# Veri Ambarları ve OLAP
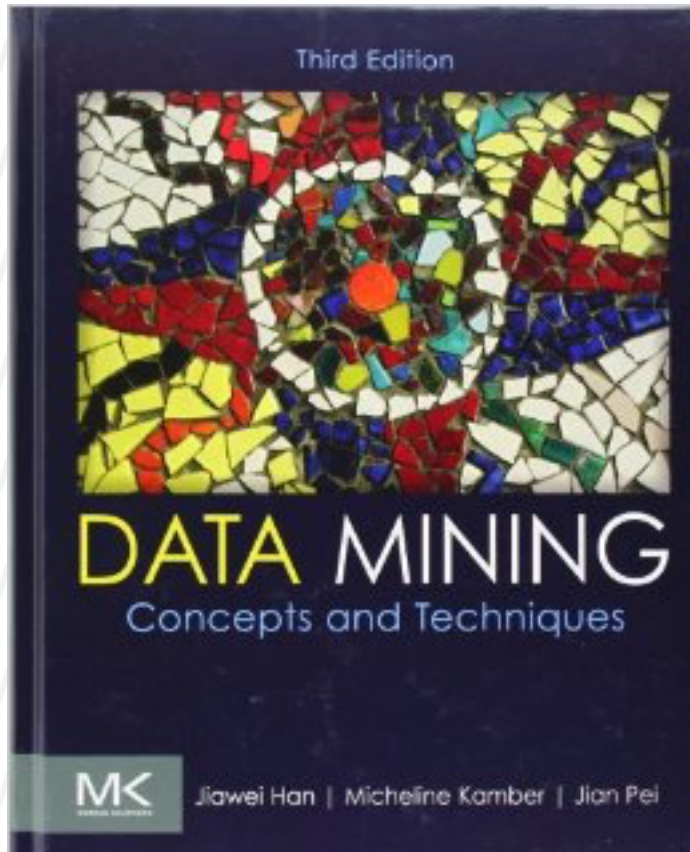
Şadi Evren ŞEKER

YouTube Kanalı: Bilgisayar Kavramları

# Kaynaklar

- **Data Mining: Concepts and Techniques, Third Edition (The Morgan Kaufmann Series in Data Management Systems)** 3rd Edition

- by Jiawei Han (Author), Micheline Kamber (Author), Jian Pei (Author)

# 90'lar ve OLAP

- Veri Ambarı Teknolojileri (OLAP'a ilk geçişler ve OLTP'lerdeki zorluklardan dolayı geçici veri ambarları oluşturma fikri)

