

CS352 Lecture - Decision Support Systems, Data Mining

11/27/06

Objectives:

1. To introduce basic OLAP concepts (cubes, rollups, ranking)
2. To introduce the notion of a data warehouse
3. To introduce the notion of a decision tree

Materials

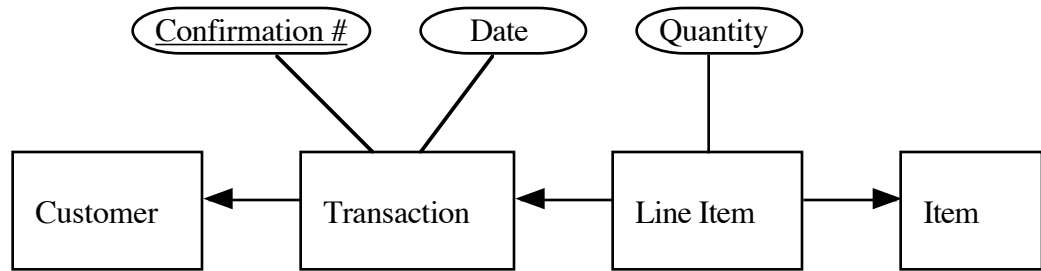
1. Projectable of Figure 18.1 from book
2. Projectable of Figure 18.3 from book
3. Projectable of Figure 18.3 showing the various summaries
4. Eight SQL queries equivalent to one group by cube query
5. Projectable of equivalent single query using cube
6. Online access to sample database; electronic form of queries for cut and paste
7. Projectable of Figure 18.8 from book
8. "Jets and Sharks" database (called westside)

I. Introduction

A. One of the largest categories of use for a database system is the storing of transactional data.

1. The term "transaction" here is used in a different sense from the way we used it earlier, though the concepts are related.
2. In this context, a transaction refers to a single interaction between a "customer" and an organization - such as:
 - a) A single deposit, withdrawal, transfer transaction at a bank.
 - b) A single purchase at a store, or from an ecommerce siteetc.
3. For large firms, thousands of transactions may take place in a single hour.
4. Each transaction results in one or more rows being added to tables in the database

Example: An ecommerce site might use a structure like this:



A new row is added to the transaction table for each purchase. Associated with this row are one or more rows in the line item table, each describing “line item” of the purchase - i.e. the purchase of some quantity of some particular item.

5. The tables used to store transactional data can be huge of course.

B. While transactional data needs to be stored for its own sake (e.g. to support shipping of an order or listing of a bank transaction on a monthly statement), it also has relevance to organizational decision-making. For example, companies often use data from sales transactions to help them adjust their product offerings to be a better match to what customers are buying.

Using transactional data in this way is called OLAP - online analytical processing.

C. Another related concept is the notion of a Data Warehouse.

1. For a large organization, the transactional data may be stored in multiple databases (e.g. at different physical branches, or in different departments)
 - a) These databases may have different schemes
 - b) No one person may have the ability to access all of these databases.
 - c) The databases may store only relatively current data
2. To support corporate decision-making, an organization may create a data warehouse - a separate database which stores data collected from a variety of sources - specifically to support decision-making activities.
 - a) Data is propagated from the various transactional databases (called data sources) to the warehouse on a regular basis (perhaps daily). (Of course, this means that the data in the warehouse is always

slightly out of date - but this is not an issue when the data is used for decision making that typically looks at historical trends over a period of months or years)

- b) The data may be transformed in various ways to accommodate differences in schemes in the various sources and to eliminate duplicate information, etc.
- c) The data may be converted to a summarized form to keep down the aggregate size.
- d) Historical data is maintained to facilitate looking at long-range trends.

3. We will not discuss this topic further.

D. OLAP assumes that the user knows what he/she wants and requests it explicitly (“give me a cross tab showing item type versus item color ...”).

- 1. Often times, though, there are important relationships present in the data which can be discovered by a process known as data mining.

