

# CSE 592 Data Mining

Instructor: Pedro Domingos

1

## Today's Program

- Logistics and introduction
- Data warehousing and OLAP

2

## Logistics

- **Instructor:** Pedro Domingos
  - Email: pedrod@cs
  - Office: Sieg 216
  - Office hours: TBA
- **TA:** David Grimes
  - Email: grimes@cs
  - Office: TBA
  - Office hours: TBA
- **Web:** [www.cs.washington.edu/592](http://www.cs.washington.edu/592)
- **Mailing list:** cse574@cs

3

## Assignments

- **Two projects**
  - Groups of two
  - First project: Clickstream mining (37.5%)
  - Second project: Collaborative filtering (25%)
- **Three homeworks**
  - Individual
  - 12.5% each

4

## Source Materials

- Jiawei Han & Micheline Kamber, *Data Mining: Concepts and Techniques*, Morgan Kaufmann, 2001.
- Tom Mitchell, *Machine Learning*, McGraw-Hill, 1997.
- Papers

5

## What is Data Mining?

- **Data mining** is the process of identifying valid, novel, useful and understandable patterns in data.
- Also known as **KDD** (**K**nowledge **D**iscovery in **D**atabases).
- "We're drowning in information, but starving for knowledge." (John Naisbett)

6

## Related Disciplines

- Machine learning
- Databases
- Statistics
- Information retrieval
- Visualization
- High-performance computing
- Etc.

7

## Applications of Data Mining

- E-commerce
- Marketing and retail
- Finance
- Telecoms
- Drug design
- Process control
- Space and earth sensing
- Etc.

8

## The Data Mining Process

- Understanding domain, prior knowledge, and goals
- Data integration and selection
- Data cleaning and pre-processing
- Modeling and searching for patterns
- Interpreting results
- Consolidating and deploying discovered knowledge
- Loop

9

## Data Mining Tasks

- Classification
- Regression
- Probability estimation
- Clustering
- Association detection
- Summarization
- Trend and deviation detection
- Etc.

10

## Inductive Learning

- **Inductive learning** or **Prediction**:
  - **Given** examples of a function  $(X, F(X))$
  - **Predict** function  $F(X)$  for new examples  $X$
- Discrete  $F(X)$ : Classification
- Continuous  $F(X)$ : Regression
- $F(X) = \text{Probability}(X)$ : Probability estimation

11

## Widely-used Approaches

- Decision trees
- Rule induction
- Bayesian learning
- Neural networks
- Genetic algorithms
- Instance-based learning
- Etc.

12

## Requirements for a Data Mining System

- Data mining systems should be
  - Computationally sound
  - Statistically sound
  - Ergonomically sound

13

## Components of a Data Mining System

- Representation
- Evaluation
- Search
- Data management
- User interface

14

## Topics for this Quarter (Slide 1 of 2)

- Data warehousing and OLAP
- Decision trees
- Rule induction
- Bayesian learning
- Neural networks
- Genetic algorithms

15

## Topics for this Quarter (Slide 2 of 2)

- Model ensembles
- Instance-based learning
- Learning theory
- Association rules
- Clustering

16

## Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- Extensions of data cubes
- From data warehousing to data mining

17

## What is a Data Warehouse?

- Defined in many different ways, but not rigorously.
  - A decision support database that is maintained separately from the organization's operational database
  - Support information processing by providing a solid platform of consolidated, historical data for analysis.
- "A data warehouse is a subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process."—W. H. Inmon
- Data warehousing:
  - The process of constructing and using data warehouses

18

## Data Warehouse—Subject-Oriented

- Organized around major subjects, such as customer, product, sales.
- Focusing on the modeling and analysis of data for decision makers, not on daily operations or transaction processing.
- Provide a simple and concise view around particular subject issues by excluding data that are not useful in the decision support process.

