No of Pages : 4 Course Code : 09XT83

Roll No:

(To be filled in by the candidate)

PSG COLLEGE OF TECHNOLOGY, COIMBATORE - 641 004

SEMESTER EXAMINATIONS, MAY/JUNE 2011

PhD(PT) - MATHEMATICS

09XT83 DATA MINING

Time: 3 Hours Maximum Marks: 100

INSTRUCTIONS:

 Answer ALL questions from PART - A and Answer any 4 questions from PART -B. Question under PART - C is Compulsory.

PART - A Marks: $10 \times 3 = 30$

- Consider a naïve Bayes classifier with 3 boolean input variables X1, X2, X3 and one Boolean output, Y. How many parameters must be estimated to train such a naïve Bayes classifier?
- 2. What is the apriori principle and how is it applied to frequent itemset generation?
- 3. Why is Gain ratio said to be a better attribute selection measure than information gain?
- What is the significance of feature selection in preprocessing stage of data mining. List the various techniques.
- Consider the following training set in the 2-dimensional Euclidean space:

X	Y	class
-1	150	-
0	1	+
0	2	SON
1	-1 ₂ G	-
1	0	+
,10%	2	+08
2	2	(500
2	350	+

What is the prediction of the 3, 5, 7-nearest-neighbor classifier at the point (1,1)?

- Differentiate between supervised and unsupervised learning with an example.
- For the dataset shown below,
 - i) Compute the support for itemsets {e}, {b,d} and {b,d,e} by treating each trasaction ID as a market basket.

No of Pages : 4 Course Code: 09XT83

PSG TECH PSG TECH ii) Compute the confidence for the association rules {b,d}→ {e} and {e} → {b,d}. Is confidence a symmetric measure?

Tra	ansaction ID	Items bought
:000	01	{a,d,e}
002	24	{a,b,c,e}
00	12	{a,b,d,e}
003	31	{a,c,d,e}
00	15	{b,c,e}
002	22	{b,d,e}
000	29	{c,d}
004	40	{a,b,c}
003	33	{a,d,e}
00:	38	{a,b,e}

- 8. Define a core object and density reachability. Is density reachability a symmetric function?
 - 9. Given 5-dimensional numeric samples A= (1, 0, 2, 5, 3) and B= (2, 1, 0, 3, -1), find
 - i.) The Euclidean distance between points
 - ii.) The city block distance
 - iii.) The Minkowski distance for p=3
 - 10. What is Jaccard coefficient? What is its significance in clustering?

PART - B

- a) Why is naïve Bayesian classification called 'naïve"? Briefly outline the major ideas of naïve Bayesian classification. (4.5)
 - b) Using the data given below perform classification using naïve Bayes classifier.

Consider the following new instance to be classified

Magazine promotion =Yes

Watch promotion = Yes

Life insurance promotion = No

Age = 45

Credit card insurance = No

Sex = ?

Magazine promotion	Watch promotion	Life insurance promotion	Age	Credit card insurance	Sex
Yes	No	No. CV	45	∠ No	Male
Yes	Yes	Yes	40	Yes	Female
No .	No No	_CNo	42	No ~G	Male
Yes	Yes	Yes Yes	30	Yes	Male
Yes	No_\	Yes	38	No	Female
_ No	No	No 2	55	No	Female
Yes	Yes	Yes	35	Yes	Male
No	S No	_ No	27	No	Male
Yes	No	No	43	No	Male
Yes	Yes	Yes	41	No	Female

No of Pages : 4 Course Code: 09XT83

12. Briefly explain the Apriori principle of frequent itemset mining. Consider the sample database shown below. Assume min support count =3. Show the step by step process PSG TECH PSG TECH for Apriori algorithm. A1 to A9 are the items purchased and 1 or 0 indicates whether the item is purchased(1) or not(0) for a given transaction.

								_	-
	A1	A2	А3	A4	A5	A6	A7	A8	A9
100	1	0	0,9	0	1	1_	0	1	0
2	0	1	0	1	0	0	0	1	0
3	0	0	0	1	18	0	1	0	0
4	0	0	1	0.5	0	0	0	0	0
5_0	0	0	0_0	0	1	1	J.C	0	0
6	0	1	ما	1	0	0 \	0	0	0
7	0	J.	0	0	0	<u></u>	1	0	10
8	0	0	0	0 /	ď	0	0	0 6	0
9	0	0	0	0	0	0	0	1	0
10	0	0	€,	1	1	0 <	SE.	0	0

- TECH PSG TECH a) Explain the fundamental difference between the Bagging and Boosting ensemble learning methods? How do these notions relate to the concept of generating a good [8.5] ensemble? What are the advantages and disadvantages of each method? Explain any two methods for evaluating the accuracy of a classifier. [4]
 - 14. The distance matrix shown below gives the distances in kilometers between some Italian cities. Use single link and complete link agglomerative clustering to group the data described by the following distance matrix. Show the dendograms.

		ВΛ	FI	МП	NA	RM	го
L	BA	0	552	977	255	4.2	SSE
	FI	662	0	293	468	268	400
L	МП	877	295	C	754	564	138
	NA	255	468	734	0	219	369
	RM	4.2	250	564	219	0	ees
	TO	996	400	138	869	669	C

No of Pages: 4 Course Code: 09XT83

YECH PSG TECH 15. NASA wants to discriminate between Martians(M) and Humans(H) based on following characteristics: Green ? (N,Y), Legs ? (2,3), Height ? (S,T), Smelly ? (N, Y) Our available training data is as follows:

Learn a decision tree using the ID3 algorithm and draw the tree.

	Species	Green	Legs	Height	Smelly	
1)	м	И	3	s	Y	
2)	М	Y	2	Т	Ħ	
3)	М	Y	3	Т	Ħ	
4)	M	N	2	3	Y	
5)	M	Y	3	Т	Ħ	
6)	Н	N	2	T	Y	
7)	Н	N	2	S	N	
8)	н	N	2	Т	n	
9)	н	Y	2	S	n	
10)	H	И	2	T	Y	

PART - C Marks: $1 \times 20 = 20$

PSG TECH PSG TECH PSG TECH

What is Association rule mining? Explain using terms as support, confidence and frequent itemset. On what factors does FP Growth algorithm perform better than Apriori for Association rule mining? For the given transaction data set construct the FP tree PSG TECH PSG TECH and give a step by step process of generating the frequent item sets. (Show the conditional pattern base and conditional FP tree for every item). Min support =2.

7H 62		nsaction ata Set	
CH PSC	TID 1 2 3 4 5 6 7 8 9	tems (a,b) (b,c,d) (a,c,d,e) (a,d,e) (a,b,c) (a,b,c) (a,b,c) (a,b,c) (a,b,c) (b,c,e)	2
PSG TEC	۸ .	Page 1	įÇ)

PSG TECH PORL

PSG TECH PSG TECH