Examen - Parcial 1 - Aprendizoje de máquina - 6CV1 Torres Abonce Luis Miguel

1. c Que estudo el Aprendizoje Automático o de Máquina (Machine Learning)? clubles son sus paradigmas (supervisado, no supervisado y por refuerto)? Puedes responder a las preguntas anteriores realizando un mapa mentol o un resumen min. 1 cuartilla. · El aprendizaje Automático es una rama de la inteligencia artificial que se centra en la construcción de sistemas que pueden aprender de los datos. En lugar de ser propoamados explicitomente, estos sistemas se diseñan poro oiprendet y mejorat con la experiención El aprendizoje outomatico estudia diversos aspectos: · Modelado y optimización. Se enfoca en desas rollos modelos matemáticos y algoritmos que pueden aprender patorones. · Generalización. Busca modelos que oprendon de manera efectiva a partir de conjuntos de datos especificos. · Evaluación y validación Evalua la eficacia y la genera lización de los modelos. Paradigmos principales vel oprendizaje automático 1. Aprendizoje supervisado. Este paradigma implimo el entrenomiento de un modelo en un conjunto de datos etiquetados. Cada ejemplo en el conjunto de da los tiene una entrada y una salida correspondiente. El objetivo del modelo es que aprenda a mapear las antradas a las salidas correctos. 2. Aprendizoje no supervisado. In este poradigmo, el modelo se entreno en un conjunto de jatos sin etiquetos El objetivo del modelo es que descubra estructoras interesoures en los datos por si mismo. 3. Aprendizaje por estuerzo. Este en modigmo implica que un agente aprendo a tomos decisiones realizando acciones en un enformo El agente recibe recompensors o penolitaciones en función de la colidad de sus actione v su objetivo es maximizar la suma total de recompensos.

2. Explica el modelo de ospacio vectorial el cual es amplomente utilizado en el aprendizaje Automática Machine Leorging

Elegi esta pregunta libre

3. ¿ Côme se calculan las medidas de precisión, especific 1290 (recall) y f1 (f-score) que permit en evaluor el rendimiento de un sistema de aprendizaje automoiticoz

	n la matriz de contus Predichot	Predicho -
Verdadero	Casos + correctamente	Casos + incorrectamente
Verdadera l negativo	Casos - correctamente	(a sos - incorrectamente

Precisión

VP = Precisión

Mide lan proporción de predicciones positivas que son correctas. Un valor Alto indica que el modelo tiene un 60jo numero de falsos negativos.

Especificidod (Recall)

VN - Recoil

Mide la proporción de casos negativos que son correctomente identificados. Un valor alto indica que el mode la tien un bajo número de falsos

F1 (f-score)

2 (Perecision · Especificidad) Precision + Especificidad

Un valor alto indica un buen equilibrio entre ambas metricas

A.C.En qué consiste la tarea de clasificación? Mencional tanto las caracteristicas principales como las diferencias en tre las técnicas de clasificación supervisos Y no supervisados. to Se trata de asignar una instancia de datos or una categoria o close específica. Caracteristicos. ·Problema de aprendizaje supervisado. El modelo a prende a partir de un conjunto de dates con ejemplos etiquetodos ·Objetivo. Encontrar un modelo que puedo predecir la clase de nuevas instancias de datos no vistas anteriormente · Tipo de variables. La variable de entrada pueden ser nos. ericos, categoricos o textuales. La variable objetivo es categorica. · Diversidad de aplicaciones, La clasificación se utiliza en dete-cción de correo no desado, diagnóstico médico, reconocimiento Facial, etc. Wassen von Tecnicos de clasificación: 3upervisadas No Supervisadas Corracteristicas No requiere datos Requiere dalos con Etiquetado de etiquetos de close Jafos e tique tados Encontrar patroles y Predecir la cloise de Objetivo estructuros en datos nuevas instancias Ejemplo de Regresión lagica, KNN, K-meons clustering algoritmos SVM Redes neuronoles agrupomiento jerarovio

5. Explica la técnia de validación cruzado de 10 iteraciones la cual sirve para medir la eficiencia de los clasificadores.

