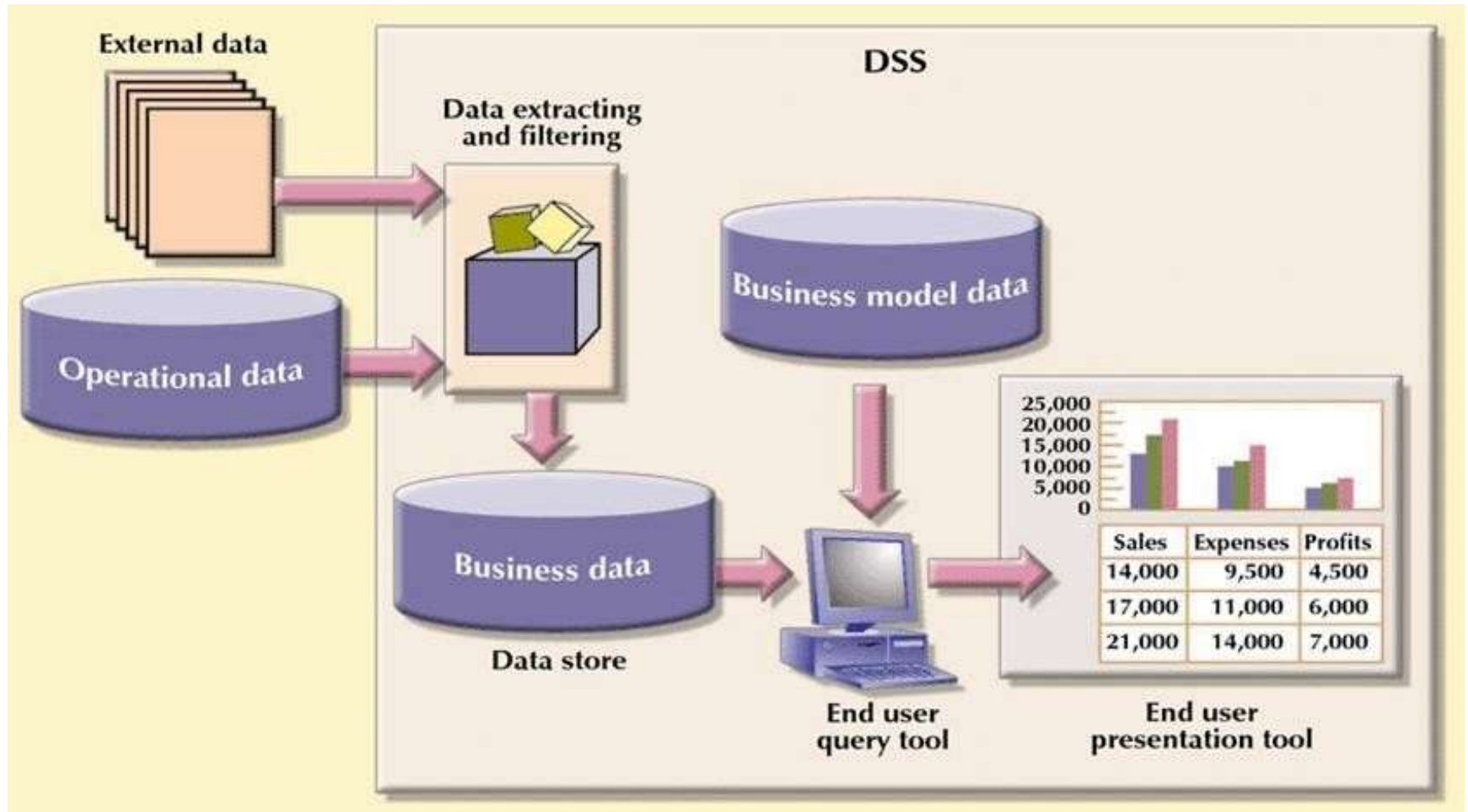# Data Warehouse & Data Mining

## Lecture 12

# Learning Outcomes

- The need for data analysis

- What a data warehouse is, how to prepare data for one, and how to implement one

