POKHARA UNIVERSITY

Year: 2021

Level: Masters Semester- Spring/Summer

Programme: M.Sc. Computer Science / MCIS

Course: Data Mining and Data Warehousing

Full Marks: 100

Time: 4 hrs

Candidates are required to give their answers in their own words as far as practicable. The figures in the margins indicate full marks.

Attempt all the questions

ON		Moster
Q.N. 1.	(a) 'Visualization/Interpretation' is one of the various steps in Data Mining, why it is so important from	Marks [6]
1.	the perspective of knowledge discovery through data-processing?	[0]
	(b) Data mining is defined as a multi-disciplinary subject, as it combines multiple technology concepts	[6]
	of Computer Science topics. List out at least four different such topics with brief explanation.	[0]
2.	What is datacube in data warehousing? What additional features and processes that datacube can	[3+3+4]
2.	provide? Explain briefly with appropriate example having 4-D cuboid and lattice structure.	= 10
3.	Consider below data about customer complaint on e-commerce portal and their service. Draw a free-	[8]
	hand sketch of Pareto diagram with clearly showing necessary calculations.	
	Code Complaint Types Count	
	G Unexpected fees 37	
	F Unresponsive customer support 142 E Lack of security 76	
	D Poor display on mobile 105 100	
	C Wrong product delivery 55 50	
	B Product different than display 90 0	
	A Slow website speed 195 G F E D C B A	
4.	Here is the sample transaction data records of a food and snacks serving stall on a particular morning of	4+4+4
	a day. Assume, the stall transaction record considers the each single bill record only, which might	= 12
	consists of either group of customers purchase or single customer purchase. Considering minimum	
	support 30% and minimum confidence of 50%, answer the followings:	
	Code Item Trx. Item List (a) Using Apriori algorithm for frequent pattern mining,	
	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	
	C Coffee T3 M,P,S,T thresholds. (You need to show all the steps of calculations.)	
	M Momo	
	Omelette T6 M,N,S,P,T antecedent (i.e. $X \land Y \rightarrow Z$).	
	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	
	S Samosa T9 B,M,T,C,O thresholds are lowered? If yes, write the new rule(s), if not	
	why, justify.	
5.	Compare the followings: (any Three)	[5*3]
	a. Data normalization using Min-max vs. Z-score method	= 15
	b. Text mining vs. web mining	
	c. Star schema vs. Galaxy schema	
	d. Trend vs Seasonality component of time series data	
6.	(a) What are the advantages of Naïve Bayes algorithm based classification?	5
	(b) Using the below data on insurance claims; find the followings:	5+5
	Vehicle Type Cover Years Claim (i) Data column of 'Year' indicates the particular customer who	
	SUV Full 6 No has been continuing with us for their vehicle insurance,	
	Compact Full 4 Yes Which is a continuous data type. Find out the best split	
	suv Full 3 Yes value on year (like Years < X and Years ≥) either using	
	Compact Partial 2 No GINI or Entropy.	
	SUV Partial 5 No (ii) Show your Naïve Bayes working based classification using	
	suv Full 2 Yes the result of (i) for new data, where Vehicle Type is 'SUV',	
	Compact Full 4 No Cover is 'Partial' and Years is '2'.	
7	Compact Partial 5 No	0
7.	What is 'curse of dimensionality'? Why does subspace clustering become interesting problem in high	8
0	dimensional data? Explain with appropriate example.	0
8.	Explain the basics of DBSCAN algorithm and compare with any portioning clustering method.	8 [2*4]
9.	Write short notes on the following (any Three):	[3*4]
	a. F_1 and F_α score	12
	b. Bootstrap method for model validation	12
	c. Outlier data identification using Box plotd. Ensemble method for classification	
		1