19

## Data Warehouse—Integrated

- Constructed by integrating multiple, heterogeneous data sources
  - relational databases, flat files, on-line transaction records
- Data cleaning and data integration techniques are applied.
  - Ensure consistency in naming conventions, encoding structures, attribute measures, etc. among different data sources
    - E.g., Hotel price: currency, tax, breakfast covered, etc.
- When data is moved to the warehouse, it is converted.

20

## Data Warehouse—Time Variant

- The time horizon for the data warehouse is significantly longer than that of operational systems.
  - Operational database: current value data.
  - Data warehouse data: provide information from a historical perspective (e.g., past 5-10 years)
- Every key structure in the data warehouse
  - Contains an element of time, explicitly or implicitly
  - But the key of operational data may or may not contain "time element".

21

## Data Warehouse—Non-Volatile

- A physically separate store of data transformed from the operational environment.
- Operational update of data does not occur in the data warehouse environment.
  - Does not require transaction processing, recovery, and concurrency control mechanisms
  - Requires only two operations in data accessing:
    - *initial loading of data* and *access of data*.

22

## Data Warehouse vs. Heterogeneous DBMS

- Traditional heterogeneous DB integration:
  - Build wrappers/mediators on top of heterogeneous databases
  - Query driven approach
    - When a query is posed to a client site, a meta-dictionary is used to translate the query into queries appropriate for individual heterogeneous sites involved, and the results are integrated into a global answer set
    - Complex information filtering, compete for resources
- Data warehouse: update-driven, high performance
  - Information from heterogeneous sources is integrated in advance and stored in warehouses for direct query and analysis

23

## Data Warehouse vs. Operational DBMS

- OLTP (on-line transaction processing)
  - Major task of traditional relational DBMS
  - Day-to-day operations: purchasing, inventory, banking, manufacturing, payroll, registration, accounting, etc.
- OLAP (on-line analytical processing)
  - Major task of data warehouse system
  - Data analysis and decision making
- Distinct features (OLTP vs. OLAP):
  - User and system orientation: customer vs. market
  - Data contents: current, detailed vs. historical, consolidated
  - Database design: ER + application vs. star + subject
  - View: current, local vs. evolutionary, integrated
  - Access patterns: update vs. read-only but complex queries

24

## OLTP vs. OLAP

	OLTP	OLAP
<b>users</b>	clerk, IT professional	knowledge worker
<b>function</b>	day to day operations	decision support
<b>DB design</b>	application-oriented	subject-oriented
<b>data</b>	current, up-to-date detailed, flat relational isolated	historical, summarized, multidimensional integrated, consolidated
<b>usage</b>	repetitive	ad-hoc
<b>access</b>	read/write index/hash on prim. key	lots of scans
<b>unit of work</b>	short, simple transaction	complex query
<b># records accessed</b>	tens	millions
<b>#users</b>	thousands	hundreds
<b>DB size</b>	100MB-GB	100GB-TB
<b>metric</b>	transaction throughput	query throughput, response

25

## Why Separate Data Warehouse?

- High performance for both systems
  - DBMS— tuned for OLTP: access methods, indexing, concurrency control, recovery
  - Warehouse—tuned for OLAP: complex OLAP queries, multidimensional view, consolidation.
- Different functions and different data:
  - missing data:** Decision support requires historical data which operational DBs do not typically maintain
  - data consolidation:** DS requires consolidation (aggregation, summarization) of data from heterogeneous sources
  - data quality:** different sources typically use inconsistent data representations, codes and formats which have to be reconciled

26

## Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- Extensions of data cubes
- From data warehousing to data mining