· Dividir el conjunto de datos en dicz portes iguales o folds.

· Repetir este proceso diez veces:

-Utilizar a Je los folds como conjuntos de entrenamiento.

-Utilizar el fold restante como conjunto de prueba.

- Entrenoir el closificador en el conjunto de entrenamiento,

- Evaluar el rendimiento del closificador en el conjunto de prueba.

· Calcular la métrica de rendimiento promedio (precisión, F1-scord de 105 10 iteraciones

6. CCual es el teorema de Bayes?

El teorema de Boyes es un formula fondamental en la teoria de la probabilidad que nos permite actualizar la probabilidad de un evento on la luz de nueva información,

P(A | B) = P(B | A) P(A)

Donde:

la hipotesis A dado que observor

· P(AIB) es la probabilidad de observor ta evidencta B si la hipotesis A es verdadera.

· P(B1A) es la probobilidad de observor la evidencia B si la hipotesis A es verdad

· P(A) es la probabilidad de observar la evidencia A . P(B) es la probabilidad de observar la evidencia B El teorema de Bayes nos permite calcular la probabilidad de una hipotesis después de naber observado alguna evidencia.

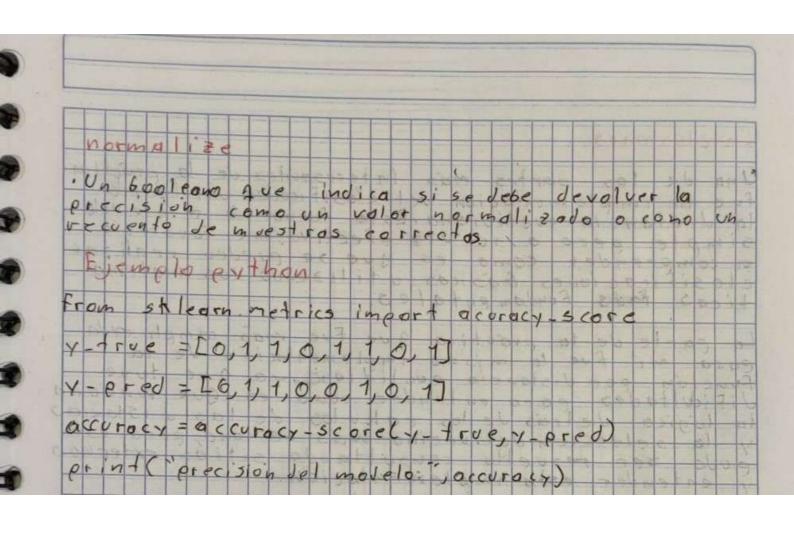
Libraria scikit-lann Make-classification una herramiento La función noke classification es una herramie poderosa en el arsenal de scikit-learn para la seneración de conjuntos de da los simulados des tinodos o problemos de elasificación. es particularmente Util en escenarios Esta función donde se requieven coepr do tos de entrenamiento y price 60 rocidamente para probar algoritmos et aprendizo, e ou bunotico, validar modelos o realizar cornentos Parametros importantes + samples · Este carometro define la contidad de mustra) en el Cinstacias de datos) que se generarán conjunto de datos si mu lados las esencial porq 104trolor el tonoño del conjunto de do fos y puede alusproblema tarse segun el escuario particula o del 'n I features · Indica el número de característicos co atributos que tendros las muestros generados, cada moestra es él conjunto de dot as simulodo estoró ve ar escatado por un vector de la racteristicas con este numero de Jimensiones Controlar el humero de característicos es crucial para simular conjuntos de da dos realistas representativos del problema · Especifico el número de closes distintos en los se eveden closificar los muestros gonerados. Este corametro esto determinado la complejivod del problend Je (lo sifi coció y sinu lodo

