

## CSE 40647/60647 Data Science (Fall 2017)

### Mid-term Exam

Instructor: Meng Jiang

(75 minutes, 100 marks, double sided reference, brief answers)

Name:

NetID:

Score:

1. [12] Introduction.

Name at least 4 steps in “Data Science Research” or called “Knowledge Discovery from Data” (KDD).

2. [18] Data processing – Measures.

- (a) [9] (Distance measures) Given two data objects and four attributes/features, we have feature vectors of the two data objects as  $(7, 4, -2, 1)$  and  $(4, 5, -1, 6)$ . Please **calculate three** specific Minkowski distance measures between the two objects and **give the measures’ names**. (Hint:  $4^2 = 16, 5^2 = 25, 6^2 = 36, 7^2 = 49, 8^2 = 64, 9^2 = 81, 10^2 = 100$ )

- (b) [9] (Correlation measures) Give one example wherein **Kulczynski measure** between two variables  $A$  and  $B$  is *more appropriate* than **Chi-square test**  $\chi^2$ : You are asked to (1) explain your variables  $A$  and  $B$ , (2) give an equation to define the Kulczynski measure, (3) explain why Kulc measure is more appropriate in this example.

3. [30] Data warehousing, OLAP, and data cube computation.  
Suppose the base cuboid of a data cube contains two cells

- $(a_1, a_2, a_3, a_4, a_5, a_6) : 1,$
- $(a_1, \mathbf{b_2}, a_3, \mathbf{b_4}, a_5, \mathbf{b_6}) : 1.$

where  $a_i \neq b_i$  for any dimension  $i \in \{2, 4, 6\}$ . Assume each dimension contains no concept hierarchy (i.e., has a single level). (Hint:  $2^3 = 8, 2^4 = 16, 2^5 = 32, 2^6 = 64$ )

- (a) [6] How many **nonempty cuboids** are there in this data cube?
- (b) [6] How many **nonempty closed cells** are there in this data cube?
- (c) [6] How many **nonempty aggregated closed cells** are there in this data cube? What are they?
- (d) [6] How many **nonempty aggregated cells** are there in this data cube?
- (e) [6] If we set **minimum support = 2**, how many **nonempty aggregated cells** are there in the corresponding **iceberg cube**?

4. [40] Frequent pattern and association rule mining.

A data set shows 100 transactions in 5 days, each being summarized as a set of items associated with the number of transactions. Let *relative minimum support* to be  $min\_sup = 0.5$  and *minimum confidence* to be  $min\_conf = 0.6$ . **Again, here we have 100 transactions, not just 5!!!**

date	items_bought	number of transactions
10/15	{a, b, c, m, p}	15
10/16	{b, e, f, p}	35
10/18	{a, c, k, p}	15
10/20	{b, e, p}	15
10/21	{a, e, g, p}	20

- (a) [10] List the frequent 1-itemset associated with their absolute counts.
- (b) [10] Draw the first frequent pattern tree (FP-tree) constructed and used in FP-Growth for the dataset. The tree is NOT for any conditional pattern base.

(c) [10] Present **all** the frequent  $k$ -itemsets for the **largest**  $k$ . Only list frequent itemsets of the largest size. The number of the largest frequent itemsets can be one, two, or many: Please list all of them.

(d) [10] Compute *relative* support and confidence on the following two rules. Are they good **association rules**? (Hint: compare with *min\_sup* and *min\_conf*.)

- i.  $pa \rightarrow b$ , i.e.,  $\{p, a\} \rightarrow \{b\}$ ;
- ii.  $p \rightarrow e$ , i.e.,  $\{p\} \rightarrow \{e\}$ .