

CSC442 Knowledge Discovery and Data Mining

Class Assignment – Monday 16th February 2015

The following questions are extracted from the text – *Data Mining: Concepts and Techniques, Second Edition*, by Jiawei Han and Micheline Kamber.

- Each of you should read chapter 3 and 4 keenly. (You all have an electronic copy of the textbook, or you can find it on the class Google folder)
- I will require each student to pick one question in Chapter 3 – Pg152. You will create 1-3 slides that you will use to explain your/response answer. These questions require a bit of thinking. I will give marks for these.
- I have randomly allocated each one of you one question.

Q3.1 - Martin

Q3.2 – Ian Brayoni

Q3.3 Olive

Q3.4 Steve Waweru

Q3.5 Joan Kirui

Q3.6 Collins

Q3.7 Osebe

The questions (from the textbook):

3.1 State why, for the integration of multiple heterogeneous information sources, many companies in industry prefer the update-driven approach (which constructs and uses data warehouses), rather than the query-driven approach (which applies wrappers and integrators). Describe situations where the query-driven approach is preferable over the update-driven approach.

3.2 Briefly compare the following concepts. You may use an example to explain your point(s).

- (a) Snowflake schema, fact constellation, starlet query model
- (b) Data cleaning, data transformation, refresh
- (c) Enterprise warehouse, data mart, virtual warehouse

3.3 Suppose that a data warehouse consists of the three dimensions time, doctor, and patient, and the two measures count and charge, where charge is the fee that a doctor charges a patient for a visit.

- (a) Enumerate three classes of schemas that are popularly used for modeling data warehouses.
- (b) Draw a schema diagram for the above data warehouse using one of the schema classes listed in (a).
- (c) Starting with the base cuboid [day, doctor, patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2004?
- (d) To obtain the same list, write an SQL query assuming the data are stored in a relational database with the schema fee (day, month, year, doctor, hospital, patient, count, charge).

3.4 Suppose that a data warehouse for Big University consists of the following four dimensions: student, course, semester, and instructor, and two measures count and avg grade. When at the lowest conceptual level (e.g., for a given student, course, semester, and instructor combination), the avg grade

measure stores the actual course grade of the student. At higher conceptual levels, avg grade stores the average grade for the given combination.

(a) Draw a snowflake schema diagram for the data warehouse.

(b) Starting with the base cuboid [student, course, semester, instructor], what specific OLAP operations (e.g., roll-up from semester to year) should one perform in order to list the average grade of CS courses for each Big University student.

(c) If each dimension has five levels (including all), such as “student < major < status < university < all”, how many cuboids will this cube contain (including the base and apex cuboids)?

3.5 Suppose that a data warehouse consists of the four dimensions, date, spectator, location, and game, and the two measures, count and charge, where charge is the fare that a spectator pays when watching a game on a given date. Spectators may be students, adults, or seniors, with each category having its own charge rate.

(a) Draw a star schema diagram for the data warehouse.

(b) Starting with the base cuboid [date, spectator, location, game], what specific OLAP operations should one perform in order to list the total charge paid by student spectators at GM Place in 2004?

(c) Bitmap indexing is useful in data warehousing. Taking this cube as an example, briefly discuss advantages and problems of using a bitmap index structure.

3.6 A data warehouse can be modeled by either a star schema or a snowflake schema. Briefly describe the similarities and the differences of the two models, and then analyze their advantages and disadvantages with regard to one another. Give your opinion of which might be more empirically useful and state the reasons behind your answer.

3.7 Design a data warehouse for a regional weather bureau. The weather bureau has about 1,000 probes, which are scattered throughout various land and ocean locations in the region to collect basic weather data, including air pressure, temperature, and precipitation at each hour. All data are sent to the central station, which has collected such data for over 10 years. Your design should facilitate efficient querying and on-line analytical processing, and derive general weather patterns in multidimensional space.

3.8 A popular data warehouse implementation is to construct a multidimensional database, known as a data cube. Unfortunately, this may often generate a huge, yet very sparse multidimensional matrix. Present an example illustrating such a huge and sparse data cube.