- İhtiyaçlar

- Online Analytical Processing

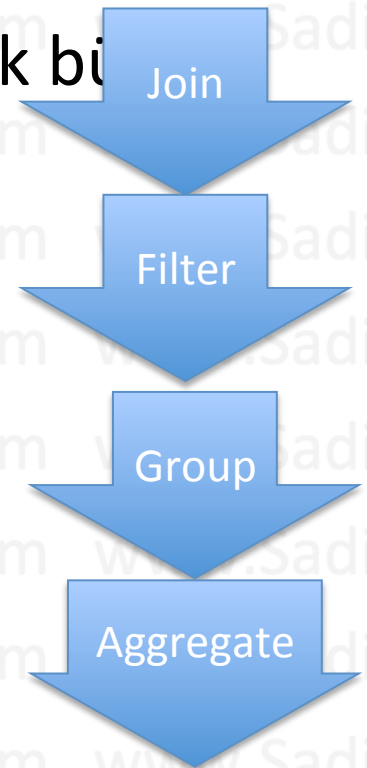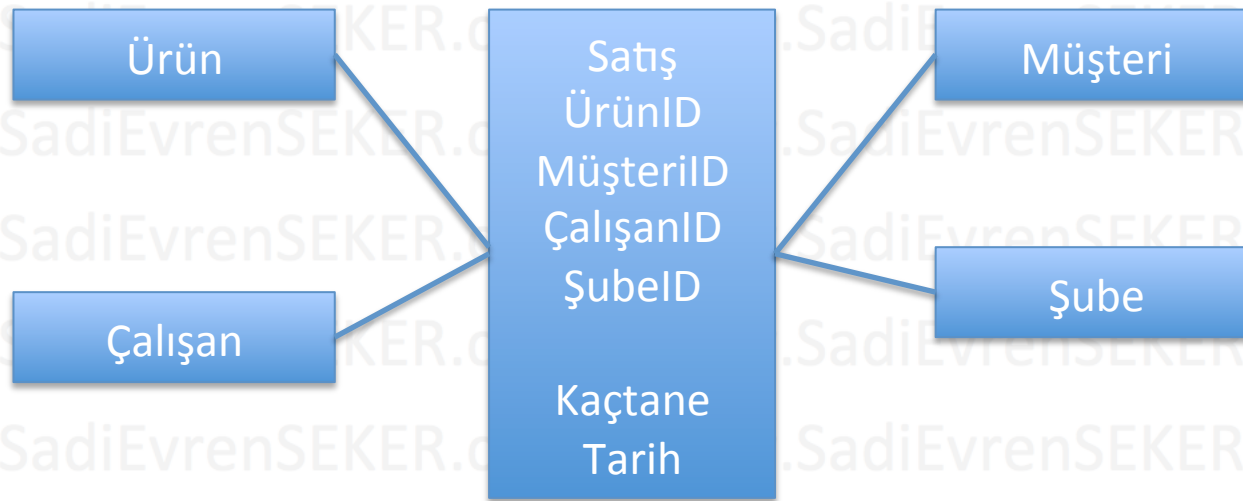- OLTP : Online Transaction Processing

# OLTP vs OLAP

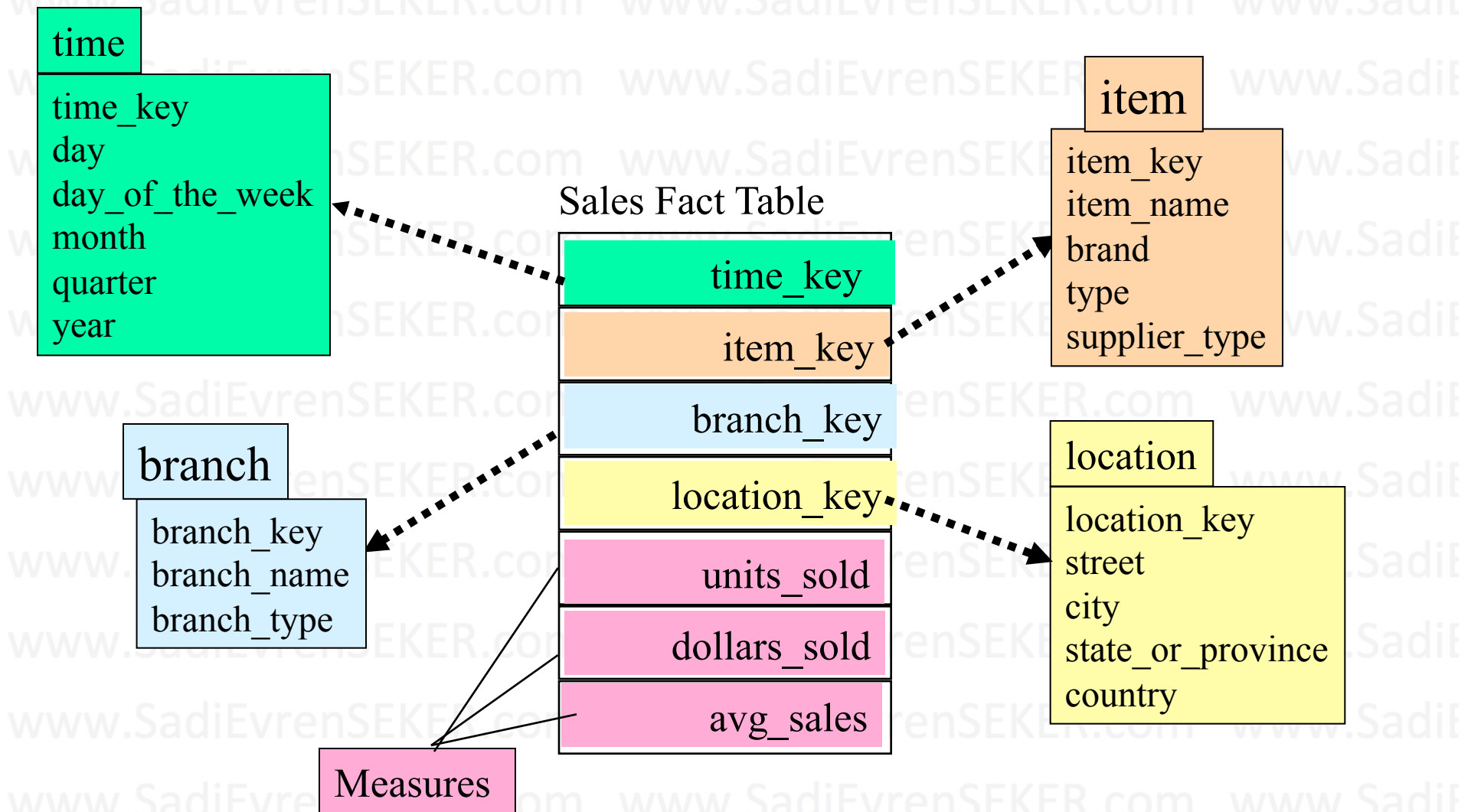|  | OLTP | OLAP |
|---|---|---|
| **Kullanıcılar** | clerk, IT professional | knowledge worker |
| **Fonksiyonlar** | day to day operations | decision support |
| **DB Tasarım** | application-oriented | subject-oriented |
| **Veri** | current, up-to-date detailed, flat relational isolated | historical, summarized, multidimensional integrated, consolidated |
| **Kullanım** | repetitive | ad-hoc |
| **Erişim** | read/write index/hash on prim. key | lots of scans |
| **İşlerin boyu** | short, simple transaction | complex query |
| **# records accessed** | tens | millions |
| **#users** | thousands | hundreds |
| **DB boyutu** | 100MB-GB | 100GB-TB |
| **Metrikler** | transaction throughput | query throughput, response |

