

$$\begin{aligned}\hat{x}(9) &= 2 \times \overset{\text{DIFF}}{x(7)} - 1,2 \times \overset{\text{DIFF}}{x(8)} \\ &= 2 \times 21 - 1,2 \times 2 = 1,6 - 21 = 14,4\end{aligned}$$

$$\begin{aligned}\hat{x}(10) &= 2 \times \overset{\text{DIFF}}{x(8)} - 1,2 \times \overset{\text{DIFF}}{x(9)} \\ &= 2 \times 2 - 1,2 \times 1,6 = 2,08\end{aligned}$$

$$23 + 1,6 + 2,08 = 26,68$$

(4) one DIFFERENZUNG ; 4 5 6 + 8
 1 6 10 13 16 19 21 23 $\leftarrow \textcircled{8}$
 5 4 3 3 3 2 2
 $\epsilon(8)$

$$\hat{x}(9) = 0,8 x(7) +$$

(S) p: n° tempos → não usamos
 d: # differenc.
 q: $\varepsilon(t-q)$

ARIMA

p: AR order

d: integration order - # differencings to stationarize time series

q: MA order

3.

(2, 3, -1, -2, 3, 3, -3, -1, 2, 4, -1, -2) RESIDUES

(1) mean forecasting error

$$MFE = \frac{\sum x_t - \hat{x}_t}{n} = \frac{7}{12} = 0,583 \quad \checkmark$$

(2) mean absolute deviation

$$MAD = \frac{\sum |x_t - \hat{x}_t|}{n} = \frac{27}{12} = 2,25 \quad \checkmark$$

(3)

++ -- → BIAS
 (acada 4)

$12/3 = 4$
 ↓
 grupos 4

(4) X

Time Series Representation

1.

$$x_1 = \langle 1, 4, 2, 2, 4, 1 \rangle \quad x_2 = \langle 1, 3, 4, 2 \rangle$$

1.1. PCL

$$(1) \bar{x}_2 = \frac{1+3+4+2}{4} = \frac{10}{4} = 2,5 \quad \sigma_{x_2} \approx 1,12$$

$$x = x_1 + x_2$$

$$\bar{x} = \frac{1+4+2+2+4+1+1+3+4+2}{10} = \frac{24}{10} = 2,4 \quad \sigma_x = 1,2$$

$$x_2' = \langle -1,34, \dots \rangle$$

↓

ponto-média

σ

$$x_2'' = \langle -1,17, \dots \rangle$$

(2)

PAA representation

↓

Piecewise aggregate approximation

$$\bar{x}_1 = \frac{14}{6} \approx 2,33$$

$$\sigma_1 = 1,247$$

$$x_1 = \langle 1, 4, 2, 2, 4, 1 \rangle \quad \downarrow \text{normalizer} \quad \frac{6}{5} = 1,2$$

$$\langle -1,07; 1,34; -0,26; -0,26; 1,34; -1,07 \rangle \quad \downarrow 5 \text{ segmentos}$$

$$\boxed{\frac{-1,07 + 1,34 \times 0,2}{1,2} \approx -0,67 \quad 1,34 \times 0,8 + (-0,26) \times 0,4}{1,2}$$

$$5 \times (-1,07) + 1 \times 1,34 = -4,01$$

$$4 \times (1,34) + 2 \times (-0,26) =$$

$$3 \times (-0,26) \times 2 =$$

$$2 \times (-0,26) + 4 \times 1,34 =$$

$$1 \times (1,34) + 5 \times (-1,07) =$$

$$\langle -\frac{4,01}{6}, \frac{4,84}{6}, -\frac{1,56}{6}, \frac{4,84}{6}, -\frac{4,01}{6} \rangle$$

$$\langle -0,67, 0,80, -0,27, 0,80, -0,67 \rangle$$

✓

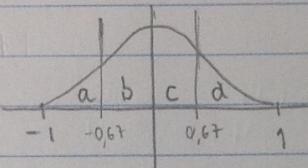
(3) SAX representation : x_2

$$\bar{x}_2 = 2,5 \quad \sigma_2 = 1,12$$

4 segmentos

normalizado: $\langle -1,34; 0,45; 1,34; -0,45 \rangle$

DISTRIBUIÇÃO normal:



$\langle a, c, d, b \rangle$

$= 0,06 + 1,34 \times 0,5$

(4) BIJECTIVE MAPPING

$\langle 0,2,3,1 \rangle \rightarrow \langle 1,3,4,2 \rangle$

CHEBYSHEV: ①

X

(5) HAMMING DISTANCE



normalized by the number of segments

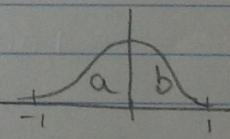
9.2.

(1) CODEBOOK SYMBOLIZATION.