Example: one often hears about medical “risk factors” for various conditions - e.g. belonging to a certain ethnic group, taking a certain medicine, or smoking, or ... Note that risk factors may include both obvious and non-obvious attributes [e.g. smoking, ethnicity].

Example: If advertisers know that a particular product is especially attractive to a particular group of customers, they may target their advertising toward that demographic in terms of choosing which programs to air their commercials on.

Example: It may turn out that people who purchase some particular product are more likely to purchase another particular product. In that case, a sales person or e-commerce system may suggest the second product to the purchaser of the first product, or a coupon printer at store checkout may print a coupon for the second product when a person purchases the first product.

- 2. Data mining can easily be the subject of graduate-level courses. We will just look at it briefly here.

II. Basic OLAP Concepts

A. Various sorts of summary information can be useful in corporate decision-making.

1. Simple summary information is often useful

Example: Total number sold of a particular item can help a firm decide whether to continue to carry it.

Such data can be obtained by a SQL query like the following:

```
select item_code, sum(quantity)
      from line_item
      group by item_code
```

2. But often times a further breakdown is useful

Example: Consider clothing items, which come in different colors. It may be that some color is selling poorly, and should be discontinued, even though other color(s) are doing well. (Suppose that the same `item_code` is used for all the colors of a given item, with `color` being an attribute of `line_item`). This would call for a query like:

```
select item_code, color, sum(quantity)
      from line_item
      group by item_code, color
```

3. Suppose the firm sells many different items in a given color. Then it might be interesting to know not only how popular a particular item is in each color, but also how popular a color is for all items. (This might lead to a decision to begin producing some item in a color that is not currently used for that item, but is popular for similar items)

The needed summary information could be obtained by using a query like the following :

```
select color, sum(quantity)
      from line_item
      group by color
```

4. Sometimes items fall into broader categories - e.g. different kinds of pants or shirts or whatever. It may be beneficial to look at summary data for broad categories, rather than detailed information for a particular code - e.g. suppose the item table stores various pieces of information for each `item_code`, including an `item_type`. Then a query like the following may be useful:

```
select item_type, color, sum(quantity)
      from line_item natural join item
      group by item_type, color
```

5. All of these kinds of information might be obtained from a cross-tab like the following:

PROJECT Figure 18.1 from book

Such a table cannot be directly produced by a SQL query, of course. However, we will meet some SQL extensions that support extracting the necessary data which can then be converted to a suitable display by additional software

6. The construction of tabulations like the above (sometimes with more than two dimensions) is referred to as Online Analytical Processing (OLAP). Since generating a cross-tab from a huge transactional database “from scratch” can involve massive computation, systems that facilitate OLAP often pre-compute certain summary data to facilitate quick responses to interactive user queries (hence the term “online”).

We will not talk about OLAP software per se, but we will talk about some of the SQL facilities that support OLAP-style operations.

B. Attributes in a transactional scheme fall into two broad categories.

1. Dimensional attributes- e.g. [for clothing] item type, color, size.

Actually, a dimensional attribute can be an attribute that is explicitly stored, or one that is computed from an explicitly stored value.

Example: We may actually store date of birth, but use age (calculated using today’s date) as a dimensional attribute.

Example: sometimes it is helpful to use ranges. For example, in an earlier homework, you worked with the employee table in the sample database, which has a “years of education column”. You were asked to write a SQL case statement that converted this value into one of three strings (GRADUATE, COLLEGE, HIGH SCHOOL). If one were analyzing some piece of information (e.g. average salary) across this dimension, using just these three categories rather than a distinct category for each different value for years of education would likely be helpful.

2. Measure attributes - a value that is summarized using some sort of aggregate operation (e.g. sum, average, maximum). It may be an attribute in the table or a value that can be computed from attributes - e.g.. [for clothing] quantity sold, monetary value (= quantity * price), or something like count() - a simple count of the number of occurrences.