# İki Kavram

- OLTP – Online Transaction Processing
  - Örneğin : Banka hesaplarındaki hareketler, bilet işlemleri
  - Genelde küçük transactionlar
  - Verinin küçük bir kısmı ile ilgili
  - Sık ve sürekli tekrarlar şeklinde çalıştırılıyor
- OLAP – Online Analytical Processing
  - Büyük transactionlar
  - Karmaşık sorgular
  - Daha büyük veriye erişim
  - Sık yapılmayan sorgular

# Temel Kavramlar

- Yıldız Şeması (Star Schema)
  - Fact Table : Sık güncellenen, çoğunlukla ekleme yapılan, ve genelde çok büyük tablolardır
  - Dimension Table: Sık güncellenmeyen, çok bü olmayan tablolar

# Star Schema

**time**

time_key
day
day_of_the_week
month
quarter
year

**item**

item_key
item_name
brand
type
supplier_type

Sales Fact Table

| time_key |
| item_key |
| branch_key |
| location_key |
| units_sold |
| dollars_sold |
| avg_sales |

**branch**

branch_key
branch_name
branch_type

**location**

location_key
street
city
state_or_province
country

Measures

# Snowflake Schema

**time**

time_key
day
day_of_the_week
month
quarter
year

**branch**

branch_key
branch_name
branch_type

**Sales Fact Table**

| time_key |
| item_key |
| branch_key |
| location_key |
| units_sold |
| dollars_sold |
| avg_sales |

Measures

**item**

item_key
item_name
brand
type
supplier_key

**supplier**

supplier_key
supplier_type

**location**

location_key
street
city_key

**city**

city_key
city
state_or_province
country

8

# Fact Constellation

**time**

- time_key
- day
- day_of_the_week
- month
- quarter
- year

**item**

- item_key
- item_name
- brand
- type
- supplier_type

Shipping Fact Table

Sales Fact Table

- time_key
- item_key
- branch_key
- location_key
- units_sold
- dollars_sold
- avg_sales