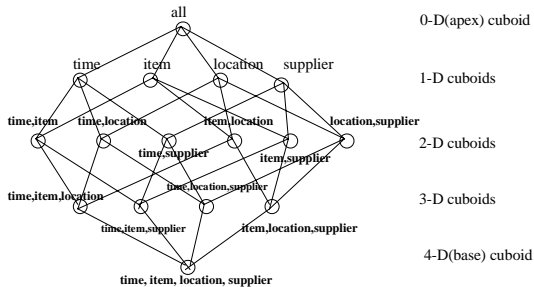
27

## From Tables and Spreadsheets to Data Cubes

- A data warehouse is based on a multidimensional data model which views data in the form of a data cube
- A data cube, such as sales, allows data to be modeled and viewed in multiple dimensions
  - Dimension tables, such as item (item\_name, brand, type), or time(day, week, month, quarter, year)
  - Fact table contains measures (such as dollars\_sold) and keys to each of the related dimension tables
- In data warehousing literature, an n-D base cube is called a base cuboid. The topmost 0-D cuboid, which holds the highest-level of summarization, is called the apex cuboid. The lattice of cuboids forms a data cube.

28

## Cube: A Lattice of Cuboids

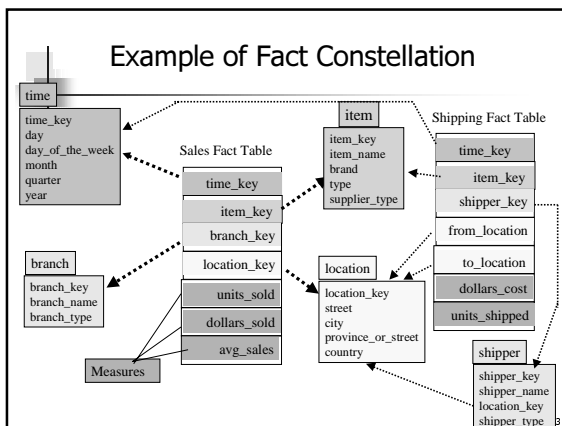
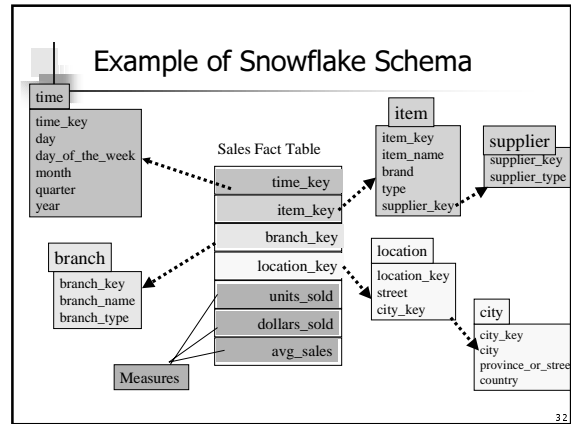
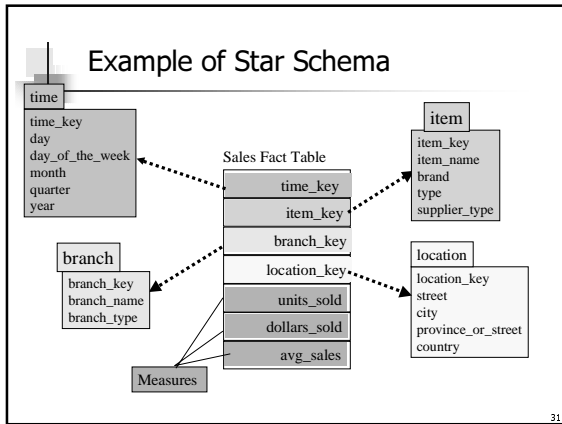


29

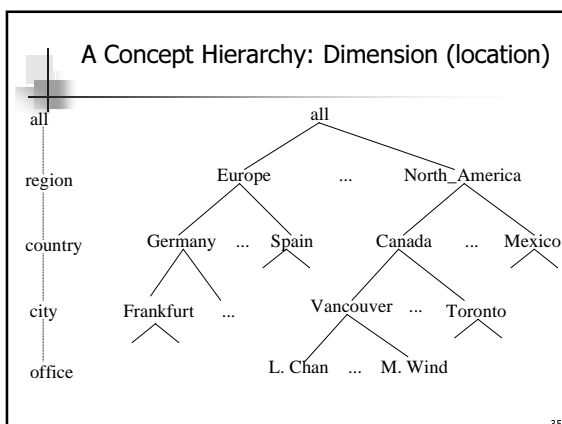
## Conceptual Modeling of Data Warehouses

