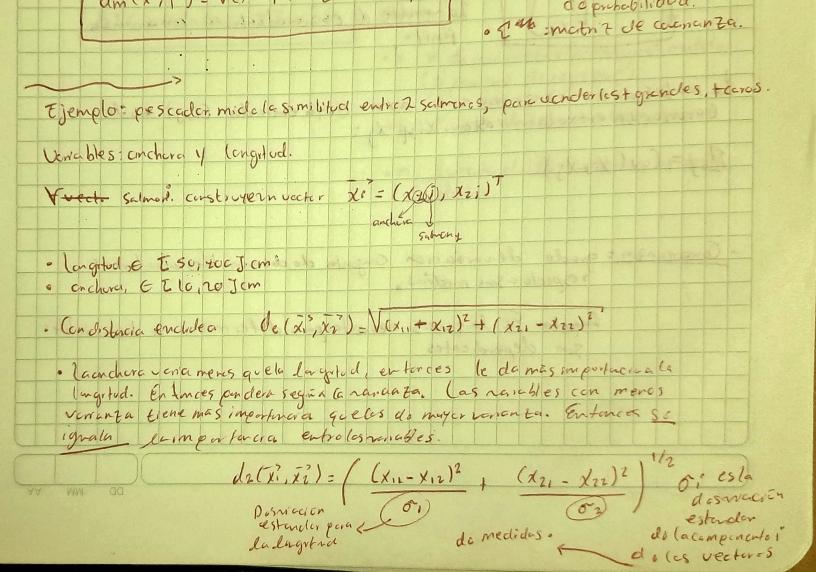
Poistencia de Mahalanobis.
· Es ena medida de distancia. (Mahalanobis, 1936).
· Sirve parci deferminar la similitud entre 2 variables aleatorias multidimen
Sychales.
· A cliferencia de la distancia euclidea, esta tiene en cuenta la corrección ents e
los variables aleafonos.
Distancia de Mahalanchis.
Distancia de vida a lanco s.
dm (x, y) = V(x-y) T Z-1(x-y) (con mismadistribución
de prehabilided.
o 214h: matriz de carnanza.
Ejemplot prescular midala similitual entra 2 salmenas, para vander lest grandes, taros
para acro or rest greates, Factor
Verrables: cinchera y longitud.
Vecto Salmois. construyern vactor Xs = (x(36)), xz;)T
X8 = (X9(1), X2))
archira
· Lancatul C. T. Sching T.
(0 9, 100) Cm
· cochera, E I lo no Jem



on retrever vectors. de (x, x2) - V (x, 2 - x2) 5 5 (x, 2 - x2) o 5: mati z chagenal Colos elementos en la draganal son: 5:5=0;2 · la carchara dependo de la Congretal. Entonces pratenes encomenta la dependercia doles namables, S so cambra per - 2: Matrit de cora, vanta. Distrovade Mahalanchis $dm(x_1,x_1)=(x_1-x_1)^T(x_1-x_2)$ = 2: Mobrit de Corculata. · pubre de cerarrenza: Matriz D quo contrene la coronanza deles entre elementes de un rector D vector aleating | え= 「x, 7 La revolte di esmadeateria concernanta

2: Matrit de cora, vanta. 5 so cambra per Disturvade Mahalanchis dm(x1,x1)=(x1-x1)[1-1(x1-x2) . 2: Mabrit de Coranionza. Mebre de commenza: Matriz D que contrene la comanza detes entre elementes de un vector D vector aleating $\chi = \chi_1 \gamma$ La revoble di Esma aleaterra Con con on fa fruita La mobil de conarranzaseia una non (aga entrada (i, j) es la Coronanta entrola monable Xi y X: Coronanta entrolanonable Xig X; 2; = (ov (x; x;) - Comaranto : grado de varsicioron respecto e sus medias: Conjuntos do des verrables aleatorras Anoila Es leque se necess ta para deferminar son dependientes Si dos remables

				T		1																1	1		-	1			_	-	-	-					
8		Con	1.	250		10	±	hc.	LIAN	100	ord	to 1	non	loca	Lin		(BI	c	-).										-							
H		·	1		10				arre-	Cie		0	17	eci																							
+		6		-	1	1					1	0	1	10.1		1	1			-10	-	wo	001	100													
1		OV	10	-	P	ral	n	nec	hr_	la (olic	luci	Y	110	ica	10	nn	9.	Ser.	20) '	1	1	0		T											
		No	,	di	ce	9	ve	fe	m	Cu	ero	e	sne	u	MEC	bel	0)	312	01	*	en	140	(c)		+			-	+	+						
		00	e	50	4	re	ne	n	(tu	at	es	el	m	efo	1.			1			-	-	-	-	1			-	+	-			-			
		1			1										0																						
		T			F	+	7										IP	Fo	rn	ال	la.	de	1	cm	fen	NC	2	de	F	ond	fer	me	CIT	000	te		
		1	j.	T	1	t	W	0		(0)	-	2	0	ci)			13	dy	e									T		T						
-		1	**) }	-	F	5.	er	10	rej	-	^	Ch	-				1	-					1						1	T	-					T
+	-	-	-	-	-	+	+	1						-			-			-		A	-	-	+	-			-	+	+	-	4		7	-	-
			1						1									10	14	:	Ħ	de	P	cro	im	et	n	59	p c	146	me	10	en	me	cle	lo.	-
						1		1												(p	0	no	h	~0	d	er	es	18	SVC	").			cle		
		T						T	1										1		1	-				1		U		T							T
			1				1	1	1	1	1	1						0	1	1	Т		200	1	-1-	A	7				e la					1	
	+		+	-		-	+	+	+	+	+	+	+	-	+				4	+	t	in	CVC	1	ue	1	LO	iav	m	Ck	C.E.	sel	Uns	as ce l	hick	-	-
	-	-	1	-			-	1	1	1	1	-			1																						
																			n		X	co	na	no	d		la	m	CE	37	70	2	0	1 1	+ (1-	
		1	1	1																1		10	L	-						1	T	1		-		VE	
1		T		1			1	1	1										-	1	-	10	TC	0.	+	+			-							-	
	-	Janes .	1															1	100	1	1					1											