Norma

lusters per- closs · Este paraimetro controla la cointidad (o clusters) de contos corclose, si se establece las muestros se distriburán uniformemente entre todas los closes, pero si se ounenta este volor, los closes queden contener subgrupos rondoun state · Si se eroporciono un volor entero se utilizaro como semilla para el generador de nomero pseudoales torios, Estan osegura que el conjunto de datas general o sea reproduíble. Codgo Python from st learn data sets import make classification incort mat plot 116, prelot as elt X y = make = classification (n samples = 1000) n features = 20 n - (| a sses = 3, n - c/v + ess - closs = 2, rand on state - 42) Pl+, figure (Figs ze = (8,6)) elt. > ca He+CX[.0], X[:,]] (=7, cmop = Virilis, = 50, alpha =0.7) elt. x label (Feature 1elt y lobel (Frature 6-) pit. + itle (closificación sintético general a) elt-color Gor (lobel= (lase) elt. show Train_test_split Esta función se utiliza cara dividir un conjunto de Jatos en des subconjuntos: uno para entrenor el made 16 (conjunto de entrenamiento) y otro para evaluaso rendimiento (conjunto de projeta). Este enfoque permite estimar como se generalizació el modelo a nuevos datos la que es fondamental para evalar so capacidod predictiva. Morma

	2 1	197	_		_	_	_	_														_		_	_	-	_	_		T
Paro	In	et	50	5	H	in	10	0	4	11	11	0	1	-								1		-	-	+	-			+
·Ar	ra.	y 5		Fr	6-1		N.	03	1-		. 4	0 80	153		-	13	50	10	1	100	-0	124	-	7.5	100	19	2		915	1
1						C	100		1	1 kg	133	1	1	La.	100		2	0	2	0	b	5	235	3 2	134	13	10		451	1
· E34	e (09,	ra	m	97	0	-	10	0	re	50	4	ta		10	3	6	F	ra	Y	3	(0		2	4	5	Ce	3)	+
De ve	Se Gar	Po	IV	15	au	1	eu		2	04	1	0 1	1-1×	15	0	e		PU	1-1	60	4	qu	2	10	47-	10		1	-	f
para		0/5		0	5	90	He	60	1 5	7	e	25	1	X	1	101	10	un t	00	23	0	0		a	Y	9 4	0	10 e	12	-
9-1-	45	e	1)													7			4							1				I
te	100	2	11	138	-	C	1	- 1	IR	7	1	1	C	1	-3		N	100	-	-7	1	10	5	1	X	N	KUR	9	-	ŀ
199 4						1	1/2	5	1		-	-		Too o										772			2	+		H
· Ind						0	non	2	0	123	10	10	0.0	100		20	5	T. H	1	-	0		3	10	3	30)	0	10	01	
011)	VN7	0	JE		VE	17	05	6	25	19	1006	10	100	4	10	J	9	1	15	0	0	0	1 +	Sing	0	1 5	30	26	336	li
ores	2 0	10	V	0	m	01	0	1	F	103	Q	h 7	1	71	en	7	r	0	0		Y	10	1	10	67		27	re	-	
ores,	n	ta		el		PE	25	ce	10	9	js		0 4	1		Co	14	ju	n	0	_ (10		201	1	25	01	191	10	L
tre	in.	5	1 3	-																							B			
. Det	ine		el	1-7	0	m	a i	10	0	1	01	-	Co	n'j	U	te		de	, 0	6	n	10	e	201	in	i'e	n	10	Si	
10 50	. e	36	ec	1	£,	00	,5	0	CO	110	U	d	0	U	10	m	di	110	-9	4	00	1	e	100	-0	m (0	4	00	-
omp	160	20	nī	0	J	61	C	0 0	2	M	0		12		6.1	- V	06	0		200		33	# 1	1	20	100		4	-10	1
ran	don	1 -	5+	gi	e			9			Đ	12 05	-	7.5		1	61	10		:	n	3		F.D		0	200	13	N. N.	14
				7																										
se en i	Cr	00	Dro	10	40	7	- 1	14		VA	ls	-	-	0	4	er	0,	5	2	U	4	;	1	20	7 1	a	(04	16	
510							21		9.0	ne 1	2	0	95			51	5	5.	40	0	5	9	56	1	00	al	ed	100	rte	2
0 105	-	0	P	C	- 0	10	2	Je	-1	6	10						-	u	0	-		C	101	J	Uh	70	-	-		
100	1	-102	4	1	1	4	2.4	10		21		-0	(ty)	Col	5	-	5	9-		100	7	4	1	200	0	0	3	7	3	
540	#	10	23	-	1	- 6	67	20	7	2	100		HIN	1		1	1	2	1	3	3		20	100	0	L	N		10	
Un	4				-	1	1		1	29		101	-	_ 1		6		I B	1	-		0	1	10	.0		140	100	-	-
leate	60	ny	200	te			25	2	0	10	-	CA	nd	05	9	1	- CO	0	6	- 0	1:	-1	0 6	7 1	0	je	0	7	9	L
sta e						*	001	1	17	+	0	0	,					0	-			V	در	+		-	w.e	10	_	9
			37																					370		1			10	
	1		1	-	-	-	-	-	1	-	1			-					-	-										-
01 2	00	-		100	7	200		500	103	21	19	100	3			1	-	51	+4		0	101		4	100	-	20	20	38	-
	1				-	1			1.09	0			T	-		100	t-	5	27	100	-		-			2/		100	12	
				3					4	0	6		5		1	-	007		100			13					Ti	- 1	12	1
																													12	1
																														F
	-	-	-	_	-	-	-	-	_	-	-		-	_	1															ı