**branch**

- branch_key
- branch_name
- branch_type

Measures

**location**

- location_key
- street
- city
- province_or_state
- country

Shipping Fact Table

- time_key
- item_key
- shipper_key
- from_location
- to_location
- dollars_cost
- units_shipped

**shipper**

- shipper_key
- shipper_name
- location_key
- shipper_type
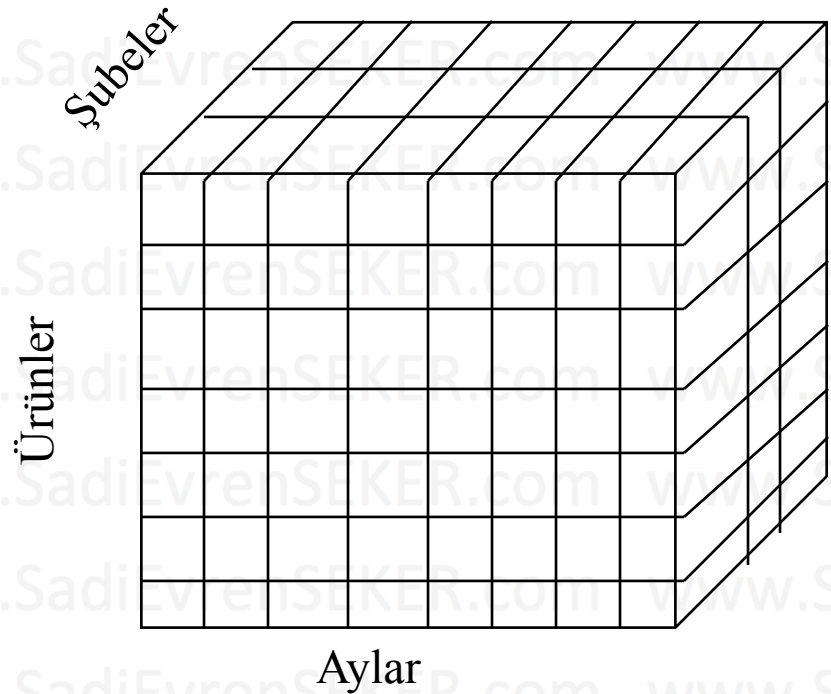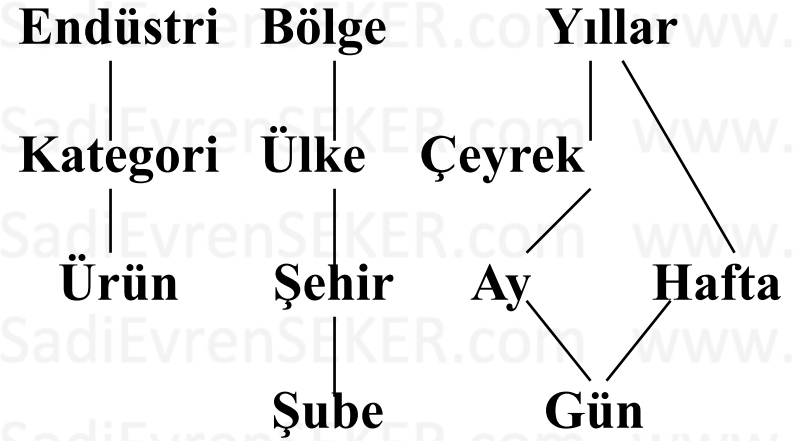
# Star Schema

- Yavaştır : Indeks oluşturulması, joinler, sorguların özel olarak çalıştırılması
- Materialized View