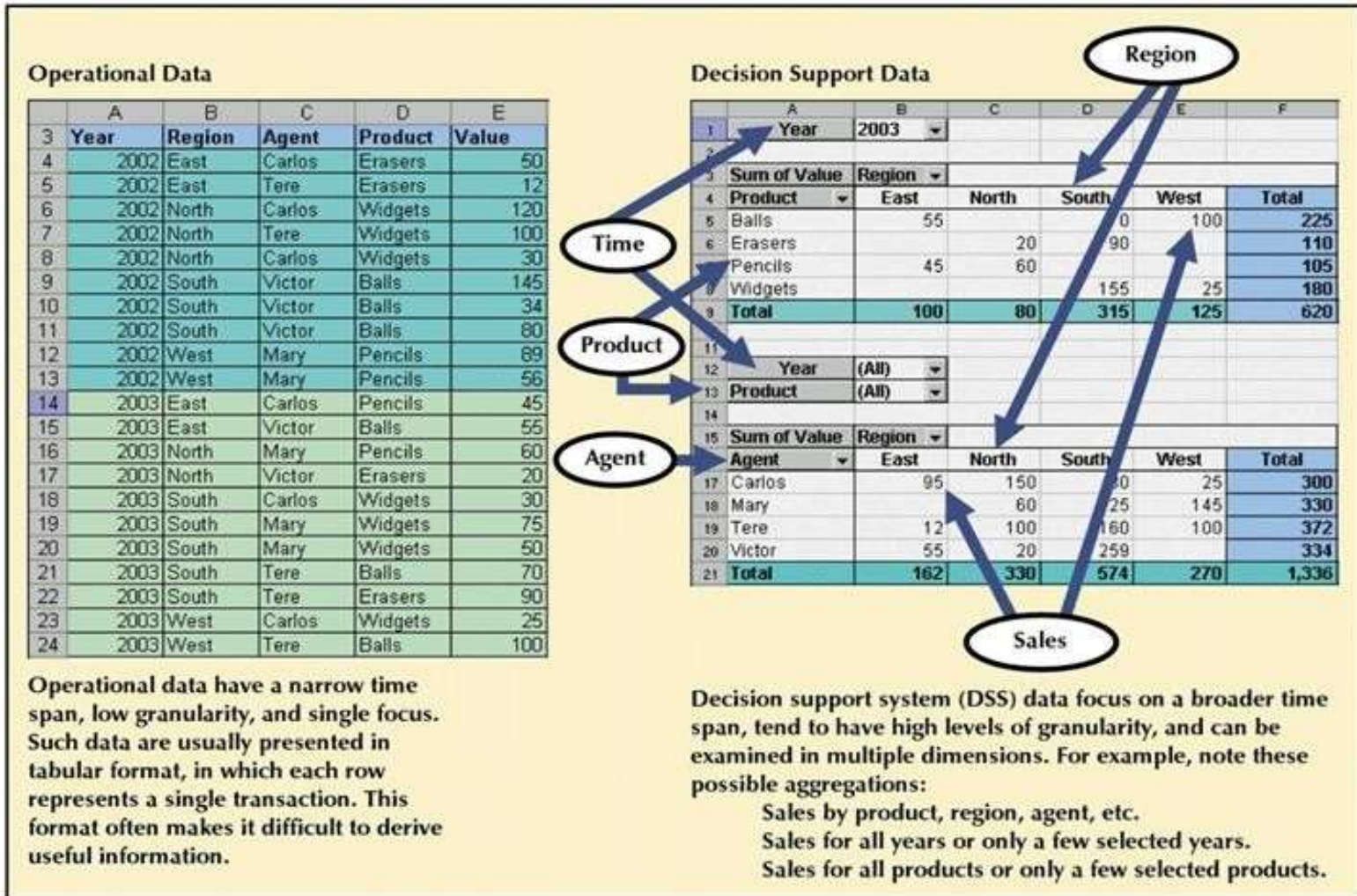- About data analytics, data mining, and predictive analytics