1	6 7																	
rom sk	earn	mod	el	5 00	Hib	2	IME	70		rai	n_	tes	1-	50	11-	1 -9	E HE	3
rom stil	eorn	· dat	0 3 8	+5	i	mpe	ort.	10	av.	IT	10							
	E 92 - E	300	1000	3/3	100		ST	100	15 0	1	10	H.	-	-	3	1	49	
+13=10	100-6	FISC	7	134	100	1	and the	- 1	-	1 3		01	10	1/1/		-	1	
x = ir is.	PATO	1 0	100	1 1	MAL	27 30	10 0		2 3	101	100	1	1	13		0	1	-0
131	arger	Tala	DE!	100	12	X	131	200		130	13	270	1	3	10	1.5	27 2	19
X-train	X-L	est .	1	-		1	and.	-			1	4 4	-	1	CX	- ot	7	0
		7		4 IN	17	- [631		ra	1-	100	1-1	911	T		17	1	1.40
print (Fe	ormo!	del	24	a for	11	4	1		I		1 5	= 0	1.6	+ 1	an	dor.	1 5	tate.
eriat C'Fo	CMO /	1-1-	JA V	1	110	dn	The Chi	m i	enle	17	2	1 1	1	19	In	, 2'	nap	ej
orint () Fo	tun.	1-1	n) v	no	10	9	Vec	97	X	110	11-		1	5 4	0 1	e		
printch	Fran	101 0	24	10	JA		Trev	a C	1.		Di	7-	Tra	in	3	10	رعم	CD S
101010	10 mg		J UN	10	- 63	9	000	we.	-	1	- 7	5	- 3	40	100	2)	300	
															1	1	-	
al Cieuro	101/	5 001	-0														100	
Es una	herr	omic	410	C	on	du	1 m	en t	e	1+1	11	200	10	50	-	50	14	k i t
The second second second					<i>a</i>	1			THE RESERVE									
	10 107 4	TITLO	DIT COL	Aug	100	114 0		100	10 31		10.	- 14				10000		
	The second second	X X			d 010		nu n	+102 1	m3 10	4.7	200	low	10		10		neo Luci	Calab
de mues	Tras	C las	ifi	cod	015	100	TTE	cho	40	te							10	-
Parame	7.505	Impe	rte	247	63	13 12		30	bill	100		-	10	4	50 4			-
and the same of th					200	1000						-	-		03/16/1	3 21		2000
E 6 6 6 6 6 8	3 3 4	2013	01/01	900	(N)	30	78	100	To be seen	0	15		-	0		019		000
y - 1 -	3 3 4	2013	0101		WALE OF THE	1	78	30	2000	90	13	22	-	0		0)1	18	000
y-tr	ve		0100		6015		78	100		2	15	22	k	0				
y - + -	ve	mets	0	re	e re	3 0 0	10	5	las.	e	lig	vet	a		Ve	rde	ade	San
Y-tr	opro	metr	0	re es	o re	300	10.	5	las	e e 5	lig	uet 3 1	a c	s e	Ve	rde	de	Cap
Y-tr	opro	mets	0 0 1	re es.	o re	300	10.	5	las	e e o	light	vet 3	a)ec	5 e	Ve 5	rd.	o de	Cap
Y-tr Este Co va array verdad	opro	nets s co	0 0 1 1 3	re es	ore y	seu e ve	10.	5 U	las	es	ligion	vet 5 l	a)	se gui	Ve 5 c +	rd. er 61	de de	cop