- Modeling data warehouses: dimensions & measures
  - Star schema:** A fact table in the middle connected to a set of dimension tables
  - Snowflake schema:** A refinement of star schema where some dimensional hierarchy is normalized into a set of smaller dimension tables, forming a shape similar to snowflake
  - Fact constellations:** Multiple fact tables share dimension tables, viewed as a collection of stars, therefore called galaxy schema or fact constellation

30



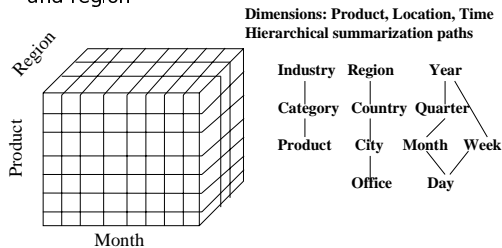
- ### Measures: Three Categories
- **distributive**: if the result derived by applying the function to  $n$  aggregate values is the same as that derived by applying the function on all the data without partitioning.
    - E.g., `count()`, `sum()`, `min()`, `max()`.
  - **algebraic**: if it can be computed by an algebraic function with  $M$  arguments (where  $M$  is a bounded integer), each of which is obtained by applying a distributive aggregate function.
    - E.g., `avg()`, `min_N()`, `standard_deviation()`.
  - **holistic**: if there is no constant bound on the storage size needed to describe a subaggregate.
    - E.g., `median()`, `mode()`, `rank()`.



- ### Specification of Hierarchies
- **Schema hierarchy**
    - day < {month < quarter; week} < year
  - **Set\_grouping hierarchy**
    - {1..10} < inexpensive

