

# Data Mining Primitives & DMQL

**Dr.M. Venkatesan**

Associate Professor

School of Computing Science and

Engineering

VIT University

# Outline

- Data mining primitives: What defines a data mining task?
- A data mining query language
- Design graphical user interfaces based on a data mining query language
- Architecture of data mining systems

# Why Data Mining Primitives and Languages?

- Finding all the patterns autonomously in a database? — unrealistic because the patterns could be too many but uninteresting
- Data mining should be an interactive process
  - User directs what to be mined
- Users must be provided with a set of **primitives** to be used to communicate with the data mining system
- Incorporating these primitives in a **data mining query language**
  - More flexible user interaction
  - Foundation for design of graphical user interface
  - Standardization of data mining industry and practice

# What Defines a Data Mining Task ?

- Task-relevant data
- Type of knowledge to be mined
- Background knowledge
- Pattern interestingness measurements
- Visualization of discovered patterns

# Task-Relevant Data (Minable View)

- Database or data warehouse name
- Database tables or data warehouse cubes
- Condition for data selection
- Relevant attributes or dimensions
- Data grouping criteria

# Types of knowledge to be mined

- Characterization
- Discrimination
- Association
- Classification/prediction
- Clustering
- Outlier analysis
- Other data mining tasks

# Background Knowledge: Concept Hierarchies

- Schema hierarchy
  - E.g., street < city < province\_or\_state < country
- Set-grouping hierarchy
  - E.g., {20-39} = young, {40-59} = middle\_aged
- Operation-derived hierarchy
  - email address: [Robert.Gyorodi@ul.ie](mailto:Robert.Gyorodi@ul.ie)  
login-name < department < university < country
- Rule-based hierarchy
  - low\_profit\_margin (X) <= price(X, P<sub>1</sub>) and cost (X, P<sub>2</sub>) and (P<sub>1</sub> - P<sub>2</sub>) < \$50

# Measurements of Pattern Interestingness

- **Simplicity**  
e.g., (association) rule length, (decision) tree size
- **Certainty**  
e.g., confidence,  $P(A|B) = \#(A \text{ and } B) / \#(B)$ , classification reliability or accuracy, certainty factor, rule strength, rule quality, discriminating weight, etc.
- **Utility**  
potential usefulness, e.g., support (association), noise threshold (description)
- **Novelty**  
not previously known, surprising (used to remove redundant rules, e.g., Canada vs. Vancouver rule implication support ratio)



# Visualization of Discovered Patterns

- Different backgrounds/usages may require **different forms of representation**
  - E.g., rules, tables, crosstabs, pie/bar chart etc.
- **Concept hierarchy** is also important
  - Discovered knowledge might be more understandable when represented at **high level of abstraction**
  - Interactive **drill up/down, pivoting, slicing and dicing** provide different perspectives to data
- Different kinds of **knowledge** require different representation: association, classification, clustering, etc.

# A Data Mining Query Language (DMQL)

- Motivation
  - A DMQL can provide the ability to support ad-hoc and interactive data mining
  - By providing a standardized language like SQL
    - Hope to achieve a similar effect like that SQL has on relational database
    - Foundation for system development and evolution
    - Facilitate information exchange, technology transfer, commercialization and wide acceptance
- Design
  - DMQL is designed with the primitives described earlier

# Syntax for DML

- Syntax for specification of
  - task-relevant data
  - the kind of knowledge to be mined
  - concept hierarchy specification
  - interestingness measure
  - pattern presentation and visualization
- Putting it all together—a DML query

# Syntax:

## Specification of Task-Relevant Data

- *use database* database\_name, or *use data warehouse* data\_warehouse\_name
- *from relation(s)/cube(s)* [*where* condition]
- *in relevance to* att\_or\_dim\_list
- *order by* order\_list
- *group by* grouping\_list
- *having* condition

# Specification of task-relevant data

**Example 4.11** This example shows how to use DMQL to specify the task-relevant data described in Example 4.1 for the mining of associations between items frequently purchased at *AllElectronics* by Canadian customers, with respect to customer *income* and *age*. In addition, the user specifies that she would like the data to be grouped by date. The data are retrieved from a relational database.

```
use database AllElectronics_db
in relevance to I.name, I.price, C.income, C.age
from customer C, item I, purchases P, items_sold S
where I.item_ID = S.item_ID and S.trans_ID = P.trans_ID and P.cust_ID = C.cust_ID
      and C.address = "Canada"
group by P.date
```



# Syntax: Kind of knowledge to Be Mined

- Characterization

Mine\_Knowledge\_Specification ::=  
    *mine characteristics*[*as pattern\_name*]  
    *analyze*measure(s)

- Discrimination

Mine\_Knowledge\_Specification ::=  
    *mine comparison*[*as pattern\_name*]  
    *for*target\_class *where*target\_condition  
    { *versus*contrast\_class\_1 *where*contrast\_condition\_1 }  
    *analyze*measure(s)

E.g. *mine comparison as purchaseGroups*

*for bigSpenders where avg(l.price) >= \$100*

*versus budgetSpenders where avg(l.price) < \$100*

*analyze count*

# Syntax: Kind of knowledge to Be Mined

- Association

Mine\_Knowledge\_Specification ::=  
*mine associations*[*as* pattern\_name]  
[*matching* <metapattern>]

E.g. mine associations as buyingHabits

matching  $P(X:\text{custom}, W) \wedge Q(X, Y) \Rightarrow \text{buys}(X, Z)$

- Classification

Mine\_Knowledge\_Specification ::=  
*mine classification*[*as* pattern\_name]  
*analyze* classifying\_attribute\_or\_dimension

- Other Patterns

clustering, outlier analysis, prediction ...