# The Need for Data Analysis

- Organizations are growing rapidly
  - Search for competitive advantage

- Managers needs to track daily transactions to evaluate how the business is performing

- Decision support system
  - Computerized tools used to extract information from data to assist managerial business decision-making

# Decision Support System

CIT6114 – Database Fundamentals

# Transforming Operational Data Into Decision Support Data



**Operational Data**

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 3 | Year | Region | Agent | Product | Value |
| 4 | 2002 | East | Carlos | Erasers | 50 |
| 5 | 2002 | East | Tere | Erasers | 12 |
| 6 | 2002 | North | Carlos | Widgets | 120 |
| 7 | 2002 | North | Tere | Widgets | 100 |
| 8 | 2002 | North | Carlos | Widgets | 30 |
| 9 | 2002 | South | Victor | Balls | 145 |
| 10 | 2002 | South | Victor | Balls | 34 |
| 11 | 2002 | South | Victor | Balls | 80 |
| 12 | 2002 | West | Mary | Pencils | 89 |
| 13 | 2002 | West | Mary | Pencils | 56 |
| 14 | 2003 | East | Carlos | Pencils | 45 |
| 15 | 2003 | East | Victor | Balls | 55 |
| 16 | 2003 | North | Mary | Pencils | 60 |
| 17 | 2003 | North | Victor | Erasers | 20 |
| 18 | 2003 | South | Carlos | Widgets | 30 |
| 19 | 2003 | South | Mary | Widgets | 75 |
| 20 | 2003 | South | Mary | Widgets | 50 |
| 21 | 2003 | South | Tere | Balls | 70 |
| 22 | 2003 | South | Tere | Erasers | 90 |
| 23 | 2003 | West | Carlos | Widgets | 25 |
| 24 | 2003 | West | Tere | Balls | 100 |

Operational data have a narrow time span, low granularity, and single focus. Such data are usually presented in tabular format, in which each row represents a single transaction. This format often makes it difficult to derive useful information.

**Decision Support Data**

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | Year | 2003 ▾ | | | | |
| 2 | | | | | | |
| 3 | Sum of Value | Region ▾ | | | | |
| 4 | Product ▾ | East | North | South | West | Total |
| 5 | Balls | 55 | | 0 | 100 | 225 |
| 6 | Erasers | | 20 | 90 | | 110 |
| 7 | Pencils | 45 | 60 | | | 105 |
| 8 | Widgets | | | 155 | 25 | 180 |
| 9 | Total | 100 | 80 | 315 | 125 | 620 |
| 10 | | | | | | |
| 11 | | | | | | |
| 12 | Year | (All) ▾ | | | | |
| 13 | Product | (All) ▾ | | | | |
| 14 | | | | | | |
| 15 | Sum of Value | Region ▾ | | | | |
| 16 | Agent ▾ | East | North | South | West | Total |
| 17 | Carlos | 95 | 150 | 30 | 25 | 300 |
| 18 | Mary | | 60 | 25 | 145 | 330 |
| 19 | Tere | 12 | 100 | 160 | 100 | 372 |
| 20 | Victor | 55 | 20 | 259 | | 334 |
| 21 | Total | 162 | 330 | 574 | 270 | 1,336 |

Decision support system (DSS) data focus on a broader time span, tend to have high levels of granularity, and can be examined in multiple dimensions. For example, note these possible aggregations:

Sales by product, region, agent, etc.
Sales for all years or only a few selected years.
Sales for all products or only a few selected products.

CIT6114 – Database Fundamentals

# Data Warehouse

- A data warehouse is an *integrated, subject-oriented, time-variant, non-volatile* database that provides support for decision-making

  - Integrated
    - The Data Warehouse is a centralized, consolidated database that integrates data retrieved from the entire organization.

  - Subject-Oriented
    - The Data Warehouse data is arranged and optimized to provide answers to questions coming from diverse functional areas within a company.

CIT6114 – Database Fundamentals

# The Data Warehouse

- ## Time Variant
  - The Warehouse data represent the flow of data through time. It can even contain projected data.