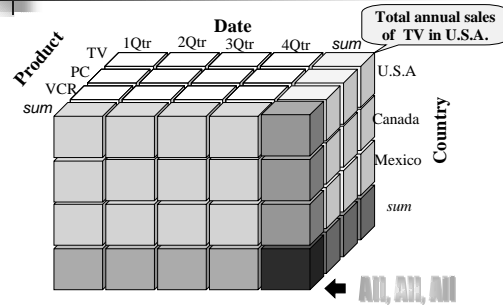
## Multidimensional Data

- Sales volume as a function of product, month, and region



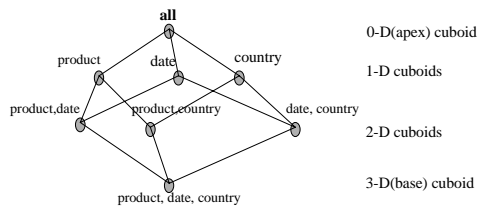
37

## A Sample Data Cube



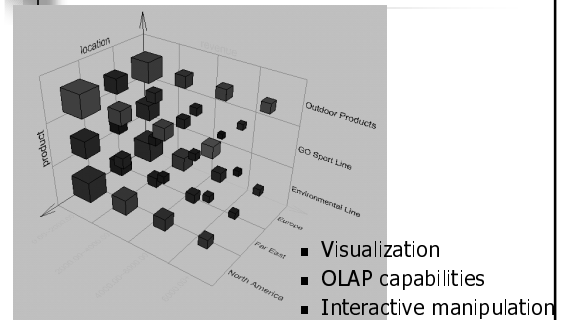
38

## Cuboids Corresponding to the Cube



39

## Browsing a Data Cube



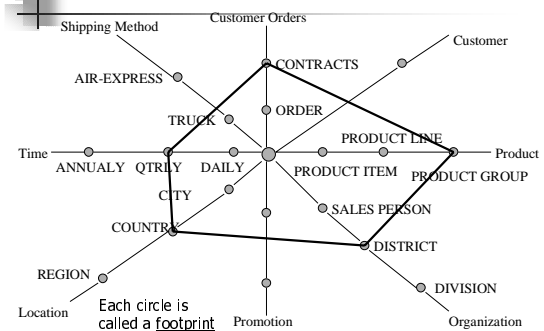
40

## Typical OLAP Operations

- Roll up (drill-up): summarize data
  - by climbing up hierarchy or by dimension reduction
- Drill down (roll down): reverse of roll-up
  - from higher level summary to lower level summary or detailed data, or introducing new dimensions
- Slice and dice:
  - project and select
- Pivot (rotate):
  - reorient the cube, visualization, 3D to series of 2D planes.
- Other operations
  - drill across: involving (across) more than one fact table
  - drill through: through the bottom level of the cube to its back-end relational tables (using SQL)

41

## A Starnet Query Model



42

## Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- Extensions of data cubes
- From data warehousing to data mining

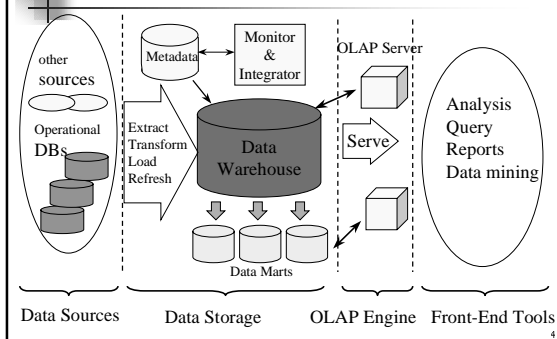
43

## Data Warehouse Design Process

- Choose the *grain* (*atomic level of data*) of the business process
- Choose a business process to model, e.g., orders, invoices, etc.
- Choose the dimensions that will apply to each fact table record
- Choose the measure that will populate each fact table record

44

## Multi-Tiered Architecture



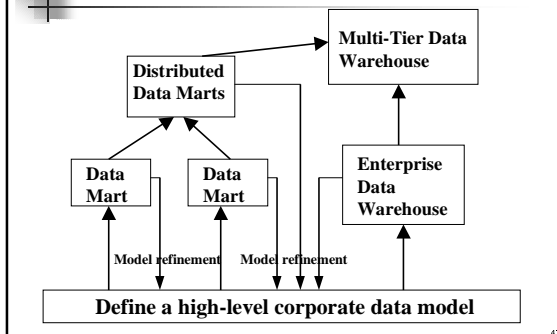
45

## Three Data Warehouse Models

- Enterprise warehouse**
  - collects all of the information about subjects spanning the entire organization
- Data Mart**
  - a subset of corporate-wide data that is of value to a specific groups of users. Its scope is confined to specific, selected groups, such as marketing data mart
    - Independent vs. dependent (directly from warehouse) data mart
- Virtual warehouse**
  - A set of views over operational databases
  - Only some of the possible summary views may be materialized

46

## Data Warehouse Development: A Recommended Approach



47

## OLAP Server Architectures

- Relational OLAP (ROLAP)**
  - Use relational or extended-relational DBMS to store and manage warehouse data and OLAP middle ware to support missing pieces
  - Include optimization of DBMS backend, implementation of aggregation navigation logic, and additional tools and services
  - greater scalability
- Multidimensional OLAP (MOLAP)**
  - Array-based multidimensional storage engine (sparse matrix techniques)
  - fast indexing to pre-computed summarized data
- Hybrid OLAP (HOLAP)**
  - User flexibility, e.g., low level: relational, high-level: array
- Specialized SQL servers
  - specialized support for SQL queries over star/snowflake schemas

48



## Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- Extensions of data cubes
- From data warehousing to data mining

49

## Efficient Data Cube Computation

- Data cube can be viewed as a lattice of cuboids
  - The bottom-most cuboid is the base cuboid
  - The top-most cuboid (apex) contains only one cell
  - How many cuboids in an n-dimensional cube with L levels?
 
$$T = \prod_{i=1}^n (L_i + 1)$$
- Materialization of data cube
  - Materialize every (cuboid) (full materialization), none (no materialization), or some (partial materialization)
  - Selection of which cuboids to materialize
    - Based on size, sharing, access frequency, etc.