# Syntax: Concept Hierarchy Specification

- To specify what concept hierarchies to use  
use hierarchy **<hierarchy>** for **<attribute\_or\_dimension>**
- We use different syntax to define different type of hierarchies
  - schema hierarchies  
define hierarchy **time\_hierarchy** on **date** as **[date,month  
quarter,year]**
  - set-grouping hierarchies  
define hierarchy **age\_hierarchy** for **age** on **customer** as  
**level1: {*young, middle\_aged, senior*} < level0: all**  
**level2: {20, ..., 39} < level1: *young***  
**level2: {40, ..., 59} < level1: *middle\_aged***  
**level2: {60, ..., 89} < level1: *senior***



# Syntax: Concept Hierarchy Specification

- operation-derived hierarchies

define hierarchy **age\_hierarchy** for **age** on **customer** as  
**{age\_category(1), ..., age\_category(5)} := cluster(default, age, 5) < all(age)**

- rule-based hierarchies

define hierarchy **profit\_margin\_hierarchy** on **item** as  
**level\_1: low\_profit\_margin < level\_0: all**  
**if (price - cost) < \$50**  
**level\_1: medium-profit\_margin < level\_0: all**  
**if ((price - cost) > \$50) and ((price - cost) <= \$250))**  
**level\_1: high\_profit\_margin < level\_0: all**  
**if (price - cost) > \$250**

# Specification of Interestingness Measures

- Interestingness measures and thresholds can be specified by a user with the statement:  
with **<interest\_measure\_name>** threshold =  
    **threshold\_value**
- **Example:**  
    **with support threshold = 0.05**  
    **with confidence threshold = 0.7**

# Specification of Pattern Presentation

- Specify the display of discovered patterns

display as **<result\_form>**

- To facilitate interactive viewing at different concept level, the following syntax is defined:

$$\begin{aligned} \text{Multilevel\_Manipulation} ::= & \textit{roll up on attribute\_or\_dimension} \\ & | \textit{drill down on attribute\_or\_dimension} \\ & | \textit{add attribute\_or\_dimension} \\ & | \textit{drop attribute\_or\_dimension} \end{aligned}$$

# Putting it all together: A DMQL query

use database **AllElectronics\_db**  
use hierarchy **location\_hierarchy** for **B.address**  
mine characteristics as **customerPurchasing**  
analyze **count%**  
in relevance to **C.age, I.type, I.place\_made**  
from **customer C, item I, purchases P, items\_sold S, works\_at W,**  
**branch**  
where **I.item\_ID = S.item\_ID and S.trans\_ID = P.trans\_ID**  
**and P.cust\_ID = C.cust\_ID and P.method\_paid = ``AmEx"**  
**and P.empl\_ID = W.empl\_ID and W.branch\_ID =**  
**B.branch\_ID and B.address = ``Canada" and I.price >=**  
**100**  
with **noise threshold = 0.05**  
display as **table**

# Other Data Mining Languages & Standardization Efforts

- Association rule language specifications
  - MSQL (Imielinski & Virmani'99)
  - MineRule (Meo Psaila and Ceri'96)
  - Query flocks based on Datalog syntax (Tsur et al'98)
- OLEDB for DM (Microsoft'2000)
  - Based on OLE, OLE DB, OLE DB for OLAP
  - Integrating DBMS, data warehouse and data mining
- CRISP-DM (CRoss-Industry Standard Process for Data Mining)
  - Providing a platform and process structure for effective data mining
  - Emphasizing on deploying data mining technology to solve business problems

# Designing Graphical User Interfaces Based on a Data Mining Query Language

- What tasks should be considered in the design GUIs based on a data mining query language?
  - Data collection and data mining query composition
  - Presentation of discovered patterns
  - Hierarchy specification and manipulation
  - Manipulation of data mining primitives
  - Interactive multilevel mining
  - Other miscellaneous information

# Data Mining System Architectures

- Coupling data mining system with DB/DW system
  - No coupling—flat file processing, not recommended
  - Loose coupling
    - Fetching data from DB/DW
  - Semi-tight coupling—enhanced DM performance
    - Provide efficient implement a few data mining primitives in a DB/DW system, e.g., sorting, indexing, aggregation, histogram analysis, multiway join, precomputation of some stat functions
  - Tight coupling—A uniform information processing environment
    - DM is smoothly integrated into a DB/DW system, mining query is optimized based on mining query, indexing, query processing methods, etc.

# Today's lecture ...

- Data mining primitives: What defines a data mining task?
- A data mining query language
- Design graphical user interfaces based on a data mining query language
- Architecture of data mining systems