

# Data Warehousing Assignment

This problem set consists of two data modeling scenarios. You will be asked to analyse the strengths and weaknesses of some design alternatives for each scenario. Short answers are fine – one or two paragraphs per question. would be an appropriate length.

## Scenario I:

In this scenario, we are interested in modeling student enrollment in Stanford courses. We would like to answer questions such as:

- Which courses are most popular? Which instructors are most popular?
- Which courses are most popular among graduate students? Undergraduates?
- Are there courses for which the assigned classrooms is too large or too small?

We are planning to have a course enrollment fact table with the grain of one row per student per course enrollment. In other words, if a student enrolls in 5 courses there will be 5 rows for that student in the fact table. We will use the following dimensions: Course, Department, Student, Term, Classroom, and Instructor. There will be a single fact measurement column, EnrollmentCount. Its value will always be equal to 1.

We are considering several options for dealing with the Instructor dimension. Interesting attributes of instructors include FirstName, LastName, Title (e.g. Assistant Professor), Department, and TenuredFlag. The difficulty is that a few courses (less than 5%) have multiple instructors. Thus it appears we cannot include the Instructor dimension in the fact table because it doesn't match the intended grain. Here are the options under consideration:

### Option A

### Option B

### Option C

Modify the Instructor dimension by adding special rows representing instructor teams. For example, CS276a is taught by Manning and Raghavan, so there will be an Instructor row representing "Manning/Raghavan" (as well as separate rows for Manning and Raghavan, assuming that they sometimes teach courses as sole instructors). In this way, the Instructor dimension becomes true to the grain and we can include it in the fact table.

Change the grain of the fact table to be one row per student enrollment per course per instructor. For example, there will be two fact rows for each student enrolled in CS276a, one that points to Manning as an instructor and one that points to Raghavan. However, each of the two rows will have a value of 0.5 in the EnrollmentCount field instead of a value of 1, in order to allow the fact to aggregate properly. (Enrollments are "allocated" equally among the multiple instructors.)

Create two fact tables. The first has the grain of one row per student enrollment per course and doesn't include the Instructor dimension. The second has the grain of one row per student enrollment per course per instructor and includes the Instructor dimension (as well as all the other dimensions). Unlike Option B, the value of EnrollmentCount will be 1 for all rows in the second fact. Tell warehouse users to use the second fact table for queries involving attributes of the instructor dimension and the first fact table for all other queries.

Please answer the following questions.

Question 1. What are the strengths and weaknesses of each option?

Ans : --

Option A:

Advantage – Since there are two different tables, There would be no confusion

Disadvantage – It would require us to do a join of both table to get our required results. So this would cause a efficiency issue. In case we have to update the instructors table and course taught by two teachers, we have to update everywhere all other entries for that particular teacher/prof.

Option B:

Advantage – Since all the data is in single table there is no need to join. It will take less time for query processing and no need to write complex queries

Disadvantage – As data stored in same table, can cause data duplication. Also can use more store for such scenario

Option C:

Advantage – It would be a better choice as we can segregate operations based on need. No join will be required as the queries can be performed independently

Disadvantage – The same data will store in two tables, will need more space for storage.

Question 2. Which option would you choose and why?

Ans :-- I would choose option C, Because It can help to perform query for each requirement.

Question 3. Would your answer to Question 2 be different if the majority of classes had multiple instructors? How about if only one or two classes had multiple instructors? (Explain your answer.)

Ans :-- Option C is good enough to handle both the cases.

Question 4. [OPTIONAL] Can you think of another reasonable alternative design besides Options A, B, and C? If

so, what are the advantages and disadvantages of your alternative design?

## Scenario II

In this scenario, we are building a data warehouse for an online brokerage company. The company makes money by charging commissions when customers buy and sell stocks. We are planning to have a Trades fact table with the grain of one row per stock trade. We will use the following dimensions: Date, Customer, Account, Security (i.e. which stock was traded), and TradeType.

The company's data analysts have told us that they have developed two customer scoring techniques that are used extensively in their analyses.

- Each customer is placed into one of nine Customer Activity Segments based on their frequency of transactions, average transaction size, and recency of transactions.
- Each customer is assigned a Customer Profitability Score based on the profits earned as a result of that customer's trades. The score can be either 1,2,3,4, or 5, with 5 being the most profitable.

These two scores are frequently used as filters or grouping attributes in queries.

For example:

- How many trades were placed in July by customers in each customer activity segment?
- What was the total commission earned in each quarter of 2003 on trades of IBM stock by customers with a profitability score of 4 or 5?

There are a total of 100,000 customers, and scores are recalculated every three months. The activity level or profitability level of some customers changes over time, and users are very interested in understanding how and why this occurs.

We are considering several options for dealing with the customer scores:

**Option A**

**Option B**

**Option C**

**Option D**

The scores are attributes of the Customer dimension, When scores change, the old score is overwritten with the new score (Type 1 Slowly Changing Dimension).

The scores are attributes of the Customer dimension. When scores change, new Customer dimension rows are created using the updated scores (Type 2 Slowly Changing Dimension).

The scores are stored in a separate CustomerScores dimension which contains 45 rows, one for each combination of activity and profitability scores. The Trades fact table includes a foreign key to the CustomerScores dimension.

The scores are stored in a CustomerScores outrigger table which contains 45 rows. The Customer dimension includes a foreign key to the outrigger table (but the fact table does not). When scores change, the foreign key column in the Customer table is updated to point to the correct outrigger row.

Please answer the following questions.

Question 5. What are the strengths and weaknesses of each option?

Ans –

Option A – This will help to save storage as only one table is being overwritten but the disadvantage will be the missing of old score as the new score will be takes place of old score

Option B – This will be better option than option A, as historical data will maintain in dim table unlike Option A. But we need to add few extra attributes like effective date to , effective date from etc to make data more meaningful. Which will be a cause to increase in data size

Option C – looks like this is the best option among all. As we have only two scoring factor and combination of them will become 45 rows. So can create a junk dim table to make is simple and reduce the complexity and storage. But if we need to add more scoring factor like previous two then we need to modify/update our junk table

Option D – This will points towards multiple join as many tables came in picture. But the advantage will be the data is segregated properly

Question 6. Which option would you choose and why?

Ans – I could choose the option C. As it will not take much tables for join and will not be required to write complex quires.

Question 7. Would your answer to Question 6 be different if the number of customers and/or the time interval between score recalculations was much larger or much smaller? (Explain your answer.)

Ans – No, if number of customer or time interval between score recalculation will change, then it will not affect the junk table of those 45 entries.

Question 8. [OPTIONAL] Can you think of another reasonable alternative design besides Options A, B, C, and D? If so, what are the advantages and disadvantages of your alternative design?