# Data Analysis and Mining

## Solutions to Practice Exercises

**18.1** query:

$$\textbf{groupby rollup}(a), \textbf{rollup}(b), \textbf{rollup}(c), \textbf{rollup}(d)$$

**18.2** We assume that multiple students do not have the same marks since otherwise the question is not deterministic; the query below deterministically returns all students with the same marks as the $n$ student, so it may return more than $n$ students.

> **select** *student*, **sum**(*marks*) **as** *total*,
>         **rank**() **over** (**order by** (*total*) **desc**) **as** *trank*
> **from** *S*
> **groupby** *student*
> **having** *trank* $\leq n$

**18.3** query:

> **select** *t1.account-number*, *t1.date-time*, *sum(t2.value)*
> **from** *transaction* as *t1*, *transaction* as *t2*
> **where** *t1.account-number* $=$ *t2.account-number* **and**
>         *t2.date-time* $<$ *t1.date-time*
> **groupby** *t1.account-number*, *t1.date-time*
> **order by** *t1.account-number*, *t1.date-time*

**18.4** query:

```
(select color, size, sum(number)
 from sales
 groupby color, size
)
union
(select color, 'all', sum(number)
 from sales
 groupby color
)
union
(select 'all', size, sum(number)
 from sales
 groupby size
)
union
(select 'all', 'all', sum(number)
 from sales
)
```

**18.5** 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 it becomes available. For a destination-driven architecture to collect data as soon as it is 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 is 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 architecure, 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.

**18.6** Consider the following pair of rules and their confidence levels :

| No. | Rule | Conf. |
|-----|------|-------|
| 1. | $\forall\ persons\quad P,\ 10000\ <\ P.salary\ \leq\ 20000\ \Rightarrow$ $P.credit = good$ | 60% |
| 2. | $\forall\ persons\quad P,\ 20000\ <\ P.salary\ \leq\ 30000\ \Rightarrow$ $P.credit = good$ | 90% |

The new rule has to be assigned a confidence-level which is between the confidence-levels for rules $1$ and $2$. Replacing the original rules by the new rule will result in a loss of confidence-level information for classifying persons, since we cannot distinguish the confidence levels of perople earning between 10000 and 20000 from those of people earning between 20000 and 30000. Therefore we can combine the two rules without loss of information only if their confidences are the same.

**18.7** 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*)