# Veri Küpleri (Data Cube)

Endüstri   Bölge          Yıllar

Kategori   Ülke   Çeyrek

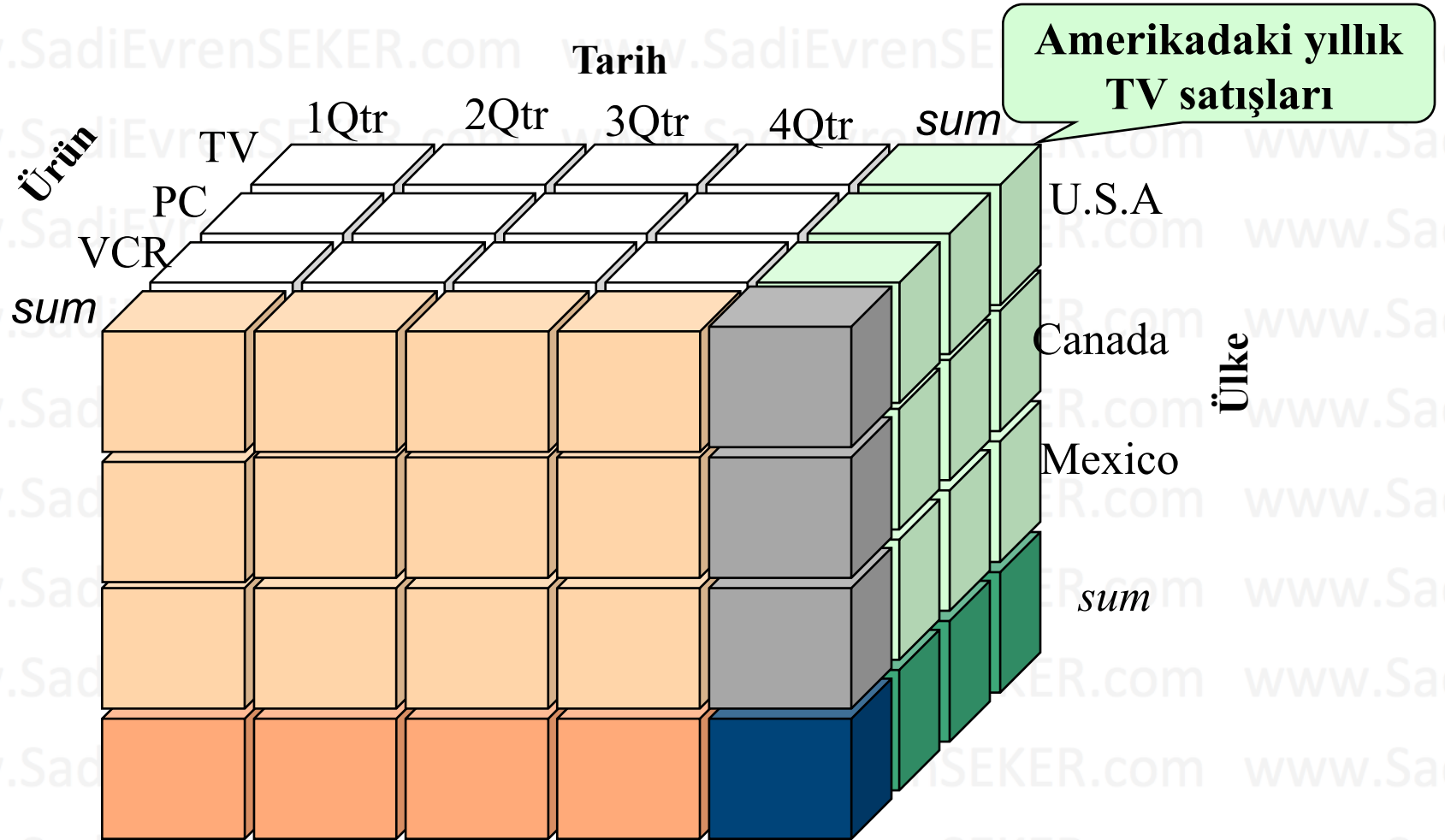Ürün   Şehir   Ay   Hafta

Şube   Gün

- Aslında Küp Değildirler
- Çok boyutlu OLAP (multidimensional OLAP) olarak da isimlendirilirler
- Fact Data hücrelerde durmaktadır