(2) SAX representation

$\langle -0,67; 0,80; -0,27; 0,80; -0,67 \rangle$

a b' a b' a



aba → MOTIF: conjunto das letras comum

support 2 → aparece 2x ✓

(3) SECOND-ORDER BITMAP

2 combinações de letras

aa	ab	→	0/4	2/4	$p_{\bar{q}} = 2 \times ab + 2 \times ba$
ba	bb		2/4	0	Δ normalizar

1. mark frequencies

2. normalize

(4) THIRD-ORDER BITMAP

aaa	aab	aba
abb	baa	bab
bba	bbb	///

aba → 2x

2

3 → pourquoi aba 2 vez es
bab 1 vez

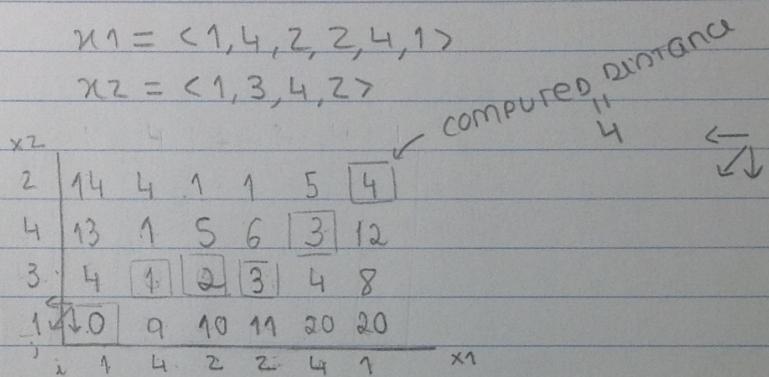
(5) lower bounding property: distances on the new space
always lower than in original space
upper-contrario

1.3

T21 → DTW distance: $\min \left\{ \sqrt{\sum_{k=1}^K w_{ik}} / K \right\}$

$$(1) \quad x_1 = \langle 1, 4, 2, 2, 4, 1 \rangle$$

$$x_2 = \langle 1, 3, 4, 2 \rangle$$



$$\gamma(i, j) = d(q_i, c_j)^2 + \min \{ \gamma(i-1, j-1), \gamma(i-1, j), \gamma(i, j-1) \}$$

(2) minimum length → 6 alignments

↳ n² elementos

(3) ✓

2.1.

(1) Fourier Transform

→ When applying the DFT over a signal, the sampling rate needs to be high enough to accurately model high frequencies

(2) → DFT decomposes a signal in terms of its sinusoids

(3) → DFT representations offer the possibility to derive tabular data mapping from complex signals and temporal information is lost

(4) → X SDFT not preferred over DFT for stationary time series
↓
for non-stationary

(S) Different types of signals are optimally described by
types of wavelets, such as Haar, norlet or
Daubechies waveforms

3. DTW is preferred over DFT:

- signals with discontinuities (abrupt changes)
- signals with drifts and trends
- when the length of stationary segments within
a signal is unknown or varying

TEMPORAL DATA mining

1.	time series
x_1	$\langle 3, 2, 6 \rangle$
x_2	$\langle 2, 5, 4, 2 \rangle$
x_3	$\langle 3, 5, 5 \rangle$
x_4	$\langle 3, 7, 6, 3 \rangle$

1.1. K-medoids - x_1 and x_4

medoid \rightarrow observation that minimizes the total distances within the cluster

DTW with square penalization

\sum

(1)

6	17	2	6	22	como 22 > 9	X
2	1	10	9	6		
x_1	3	1	5	6	7	$\sqrt{1+1+2+6+22} = 1,731$
x_2	2	5	4	2	5	

(2)

6	10	2	3	4	3	10	7	8	3	25	9	9
2	1	9	18	6	4	2	6	9	25	25	5	5
x_1	3	10	9	18	7	26	5	14	31	16	4	8
x_3	3	5	5	x_4	3	1	5	6	7	x_4	3	0
				x_2	2	5	4	2		x_3	3	5

$$\sqrt{0+1+2+3} = 0,612$$

$$\sqrt{1+5+5+6+7+8} = 0,94$$

$$\sqrt{0+4+5+9} = 1,06$$

(3)

(4) K-means is not adequate method to be applied on time series since the mean operation can conceal waveforms in the presence of misalignments

(5) K-means algorithm puts data in hypersphere distance spaces

(6) EM algorithm puts data in hyperelipsoid distance spaces

1.2.

(1) AGGLOMERATIVE + SINGLE LINK

	x1	x2	x3	x4	
x1	0	22	3	11	✓ usando DTW
x2	22	0	11	8	metodo
x3	3	11	0	9	
x4	11	8	9	0	

22+3 → x1(x2)
x1=x2 x2=x3 minima distancia
 ↓ ↓ ↓
 22+41 11+3

x1 x2 x3 x4
[]

	x_1	x_2	x_3	x_4
x_1	0	8	3	5
x_2	8	0	5	6
x_3	3	5	0	5
x_4	5	6	5	0

↓
MATRIZ (STOR)

TEMPORAL DATA
MINING

1.1.