50

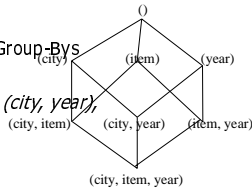
## Cube Operation

- Transform it into SQL-like language (with new operator cube by, introduced by Gray et al.'96)

```
SELECT item, city, year, SUM (amount)
FROM SALES
```

```
CUBE BY item, city, year
```

- Need compute the following Group-Bys
  - (item, city, year),
  - (item, city), (item, year), (city, year),
  - (item), (city), (year)
  - ()



51

## Efficient Processing of OLAP Queries

- Determine which operations should be performed on the available cuboids:
  - transform drill, roll, etc. into corresponding SQL and/or OLAP operations, e.g. dice = selection + projection
- Determine to which materialized cuboid(s) the relevant operations should be applied.
- Exploring indexing structures and compressed vs. dense array structures in MOLAP

52

## Metadata Repository

- Meta data is the data defining warehouse objects. It has the following kinds
  - Description of the structure of the warehouse
    - schema, view, dimensions, hierarchies, derived data defn, data mart locations and contents
  - Operational meta-data
    - data lineage (history of migrated data and transformation path), currency of data (active, archived, or purged), monitoring information (warehouse usage statistics, error reports, audit trails)
  - The algorithms used for summarization
  - The mapping from operational environment to the data warehouse
  - Data related to system performance
  - Business data
    - business terms and definitions, ownership of data, charging policies

53

## Data Warehouse Back-End Tools and Utilities

- Data extraction:
  - get data from multiple, heterogeneous, and external sources
- Data cleaning:
  - detect errors in the data and rectify them when possible
- Data transformation:
  - convert data from legacy or host format to warehouse format
- Load:
  - sort, summarize, consolidate, compute views, check integrity, and build indices and partitions
- Refresh
  - propagate the updates from the data sources to the warehouse

54

## Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- Extensions of data cubes
- From data warehousing to data mining

55

## Discovery-Driven Exploration of Data Cubes

- Hypothesis-driven: exploration by user, huge search space
- Discovery-driven (Sarawagi et al. '98)
  - pre-compute measures indicating exceptions, guide user in the data analysis, at all levels of aggregation
  - Exception: significantly different from the value anticipated, based on a statistical model
  - Visual cues such as background color are used to reflect the degree of exception of each cell
  - Computation of exception indicator (modeling fitting and computing SelfExp, InExp, and PathExp values) can be overlapped with cube construction

56

## Examples: Discovery-Driven Data Cubes

Item		all											
region		all											
Sum of sales		month											
Total		Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
Avg sales		1%	-1%	0%	1%	3%	-1	-9%	-1%	2%	-4%	3%	
Item	month	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
Sony lbr printer		9%	8%	2%	-5%	14%	4%	0%	21%	-13%	15%	-11%	
Sony color printer		0%	0%	3%	2%	4%	-10%	-13%	0%	4%	-6%	4%	
HP lbr printer		0%	1%	2%	3%	8%	0%	-12%	9%	3%	-5%	0%	
HP color printer		0%	0%	-2%	1%	0%	-1%	-7%	-2%	1%	-5%	1%	
IBM home computer		1%	-2%	-1%	-1%	3%	3%	-10%	16%	1%	4%	-1%	
IBM laptop computer		0%	0%	-1%	3%	4%	2%	-10%	2%	0%	9%	3%	
Toshiba home computer		-2%	-5%	1%	1%	-1%	1%	5%	-3%	-5%	-1%	-1%	
Toshiba laptop computer		1%	0%	3%	0%	2%	-2%	-5%	3%	2%	-1%	0%	
Logitech mouse		3%	-2%	-1%	0%	4%	0%	-11%	2%	1%	-4%	0%	
Epson-wp mouse		0%	0%	2%	3%	1%	-2%	-2%	5%	0%	-5%	8%	