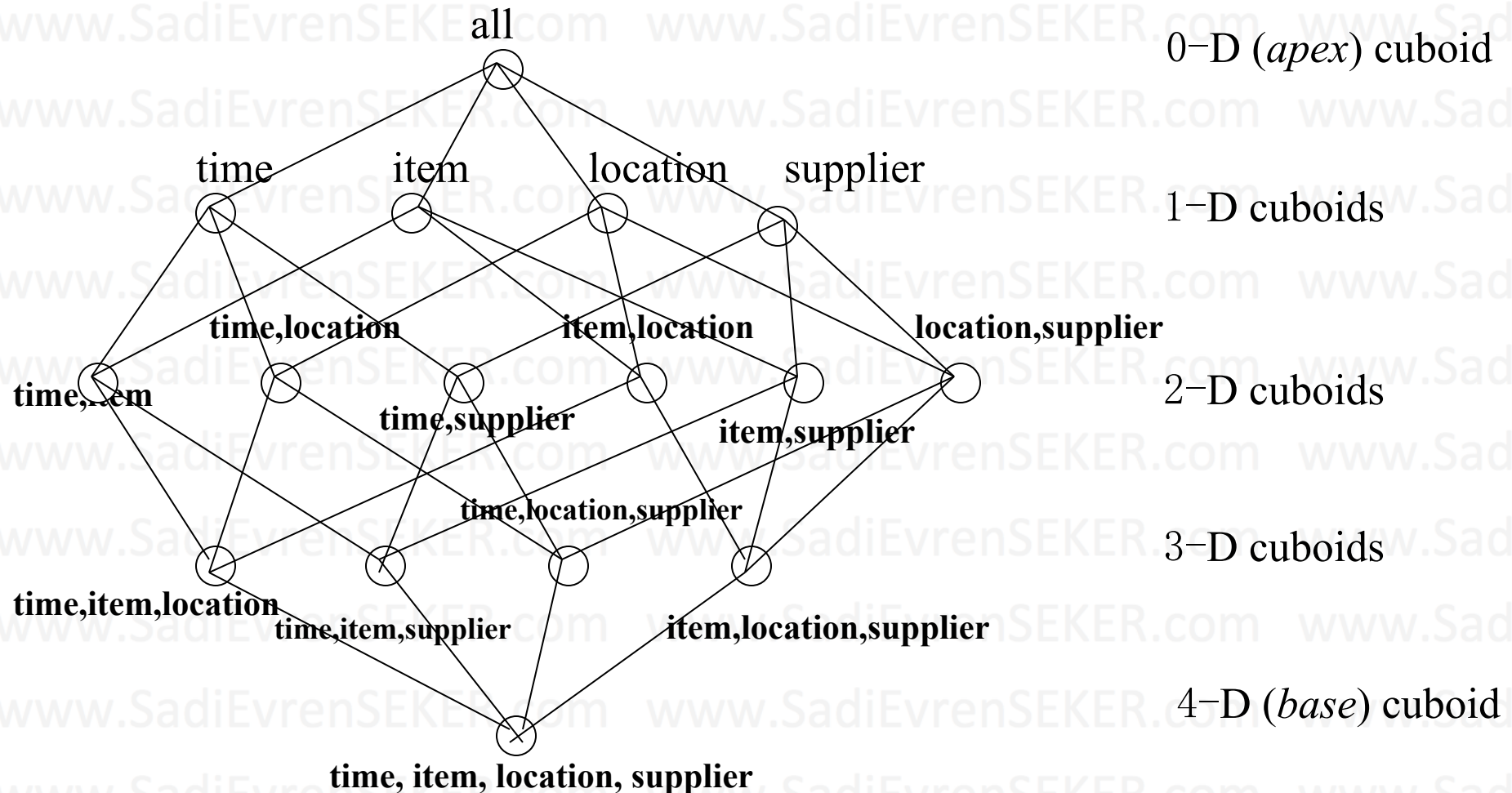- Slide, Edge ve Corner üzerinde aggregated data tutulmaktadır

Şubeler

Ürünler

Aylar

# Örnek Veri Küpü



Amerikadaki yıllık TV satışları

# Bazı Aggregate Taktikleri

- Dimension attribute şayet key değilse genelde aggregate edilir
  - 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()