- ## Non-Volatile
  - Once data enter the Data Warehouse, they are never removed.
  - The Data Warehouse is always growing.

**TABLE 13.5**  A Comparison of Data Warehouse and Operational Database Characteristics

| CHARACTERISTIC | OPERATIONAL DATABASE DATA | DATA WAREHOUSE DATA |
|---|---|---|
| Integrated | Similar data can have different representations or meanings. For example, Social Security numbers may be stored as ###-##-#### or as #########, and a given condition may be labeled as T/F or 0/1 or Y/N. A sales value may be shown in thousands or in millions. | Provide a unified view of all data elements with a common definition and representation for all business units. |
| Subject-oriented | Data are stored with a functional, or process, orientation. For example, data may be stored for invoices, payments, and credit amounts. | Data are stored with a subject orientation that facilitates multiple views of the data and facilitates decision making. For example, sales may be recorded by product, by division, by manager, or by region. |
| Time-variant | Data are recorded as current transactions. For example, the sales data may be the sale of a product on a given date, such as $342.78 on 12-MAY-2004. | Data are recorded with a historical perspective in mind. Therefore, a time dimension is added to facilitate data analysis and various time comparisons. |
| Nonvolatile | Data updates are frequent and common. For example, an inventory amount changes with each sale. Therefore, the data environment is fluid. | Data cannot be changed. Data are added only periodically from historical systems. Once the data are properly stored, no changes are allowed. Therefore, the data environment is relatively static. |

CIT6114 – Database Fundamentals

# Role of Data Warehouse

- Focal point for decision support systems
- Support decision making by allowing users to :
  - *drill-down for a more detailed information,*
  - *roll-up to view summarized information,*
  - *slice and dice a dimension for selection of a specific item of interest and*
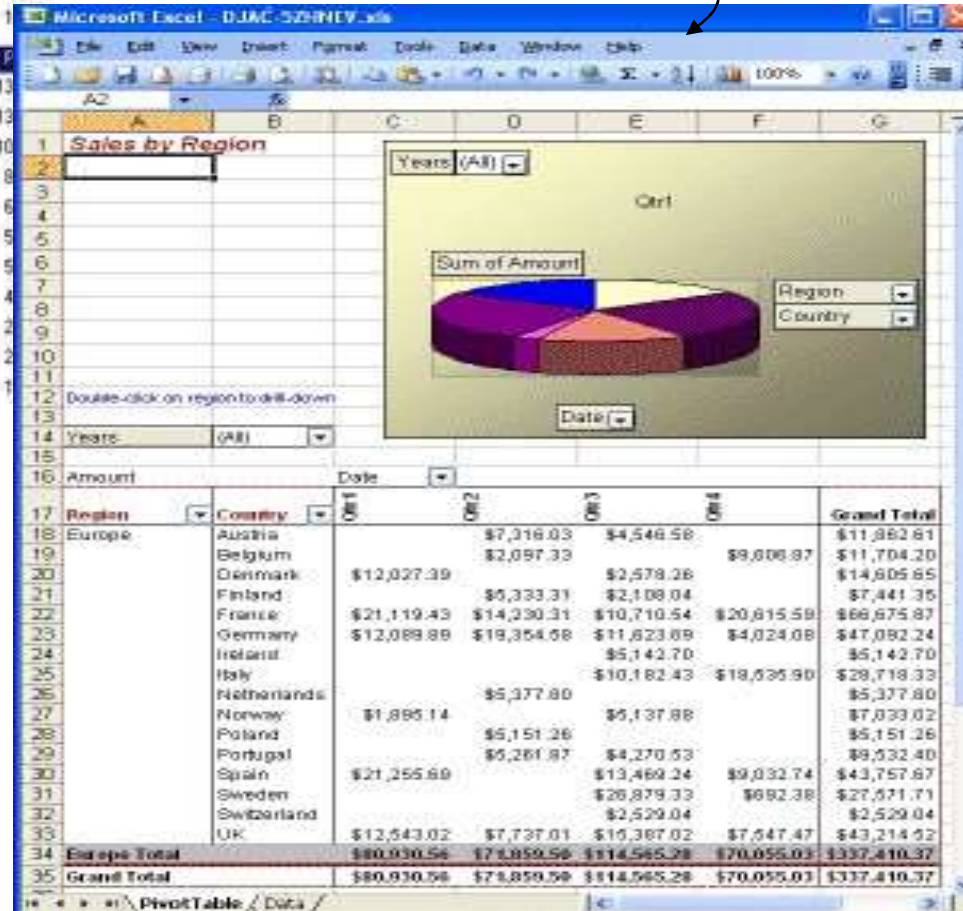  - *pivot to re-orientate the view of multidimensional data.*

