

# CHAPTER 11



## Data Analytics

### Practice Exercises

- 11.1** Describe benefits and drawbacks of a source-driven architecture for gathering of data at a data warehouse, as compared to a destination-driven architecture.

**Answer:**

In a destination-driven architecture for gathering data, data transfers from the data sources to the data warehouse are based on demand from the warehouse, whereas in a source-driven architecture, the transfers are initiated by each source.

The benefits of a source-driven architecture are

- Data can be propagated to the destination as soon as they become available. For a destination-driven architecture to collect data as soon as they are available, the warehouse would have to probe the sources frequently, leading to a high overhead.
- The source does not have to keep historical information. As soon as data are updated, the source can send an update message to the destination and forget the history of the updates. In contrast, in a destination-driven architecture, each source has to maintain a history of data which have not yet been collected by the data warehouse. Thus storage requirements at the source are lower for a source-driven architecture.

On the other hand, a destination-driven architecture has the following advantages.

- In a source-driven architecture, the source has to be active and must handle error conditions such as not being able to contact the warehouse for some time. It is easier to implement passive sources, and a single active warehouse. In a destination-driven architecture, each source is required to provide only a basic functionality of executing queries.

- The warehouse has more control on when to carry out data gathering activities and when to process user queries; it is not a good idea to perform both simultaneously, since they may conflict on locks.

**11.2** Draw a diagram that shows how the *classroom* relation of our university example as shown in Appendix A would be stored under a column-oriented storage structure.

**Answer:**

The relation would be stored in three files, one per attribute, as shown below. We assume that the row number can be inferred implicitly from position, by using fixed-size space for each attribute. Otherwise, the row number would also have to be stored explicitly.

<i>building</i>
Packard
Painter
Taylor
Watson
Watson

<i>room_number</i>
101
514
3128
100
120

<i>capacity</i>
500
10
70
30
50

**11.3** Consider the *takes* relation. Write an SQL query that computes a cross-tab that has a column for each of the years 2017 and 2018, and a column for **all**, and one row for each course, as well as a row for **all**. Each cell in the table should contain the number of students who took the corresponding course in the corresponding year, with column **all** containing the aggregate across all years, and row **all** containing the aggregate across all courses.

**Answer:**

- 11.4** Consider the data warehouse schema depicted in Figure 11.2. Give an SQL query to summarize sales numbers and price by store and date, along with the hierarchies on store and date.

**Answer:**

query:

```
select store-id, city, state, country,
       date, month, quarter, year,
       sum(number), sum(price)
from sales, store, date
where sales.store-id = store.store-id and
      sales.date = date.date
groupby rollup(country, state, city, store-id),
         rollup(year, quarter, month, date)
```

- 11.5** Classification can be done using *classification rules*, which have a *condition*, a *class*, and a *confidence*; the confidence is the percentage of the inputs satisfying the condition that fall in the specified class.

For example, a classification rule for credit ratings may have a condition that salary is between \$30,000 and \$50,000, and education level is graduate, with the credit rating class of *good*, and a confidence of 80%. A second rule may have a condition that salary is between \$30,000 and \$50,000, and education level is high-school, with the credit rating class of *satisfactory*, and a confidence of 80%. A third rule may have a condition that salary is above \$50,001, with the credit rating class of *excellent*, and confidence of 90%. Show a decision tree classifier corresponding to the above rules.

Show how the decision tree classifier can be extended to record the confidence values.

**Answer:**

FILL IN

- 11.6** Consider a classification problem where the classifier predicts whether a person has a particular disease. Suppose that 95% of the people tested do not suffer from the disease. Let *pos* denote the fraction of *true positives*, which is 5% of the test cases, and let *neg* denote the fraction of *true negatives*, which is 95% of the test cases. Consider the following classifiers:

- Classifier  $C_1$ , which always predicts negative (a rather useless classifier, of course).
- Classifier  $C_2$ , which predicts positive in 80% of the cases where the person actually has the disease but also predicts positive in 5% of the cases where the person does not have the disease.

- Classifier  $C_3$ , which predicts positive in 95% of the cases where the person actually has the disease but also predicts positive in 20% of the cases where the person does not have the disease.

For each classifier, let  $t_{pos}$  denote the *true positive* fraction, that is the fraction of cases where the classifier prediction was positive, and the person actually had the disease. Let  $f_{pos}$  denote the *false positive* fraction, that is the fraction of cases where the prediction was positive, but the person did not have the disease. Let  $t_{neg}$  denote *true negative* and  $f_{neg}$  denote *false negative* fractions, which are defined similarly, but for the cases where the classifier prediction was negative.

- Compute the following metrics for each classifier:
  - Accuracy*, defined as  $(t_{pos} + t_{neg})/(pos + neg)$ , that is, the fraction of the time when the classifier gives the correct classification.
  - Recall* (also known as *sensitivity*) defined as  $t_{pos}/pos$ , that is, how many of the actual positive cases are classified as positive.
  - Precision*, defined as  $t_{pos}/(t_{pos} + f_{pos})$ , that is, how often the positive prediction is correct.
  - Specificity*, defined as  $t_{neg}/neg$ .
- If you intend to use the results of classification to perform further screening for the disease, how would you choose between the classifiers?
- On the other hand, if you intend to use the result of classification to start medication, where the medication could have harmful effects if given to someone who does not have the disease, how would you choose between the classifiers?

**Answer:**

FILL