Y-tr Este Co va corray verdod	opro	mets s r c rara	0 0 1 1 3	re es	ore V	sei e ve	10.	5 4 4 6 4	las	e	liga	vet 3 1	a Dec	s e gui	Ve S c t	rd. et 6,	de de	rop h
Y-tr Este (o va array verdod Jatos.	o a ro	nets s c care	0 0 1 1 3	re es tol	9 9	s e u	10.	5 u	las	es	ligas	vet 5 1	a) e o o o o o o o o o o o o o o o o o o	5 6 9 01	Ve S	rd. et 6, 0	de de	cas
Y-tr Este Co va corray verdod	o a ro	mets na cara	0 0 1 1 3	re es. toi	ore V	s e i	la	5 46 0	las	e	lig	vet 3 1	000	s e gu	Ve 5	rd.	de	cas
Y-tr Este (o va verdad Jatos. Y-pr	o o ro	cara	1 3	101	9	ne no	1 85	16 0	149	2 n	03	e	11	701	1	0	5 de	1 to
Y-tr Este (o va verdad Jatos. Y-pr	o o ro	cara	1 3	101	9	ne no	1 85	16 0	149	2 n	03	e	11	701	1	0	5 de	1 to
Y-tr Este Co va verdad Vatos.	o o ro	cara	1 3	101	9	ne no	1 85	16 0	149	2 n	03	e	11	701	1	0	5 de	1 to
Y-tr Este (o va verdad verdad Vatos.	o o ro	cara	1 3	101	9	ne no	1 85	16 0	149	2 n	03	e	11	701	1	0	5 de	1 to
Y-tr Este (o va verdad verdad Vatos.	o o ro	cara	1 3	101	9	ne no	1 85	16 0	149	2 n	03	e	11	701	1	0	5 de	1 to
Y-tr Este (o va array verdod Jatos.	o o ro	cara	1 3	101	9	ne no	1 85	16 0	149	2 n	03	e	11	701	1	0	5 de	1 to



8. Implementación de un clasificador Novive Bayes en Python Explica ¿ para que sirven fit() y predict()? A nota una descripción de sus parametros más importantos Muestra un ejemplo que los Stilice en python entrenamiento. Este conjunto de datos debe contener las variables independientes catributos) y la variable de pendiente (close) opredict (). Utiliza el modelo entrenado pora predecir la clase de nuevos instancios de datos. Porametros importantes · data: El conjunto de datos utilizado para entrenor el modelo o para realizar los predecciones. · Larget: La variable dependiente (clase) del conjunto de datas. ·alpho: Un parametro de regularorzación que ayuda a evitor el 306reojuste. · priors: Los probobilades a priori de las diferentes clases. Ejemplo de uso: from sklearn, naive-boyes import Gaussian NB data = pd. read_csv ("data.csv") X = Jato. drop("clase", ax is=1) y= data [close] mode lo = Gaussian NBC) modelo, fi+(X, y) predicciones = modelo, predict(x) print (clossification-report (x, predictiones))