Drill-down and Roll-up

Pivot table

Slice & dice

CIT6114 – Database Fundamentals

# The ETL Process



Operational data

Data warehouse

Transformation

Extraction

Loading

- Filter
- Transform
- Integrate
- Classify
- Aggregate
- Summarize

- Integrated
- Subject-oriented
- Time-variant
- Nonvolatile

SOURCE: Course Technology/Cengage Learning

CIT6114 – Database Fundamentals

# Twelve Rules That Define a Data Warehouse

- Data warehouse and operational environments are **separated**

- Data warehouse data are **integrated**

- Data warehouse contains **historical data** over long time

- Data warehouse data are snapshot **data captured at given point in time**

- Data warehouse data are **subject-oriented**

# Twelve Rules That Define a Data Warehouse (cont'd.)

- Data warehouse data are mainly **read-only**
  - Periodic batch updates from operational data
  - No online updates allowed
- Data warehouse development life cycle **differs** from classical systems development
- Data warehouse contains data with **several levels of detail**:
  - Current detail data, old detail data, lightly summarized data, and highly summarized data

# Twelve Rules That Define a Data Warehouse (cont'd.)

- Read-only transactions to **very large data sets**

- Data warehouse environment **traces data sources, transformations, and storage**

- Data warehouse's **metadata are critical** component of this environment

- Data warehouse contains **chargeback mechanism** for resource usage

  - Enforces optimal use of data by end users – users are charged when involving data warehouse processing.

# Implementing a Data Warehouse

- Numerous constraints, including:
  - Available funding
  - Management's view of role played by an IS department
    - Extent and depth of information requirements
  - Corporate culture

- No single formula can describe perfect data warehouse development

# The Data Warehouse as an Active Decision Support Framework

- ## Data warehouse:

  - *Is not a static database*

  - *Is a dynamic framework for decision support that is always a work in progress*

CIT6114 – Database Fundamentals

**FIGURE 13.21** Data warehouse design and implementation road map

Initial data gathering
- Identify and interview key users
- Define main subjects
- Identify operational data model
- Define ownership of data
- Define frequency of use and update
- Define end-user interface
- Define outputs

Design and mapping
- Design star schema
- Facts, dimensions, attributes
- Create star schema diagrams
- Attribute hierarchies
- Map to relational tables
- Naming conventions

Loading and testing
- Prepare for loading
- Define initial and update processes
- Define transformation
- Map from operational data
- Integrate and transform
- Load data, index data, and validate data
- Verify metadata and star schemas

Building and testing
- Training in development environment
- Build menus
- Customize query tools
- Build required queries
- Lay out outputs
- Test interfaces and results
- Optimize for speed and accuracy
- End-user prototyping and testing

