















## Definitions Classification techniques Training set Decision trees Collection of labeled data objects used to learn Classification rules the classification model Association rules Test set Neural Networks Collection of labeled data objects used to Naïve Bayes and Bayesian Networks validate the classification model k-Nearest Neighbours (k-NN) Support Vector Machines (SVM) • ... DMG DMG

























































































































## Data mining: classification





























## Data mining: classification

































🌍 E	Bayes	ian clas	sifica	tion:	Exa	ample	Э
	Outlook	Temperature	Humidity	Windy	Class		
	sunny	hot	high	false	Ν		
	sunny	hot	high	true	Ν		
	overcast	hot	high	false	Р		
	rain	mild	high	false	Р		
	rain	cool	normal	false	Р		
	rain	cool	normal	true	Ν		
	overcast	cool	normal	true	Р		
	sunny	mild	high	false	Ν		
	sunny	cool	normal	false	Р		
	rain	mild	normal	false	Р		
	sunny	mild	normal	true	Р		
	overcast	mild	high	true	Р		
	overcast	hot	normal	false	Р		
	rain	mild	high	true	N		
<sup>B</sup> MG	From: Han, Kamber, "Data mining; Concepts and Techniques", Morgan Kaufmann 2006				87		

	Bayesian cla	assification: E	Example		
	outlook				
	P(sunny p) = 2/9	P(sunny n) = 3/5	P(p) = 9/14		
	P(overcast p) = 4/9	P(overcast n) = 0	P(n) = 5/14		
	P(rain p) = 3/9	P(rain n) = 2/5	- ()		
	temperature				
	P(hot p) = 2/9	P(hot n) = 2/5			
	P(mild p) = 4/9	P(mild n) = 2/5			
	P(cool p) = 3/9	P(cool n) = 1/5			
	humidity				
	P(high p) = 3/9	P(high n) = 4/5			
	P(normal p) = 6/9	P(normal n) = 2/5			
	windy				
	P(true p) = 3/9	P(true n) = 3/5			
	P(false p) = 6/9	P(false n) = 2/5			
DB	DMG From: Han, Kamber, Data mining: Concepts and Techniques", Morgan Kaufmann 2006				

















































## Data mining: classification







110





















Class specific measures			
<ul> <li>For a binary classification problem</li> <li>on the confusion matrix, for the positive class</li> </ul>			
Precision (p) = $\frac{a}{a+c}$ Recall (r) = $\frac{a}{a+b}$ F - measure (F) = $\frac{2rp}{r+p} = \frac{2a}{2a+b+c}$			
DMG From: Tan,Sleinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006	121		





No. of Contraction	Hc	w to	build a	ROC curve     Use classifier that produces		
	Instance	P(+ A)	True Class	posterior probability for each		
	1	0.95	+	test instance P(+ A)		
	2	0.93	+	<ul> <li>Sort the instances according to</li> </ul>		
	3	0.87	-	P(+ A) in decreasing order		
	4	0.85	-	Apply threshold at each unique		
	5	0.85	-	value of P(+IA)		
	6	0.85	+	<ul> <li>Count the number of TP_EP</li> </ul>		
	7	0.76	-	TN_EN at each threshold		
	8	0.53	+	- TP rate		
	9	0.43	-	TPR = TP/(TP+FN)		
	10	0.25	+	= FP rate		
				FPR = FP/(FP + TN)		
Ľ	DMG From: Tan, Steinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006					