C. The term “cube” refers to a structure that represents an aggregation of the value of some measure attribute for each possible combination of values of some dimensional attributes, with the number of dimensions of the cube equal to the number of dimensional attributes used. For each dimension, there is usually also a summary value (“all”) which represents the aggregation of the values for all the possible values of the dimensional attribute.

Example: PROJECT Figure 18.3

1. Of course, for a cube of three or more dimensions, it is not possible to meaningfully display all the data at one time. Instead, the user can ask to see a particular two-dimensional slice of the cube, formed by using specific values for all but two of the attributes.

Example: for the cube just projected, the user might ask to see a slice showing all the different values for color and item_name for a particular size (e.g. medium).

PROJECT

2. Note that a cube is quite different from a relational table. In a table, the columns are specified when the table is created. In a cube, the “columns” are determined by the number of different values that occur in the database for a particular attribute.

Example: if the database included data for striped clothing as well as dark, pastel, and white, then the cube would have one more vertical level.

D. It turns out that cubes cannot be constructed using features of SQL we have covered thus far.

1. We could, however, create some of the data by a query like the following - assuming a table that has columns item_name, color, size (dimension attributes) and number (measure attribute):

```
select item_name, color, size, sum(number)
      from sales
      group by item_name, color, size
```

a) What data needed for the cube would be missing?

ASK

The “all” faces

PROJECT result of above query, then the various summary portions that are missing.

b) Of course, we could use additional SQL queries to generate the missing summaries - e.g.item name versus color for all sizes could be produced by:

```
select item_name, color, sum(number)
  from sales
  group by item_name, color;
```

c) How many queries in all would be needed to get all the different ways of slicing the data (including the original query)?

ASK

8

PROJECT queries

```
select item_name, color, size, sum(number)
  from sales
  group by item_name, color, size;
```

```
select item_name, color, sum(number)
  from sales
  group by item_name, color;
```

```
select item_name, size, sum(number)
  from sales
  group by item_name, size;
```

```
select color, size, sum(number)
  from sales
  group by color, size;
```

```
select item_name, sum(number)
  from sales
  group by item_name;
```

```

select color, sum(number)
  from sales
  group by color;

select size, sum(number)
  from sales
  group by size;

select sum(number)
  from sales;

```

2. However, SQL:1999 added a cube operator which makes it easy to create such structures.

The book gives a SQL query that would construct the needed data (which would then have to be displayed appropriately):

```

select item_name, color, size, sum(number)
  from sales
  group by cube(item_name, color, size)

```

This produces the same result as we would get from the above queries

3. An example of using SQL cube.

- a) Recall the employee table in the sample database we have used on various homework problems.

```

DEMO select * from employee

```

- b) Now suppose we want to explore how average salary (a measure attribute) varies with various dimensional attributes (such as department, length of service, job title, education level, gender, age).

We could, for example, use a query like:

```

select sex, avg(salary)
  from employee
  group by sex;

```

To see how average salary varies by gender.

DEMO

- c) We could see how average salary varies over several different attributes (job, education level, gender) by using a cube query like:

```

select job, edlevel, sex, avg(salary)
  from employee
  group by cube(job, edlevel, sex)
  order by job, edlevel, sex;

```

DEMO

(BTW: Note the use of order by to control the final order of printing. DEMO w/o this).

- d) Actually, it might be better to look at the edlevel dimension by broad groups (as in homework 2) rather than by individual values. To do this, we could use something like:

```
select job, education, sex, avg(salary)
  from (select job, case
          when edlevel >= 18 then 'GRADUATE'
          when edlevel >= 16 then 'COLLEGE'
          else 'HIGH SCHOOL'
        end as education, sex, salary
        from employee) as e
  group by cube(job, education, sex)
  order by job, education, sex;
```

(BTW: Is there something that becomes noticeable when the data is presented this way that was not obvious before?)