Item		IBM home computer											
Avg sales		month											
region		Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
North		-1%	-3%	-1%	0%	3%	4%	-7%	1%	0%	-3%	-3%	
South		-1%	1%	-9%	6%	-1%	39%	9%	-34%	4%	1%	7%	
East		-1%	-2%	2%	-3%	1%	18%	-2%	11%	-3%	-2%	-1%	
West		4%	0%	-1%	-3%	5%	-1%	-18%	8%	5%	-8%	1%	

57

## Complex Aggregation at Multiple Granularities: Multi-Feature Cubes

- Multi-feature cubes (Ross, et al. 1998): Compute complex queries involving multiple dependent aggregates at multiple granularities
- Ex. Grouping by all subsets of {item, region, month}, find the maximum price in 1997 for each group, and the total sales among all maximum price tuples
 

```
select item, region, month, max(price), sum(R.sales)
from purchases
where year = 1997
cube by item, region, month: R
such that R.price = max(price)
```
- Continuing the last example, among the max price tuples, find the min and max shelf life, and find the fraction of the total sales due to tuple that have min shelf life within the set of all max price tuples

58

## Data Warehousing and OLAP

- What is a data warehouse?
- A multi-dimensional data model
- Data warehouse architecture
- Data warehouse implementation
- Extensions of data cubes
- From data warehousing to data mining

59

## Data Warehouse Usage

- Three kinds of data warehouse applications
  - Information processing
    - supports querying, basic statistical analysis, and reporting using crosstabs, tables, charts and graphs
  - Analytical processing
    - multidimensional analysis of data warehouse data
    - supports basic OLAP operations, slice-dice, drilling, pivoting
  - Data mining
    - Knowledge discovery from hidden patterns
    - supports associations, constructing analytical models, performing classification and prediction, and presenting the mining results using visualization tools.
- Differences among the three tasks

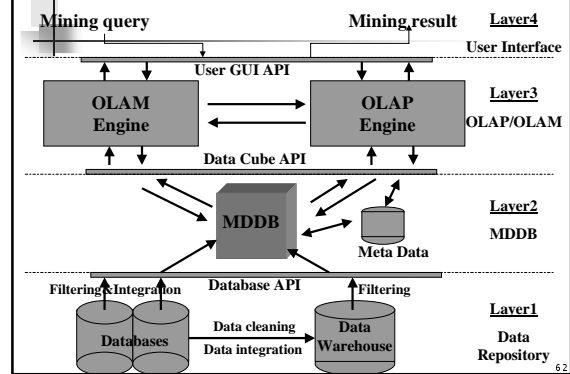
60

## From Online Analytical Processing to Online Analytical Mining (OLAM)

- Why online analytical mining?
  - High quality of data in data warehouses
    - DW contains integrated, consistent, cleaned data
  - Available information processing structure surrounding data warehouses
    - ODBC, OLEDB, Web accessing, service facilities, reporting and OLAP tools
  - OLAP-based exploratory data analysis
    - mining with drilling, dicing, pivoting, etc.
  - On-line selection of data mining functions
    - integration and swapping of multiple mining functions, algorithms, and tasks.
- Architecture of OLAM

61

## An OLAM Architecture



62

## Summary

- Data warehouse
  - A subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process
- A multi-dimensional model of a data warehouse
  - Star schema, snowflake schema, fact constellations
  - A data cube consists of dimensions & measures
- OLAP operations: drilling, rolling, slicing, dicing and pivoting
- OLAP servers: ROLAP, MOLAP, HOLAP
- Efficient computation of data cubes
  - Partial vs. full vs. no materialization
  - Multiway array aggregation
  - Bitmap index and join index implementations
- Extensions of data cubes
  - Discovery-drive and multi-feature cubes
  - From OLAP to OLAM (on-line analytical mining)

63