



USN М

M S RAMAIAH INSTITUTE OF TECHNOLOGY

(AUTONOMOUS INSTITUTE, AFFILIATED TO VTU) **BANGALORE - 560 054**

SEMESTER END EXAMINATIONS -JANUARY 2016

Course & Branch : B.E.- Information Science & Engg.

Semester : VII

Subject

Data Mining

Max. Marks: 100

Subject Code

Duration

: 3 Hrs

Instructions to the Candidates:

· Answer one full question from each unit.

: IS711

UNIT - I

1 Discuss challenges that motivate development of data mining. C01 (10)

Consider the following set of frequent 3 item sets:

C01 (04)

 $\{1,2,3\},\{1,2,4\},\{1,2,5\},\{1,3,4\},\{1,3,5\},\{2,3,4\},\{2,3,5\},\{3,4,5\}$ Assume that there are only five items in the data set.

i) List all candidate 4-itemsets obtained by a candidate generation procedure using the F k-1 x F1

ii) List all candidate 4-itemsets obtained by the candidate generation procedure in Apriori.

Discuss at least three ways of dealing with missing data listing and their C01 advantages and disadvantages.

2 Consider the following transactional table to find out the frequent tem sets C01 (08) using Apriori algorithm with the minimum support =20%

TID	List of item- ids
T100	I1,I2,I5
T200	I2,I4
T300	I2,I3
T400	I1,I2,I4
T500	I1,I3
T600	12,13
T700	I1,I3
T800	I1,I2,I3,I5
T900	I1,I2,I3

Discuss the three issues on measurement and data collection aspect of data (06)quality

Indicate the type of the data preprocessing task and elaborate.

C01 (06)

 Replace all the transactions at a given store with a single storewide transaction

ii) Transform a continuous attribute to a categorical attribute



IS711

(06)

UNIT - II

O JANUAR DE LA CONTRACTION DEL CONTRACTION DE LA CONTRACTION DE LA

The following contingency table summarizes supermarket transaction data C02

	Hot dogs	hotdogs	Σrow
hambur gers	2000	500	2500
hambur gers	1000	1500	2500
Σrow	3000	2000	5000

- a) Suppose that the transaction rule "hot dogs →hamburgers" is mined. Given a minimum support threshold of 25% and a minimum confidence threshold of 50%
- b) Based on the given data, is the purchase of hot dogs independent of the purchase of hamburgers? If not, what kind of correlation relationship exists between the two?
- b) Discuss the various methodologies for applying association analysis to C02 (06) continuous data.
- c) Define maximal frequent item sets. Illustrate with an example maximal C02 (08) frequent item sets.
- 4 a) Define the following terms C02 (04) i)Discretization-based approach ii) Scaling Property
 - b) Describe the inversion and null addition properties of the objective C02 (06) measures. Indicate a measure for each that satisfies the properties.
 - c) a)Draw a contingency table for each of the following rules using the C02 (10) transactions shown below

Tid	1	2	3	4	5	6	7	8	9	10
Ite ms	{a,b ,d,e }	{b, c,d }	{a,b, d,e}	{a,c ,d,e }	{b, c,d ,e}	{b,d ,e}	{c,d }	{a,b,c }	{a,d, d}	{b,d }

Rules $\{b\}\rightarrow\{c\},\{a\}\rightarrow\{d\},\{b\}\rightarrow\{d\},\{e\}\rightarrow\{c\},\{c\}\rightarrow\{a\}$

- a) Use the contingency tables in part (a) to compute the rank the rules in decreasing order according to the following measures
- i)Support
- ii)Confidence
- iii)interestfactor(X \rightarrow Y) = $\frac{P(X|Y)}{P(X)}P(Y)$

UNIT - III

5 a) Consider the data set shown below

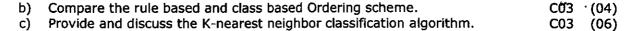
Instance Α В C Class 0 0 1 2 1 0 1 + 3 0 1 0 4 1 0 0 _ 5 1. Ō 1 6 0 0 1 + 1 1 7 0 -8 0 0 0 9 0 0 1 + 10 1 1

C03 (10)





- i) Estimate the conditional probabilities for P(A=1|+), P(B=1|+), P(C=1|+), P(A=1|-), P(B=1|-) and P(C=1|-).
- ii) Use the conditional probabilities in part(a) to predict the class label for a test sample (A=1,B=1,C=1)using the Naïve Bayes Approach.
- iii) Compare P(A=1),P(B=1) and P(A=1,B=1).State the relationships between A and B



6 a) Consider the training examples shown below in Table for a binary C03 (08) classification problem

Insta nce	a1	a2	a3	Target Class
1	T	T	1.0	+
2	T	T	6.0	+
3	T	F	5.0	-
4	F	F	4.0	+
5	F	Т	7.0	-
6	F	Т	3.0	-
7	F	F	8.0	_
8	T	F	7.0	+
9	F	Т	5.0	-

i)What is the entropy of this collection of training examples with respect to the positive class?

ii)What are the information gains of a1 and a2 relative to these training examples?

- b) Provide the Sequential Covering algorithm that is used to build a rule based C03 (06) classifier.
- c) What is boosting? State why it may improve the accuracy of decision tree C03 (06) induction.

UNIT - IV

7. a) Explain any three different types of clusters.

C04 (06)

(06)

b) Use the similarity matrix in the following table to perform single and complete link hierarchical clustering. Show your results by drawing a dendrogram. The dendrogram should clearly show the order in which the points are merged.

	P1	P2	P3	P4	P5
P1	1.00	0.10	0.41	0.55	0.35
P2	0.10	1.00	0.64	0.47	0.98
Р3	0.41	0.64	1.00	0.44	0.85
P4	0.55	0.47	0.44	1.00	0.76
P5	0.35	0.98	0.85	0.76	1.00

c) Explain the DBSCAN algorithm. Illustrate with an example.

C04 (08)



IS711

8. a) Highlight the strength and weakness of K-means algorithm. C04 (06)
b) Given a table where patients are described by binary attributes. Excluding C04 (06)
name, all other attributes are asymmetric binary variables.

Name	Attr1	Attr2	Attr3	Attr4	Attr5	Attr6
J	1	0	1	0	0	0
M	1	0	1	0	1	0
k	1	1	0	0	0	0

Compute the distance between each pair of the 3 objects. Provide suitable explanation.

c) Explain the single link or MIN version of Hierarchical clustering with an C04 (08) example.

UNIT - V

9	a)	List the challenges posed by the web for effective resource and knowledge discovery. Describe ways of resolving these challenges.	C05	(10)
	b) c)	Write a note on mining the web page Layout structure. Discuss the role of data mining in financial sector.	C05 C05	(04) (06)
10	a) b)	Discuss in detail the application of data mining for the retail industry. Write a short notes on i) Web usage mining	C05 C05	(06) (08)
	c)	 ii) Spatial Data mining What are the basic measures for text retrieval? Discuss the Text retrieval methods in detail. 	C05	(06)