ASK

- e) Actually, another issue may be involved here that one can see when looking at experience as well as formal education

```
select education, experience, sex, avg(salary)
  from (select case
          when edlevel >= 18 then 'GRADUATE'
          when edlevel >= 16 then 'COLLEGE'
          else 'HIGH SCHOOL'
        end as education,
        case
          when hiredate < '01-01-1960' then 'OVER 20'
          when hiredate < '01-01-1970' then '10-20'
          when hiredate < '01-01-1975' then '5-10'
          else 'NEW'
        end as experience,
        sex,
        salary
        from employee) as e
  group by cube(education, experience, sex)
  order by education, experience, sex;
```

E. SQL also supports another similar OLAP operation called `rollup`.

1. The difference is basically this: with `cube`, if you specify n dimensions, you form 2^n groups, representing all the possible combinations of the various dimensions and “all”. With `rollup`, you get $n+1$ groups - all the dimensions, all the dimensions except the last, all the dimensions except the last and second to the last ...
2. We can illustrate this using a query similar to the one we just did with `cube`.

```
select job, education, sex, avg(salary)
  from (select job, case
          when edlevel >= 18 then 'GRADUATE'
          when edlevel >= 16 then 'COLLEGE'
          else 'HIGH SCHOOL'
        end as education, sex, salary
        from employee) as e
  group by rollup(job, education, sex)
  order by job, education, sex;
```

DEMO

In this case, the groups we get are based on all three dimensions; just `workdept` and `job`; just `workdept`, and an overall summary

F. SQL also supports OLAP operations related to ranking

1. For example, suppose we were interested in looking at the relationship between salary and education level. We might formulate a query like the following:

```
select firstname, lastname, salary,
       rank() over (order by edlevel desc) as edrank
  from employee
  order by salary desc;
```

DEMO (Note that “1” means the greatest amount of formal education, etc.)

2. Do you notice anything funny about the values reported for `EDRANK`?

ASK

Several values appear multiple times (e.g. 3, 12) and other values don't appear at all (e.g. 4, 13 ...)

This is because, when there are ties, the rank() function gives all the tied values the same rank, and then skips over as many rank values as there are duplicates (e.g. since four people are tied for third, the next person is given rank 7)

3. An alternate function called dense_rank does not skip rank numbers

```
select firstnme, lastname, salary,  
       dense_rank() over (order by edlevel desc) as edrank  
from employee  
order by salary desc;
```

DEMO

4. There are other facilities available as well which the book discusses - e.g. GROUPING and PARTITION by.

III. Data Mining

A. Data mining refers to systematically analyzing large collections of data in order to find useful patterns. It differs from OLAP in that OLAP requires a human user to explicitly specify the groupings to be considered, while data mining looks for patterns a human user might not anticipate.

B. In data mining, we are looking for rules of the general form

if these conditions are present in the data, then this conclusion is likely

1. Such rules are seldom 100% accurate, of course.
2. But they are valuable to the extent that they are generally true.

C. One important use of data mining is for prediction. For example, a bank might want to predict whether or not an individual is a good credit risk. Data mining techniques approach a problem like this as follows:

1. We begin with a collection of historical data which includes the value of various variables that were known at the time an individual applied for credit plus an indication of how the applicant actually performed as a credit recipient (did the applicant default on a loan, was the applicant habitually late with payments ...)

This set of historical data is known as the training instances.

2. The training instances are mined to find correlations between variables known at application time and the actual credit-worthiness of an applicant. The goal of the mining process is to discover rules which, given variables known at application time, will successfully classify an individual as an excellent, good, average, or poor credit risk. - e.g

if this pattern is present in the input data, then it is likely that this individual will be a _____ credit risk

A desirable rule is one which, when applied to the various training instances, with high reliability is able to predict what actually transpired.

3. These rules are then used to “score” future applications - assuming that the same patterns that were present in the historical data will also be true in the future.
4. A technique that is often used for discovering such rules is based on the notion of a decision tree.

PROJECT Figure 18.8

D. Another use of data mining is to discover associations.

1. As an example, book sellers like Amazon look for associations between books people bought - i.e. they are seeking to discover rules of the form “a person who bought X is also likely to buy Y”. How do they use this kind of information?