Rollout and feedback
- Roll out system
- Get end-user feedback
- System maintenance
- System expansion
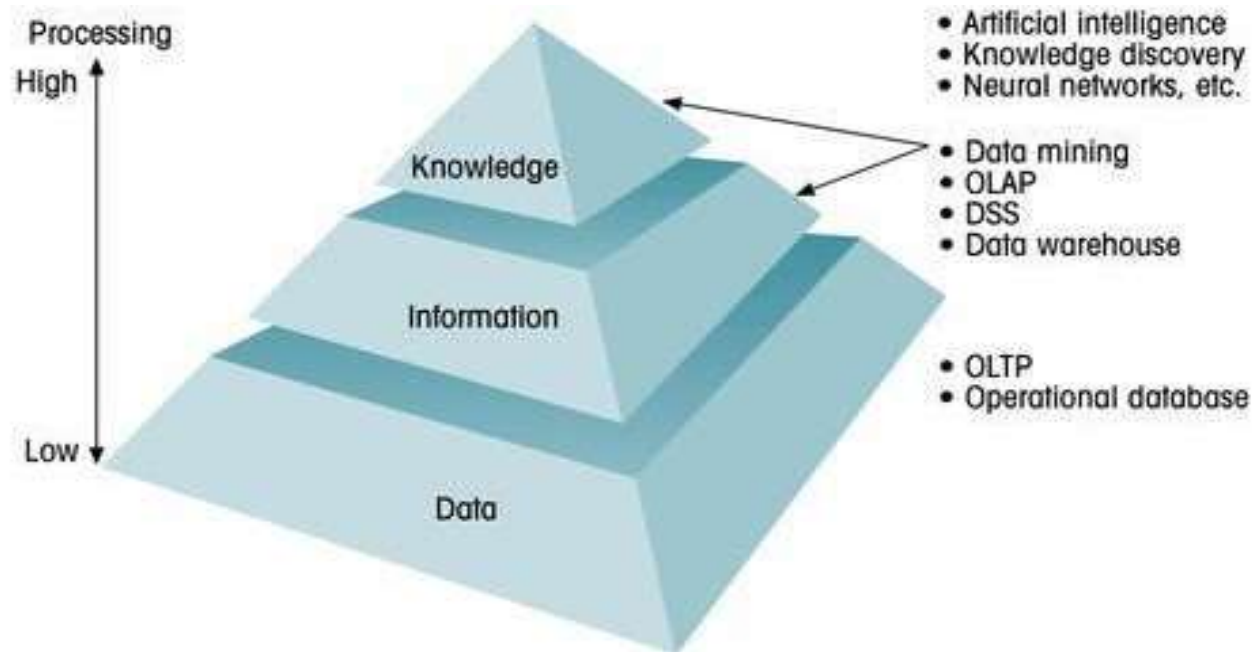
CIT6114 – Database Fundamentals

# Data Mart

- A data mart is a **small**, **single-subject data warehouse** subset that provides decision support to a small group of people.

- Data Marts can serve as a <span style="color:red">test vehicle</span> for companies exploring the potential benefits of Data Warehouses.

- <span style="color:red">Data Marts</span> - address local or departmental problems

- <span style="color:blue">Data Warehouse</span> - involves a company-wide effort to support decision-making at all levels in the organization.

CIT6114 – Database Fundamentals

# Data Mining

- Data mining tools **automatically** search the data for anomalies and identify possible relationships, thereby identifying problems that have not yet been identified by the end user.

- Requires minimal end-user intervention

# Data Mining (cont')



Processing
High

Low

Knowledge

Information

Data

- Artificial intelligence
- Knowledge discovery
- Neural networks, etc.

- Data mining
- OLAP
- DSS
- Data warehouse

- OLTP
- Operational database

Data-mining tools use advanced techniques from knowledge discovery, artificial intelligence, and other fields to obtain "knowledge" and apply it to business needs. Knowledge is then used to make predictions of events or forecasts of values such as sales returns, etc. Several OLAP tools have integrated at least some of these data-mining features in their products.

FIGURE 13.22 ▪ EXTRACTION OF KNOWLEDGE FROM DATA

# Phases in Data Mining

1. *Data Preparation*

   - Identify, collect and consolidating data for analysis

2. *Data Analysis*

   - Identify and explore data with the goal of discovering useful information

3. *Knowledge Acquisition*

   - Select the appropriate modeling or knowledge acquisition algorithms.
   - Examples: neural networks, decision trees, etc.

4. *Prognosis*

   - Predict future behavior and forecast business outcomes using the data mining findings.

CIT6114 – Database Fundamentals