(1) min : 3

x

max : 8

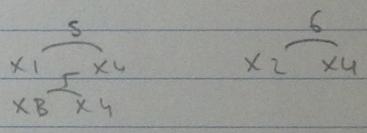
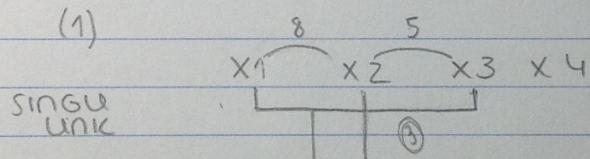
(2)

$x_1 - x_3$

$x_4 - x_2$

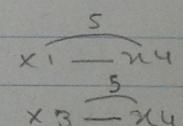
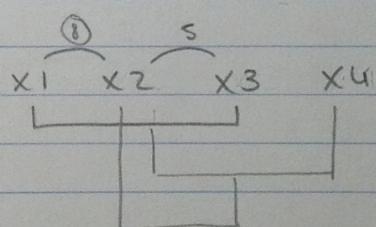
1.2. AGGLOMERATIVE CLUSTERING

(1)



(2)

compute
link



escolher mínimo
entre os máximos

(3) SILHOUETTE SCORE

$d(u_1, u_2, u_3, u_4)$

$d(u_1, u_4) = 5$

$d(u_1, u_2) = 8$

$d(u_1, u_4) = 6$

$d(u_1, u_3) = 3$

$d(u_1, u_4) = 5$

$d(u_2, u_3) = 5$

$$a_{x_1} = \frac{8+3}{2} = \frac{11}{2} = 5,5$$

$$b_{x_1} = 5$$

$$s = \frac{b - a}{a} - 1 = \frac{5}{5,5} - 1 = -0,09$$

$$a \times 2 = \frac{8+5}{2} = \frac{13}{2} = 6,5 \quad b \times 2 = 6 \quad s = \frac{6}{6,5} - 1 = -0,0769$$

$$a \times 3 = \frac{3+5}{2} = \frac{8}{2} = 4 \quad b \times 3 = 5 \quad s = 1 - \frac{4}{5} = 0,2$$

$$a \times 4 = 0 \quad b \times 4 = \frac{5+6+5}{3} = \frac{16}{3} \quad s = 1$$

$$\frac{1 + 0,2 - 0,0769 - 0,09}{4} = 0,26 \quad ??$$

2.	Time series	z	DTW to measure distances with Manhattan distance ↓ l0 function between the feature vectors on two time points.
x_1	$\langle \begin{matrix} y_1 \\ y_2 \end{matrix} \rangle = \langle \begin{matrix} 3,2,2 \\ 2,3,1 \end{matrix} \rangle$	1,2	
x_2	$\langle \begin{matrix} y_1 \\ y_2 \end{matrix} \rangle = \langle \begin{matrix} 2,3,2,3 \\ 2,3,4,1 \end{matrix} \rangle$	0,9	
x_3	$\langle \begin{matrix} y_1 \\ y_2 \end{matrix} \rangle = \langle \begin{matrix} 2,2,2 \\ 4,4,1 \end{matrix} \rangle$	0,6	

2.1.

y_1	2	1	2	1	2
	2	1	2	1	2
x_1	3	1	1	2	2
y_2		2	3	2	3

y_2	1	2	2	3	1
	3	1	0	1	3
x_1	2	0	1	3	4
x_2	2	3	4	1	

(2) ✓

(3) DTW não parametriza !!

(4)

|| Regressão
de UNR
com Time
series :)

2.2. Lazy regression - KNN

(1) uniform weights \rightarrow mean operator

$$(x_1 - x_2) \leftarrow \text{mais próximo} \rightarrow \hat{z}_1 = 0,9$$

\diagdown
 x_3

$$x_2 - x_1 \quad \hat{z}_2 = ?$$

\diagdown
 x_3

$$\begin{aligned}
 (2) \text{RMSE} &= \sqrt{\frac{1}{n} \sum (z_i - \hat{z}_i)^2} = \\
 &= \sqrt{\frac{1}{3} [(1,2 - 0,9)^2 + (0,9 - 1,2)^2 + (0,6 - 0,9)^2]} \\
 &= \sqrt{\frac{1}{3} (0,09 + 0,09 + 0,09)} = 0,3 \quad \checkmark
 \end{aligned}$$

$$\text{MAPE} = \frac{1}{3} \cdot (0,25 + 0,33 + 0,5) = 0,36$$

$$(3) \quad s = 1 - \frac{d}{d_{\max}}$$

$$s_{11} = 1 - 0 = 1$$

$$s_{12} = 1 - \frac{1}{10} =$$

(4) \hat{z}