ASK

- a) Most of their pages describing products include a section near the bottom that reads “Customers who bought this item also bought” (The hope is that, if you are interested in some item - even if you don’t buy it, you may be interested in these other items as well, and may buy one of them).
 - b) When a regular customer connects to Amazon, they furnish a section labelled “Today’s recommendations for you” - which lists titles Amazon thinks a customer might buy based on previous purchases that customer has made.
2. Companies like grocery stores, drug stores, etc. look for associations between products a customer purchases. This may affect placement of the products in the store, as well as generation of checkout coupons.

3. While associations within a single visit are important, companies are also interested in associations over time - e.g. how does what associations are there between what the customer buys on one visit to the store and what the customer buys on an other visit.

a) One way of seeing how important this is to various businesses is the proliferation of various “frequent shopper programs”. Many stores will let you register and offer significant discounts if you consistently use your frequent shopper tag when you make purchases.

b) For the most part, such programs are not designed to generate mailing lists for email or snail mail. Rather, they are used to facilitate data mining, by allowing the company to build a record of purchases a given customer makes over time. (Absent such a program, the store would have records of myriads of transactions, but no way to associate a given transaction with previous transactions serving the same customer.)

4. Associations may be expressed as rules of the form

some pattern => some other pattern

(e.g. for a grocery store: includes(T, bread) => includes(T, milk)) - i.e. “if some transaction T includes bread, then it also includes milk”

a) The left hand side of the rule is called the antecedent (e.g., in the above, includes(T, bread))

b) The right hand side of the rule is called the consequent (e.g. in the above includes(T, milk))

c) For each possible rule, two measures can be calculated. For any given problem, we are only interested in rules for which both measures are sufficiently high (where “sufficiently high” depends on the application)

(1) The support for a rule is the percentage of all cases in which both the antecedent and the consequent are present. Rules with low support are uninteresting for two reasons:

(a) The evidence for them is weak, and may be statistically insignificant

Example: the support for the rule

name(s, bjork) => teaches(s, cs) is less than 0.01 (less than 1% of Gordon employees have name Bjork and teach CS)
Even though the confidence is high (currently 50% of the Gordon employees whose name is Bjork teach CS!), the rule is uninteresting because it is not statistically significant

(b) The number of cases the rule applies to is so small that it doesn't warrant attention on the company's part

(2) The confidence for a rule is the percentage of cases where the antecedent is true where the consequent is also true. Only rules with high confidence are of significant predictive value

Example: the confidence for the rule

major(S, cs) => gender(S, male) is - unfortunately - quite high.
(As of the time I prepared this lecture, 18 of the 20 CS majors at Gordon were males - so the confidence would be 0.9)

5. An example: we will use a database describing hypothetical members of the Jets and Sharks gangs of Westside story fame.

a) Connect to westside database

b) `select * from gangsters;`

c) Consider the association between age and gang membership. It turns out that - in this particular set of data - the Jets tend to be younger than the Sharks. Consider the following pair of rules:

`age(X, in20s) => gang(X, jets)`
`gang(X, jets) => age(X, in20s)`

To calculate support and confidence for these rules, we could use a query like the following:

```
select gang, age, count(*)
  from gangsters
  group by cube(gang, age)
  order by gang, age;
```

DEMO

- d) It turns out that 9 of the 27 gangsters are in their 20's and Jets - so the support for both rules is 0.33
- e) Of the 10 gangsters in their 20's, 9 are jets, so the confidence of the first rule is $9/10$ or 0.9. (So if you were a police officer and you saw a young gangster, you'd be on pretty safe ground to assume he was a Jet!)
- f) Of the 15 Jets, 9 are in their 20's, so the confidence of the second rule is $9/15$ or 0.6.

E. Data mining turns out to be a place where there is considerable overlap between the areas of artificial intelligence and database systems - in fact, you see this as a topic in both kinds of textbooks, and the database we just used is also used a standard example in teaching about neural networks in AI!