# Cube: A Lattice of Cuboids

all

0–D (*apex*) cuboid

time          item          location          supplier

1–D cuboids

time,location          item,location          location,supplier

time,item

time,supplier

item,supplier

2–D cuboids

time,location,supplier

3–D cuboids

time,item,location

time,item,supplier          item,location,supplier

4–D (*base*) cuboid

time, item, location, supplier

14

# Typical OLAP Operations

- Roll up (drill-up): özetleme
  - *by climbing up hierarchy or by dimension reduction*
- Drill down (roll down): detaylandırma
  - *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)*
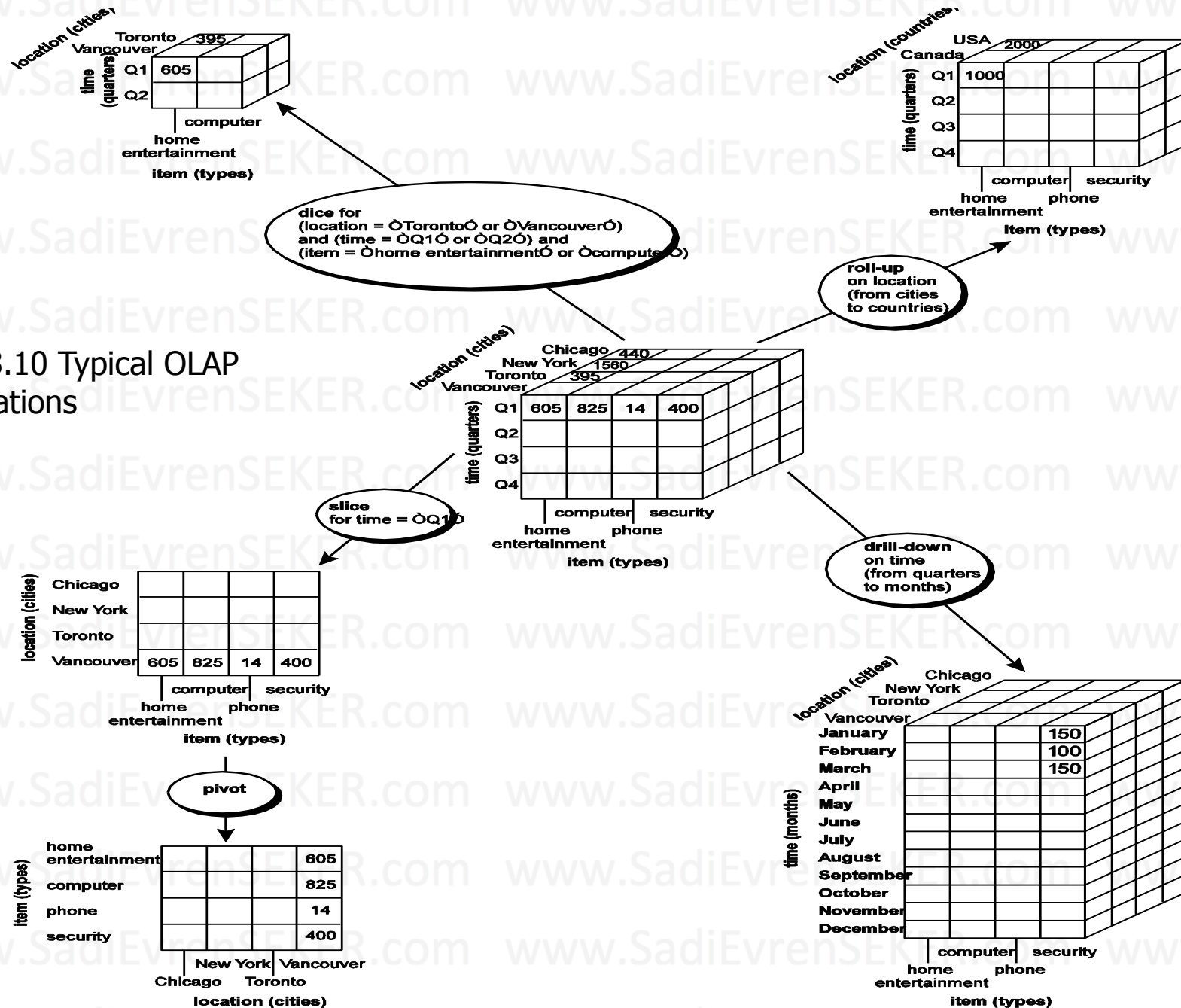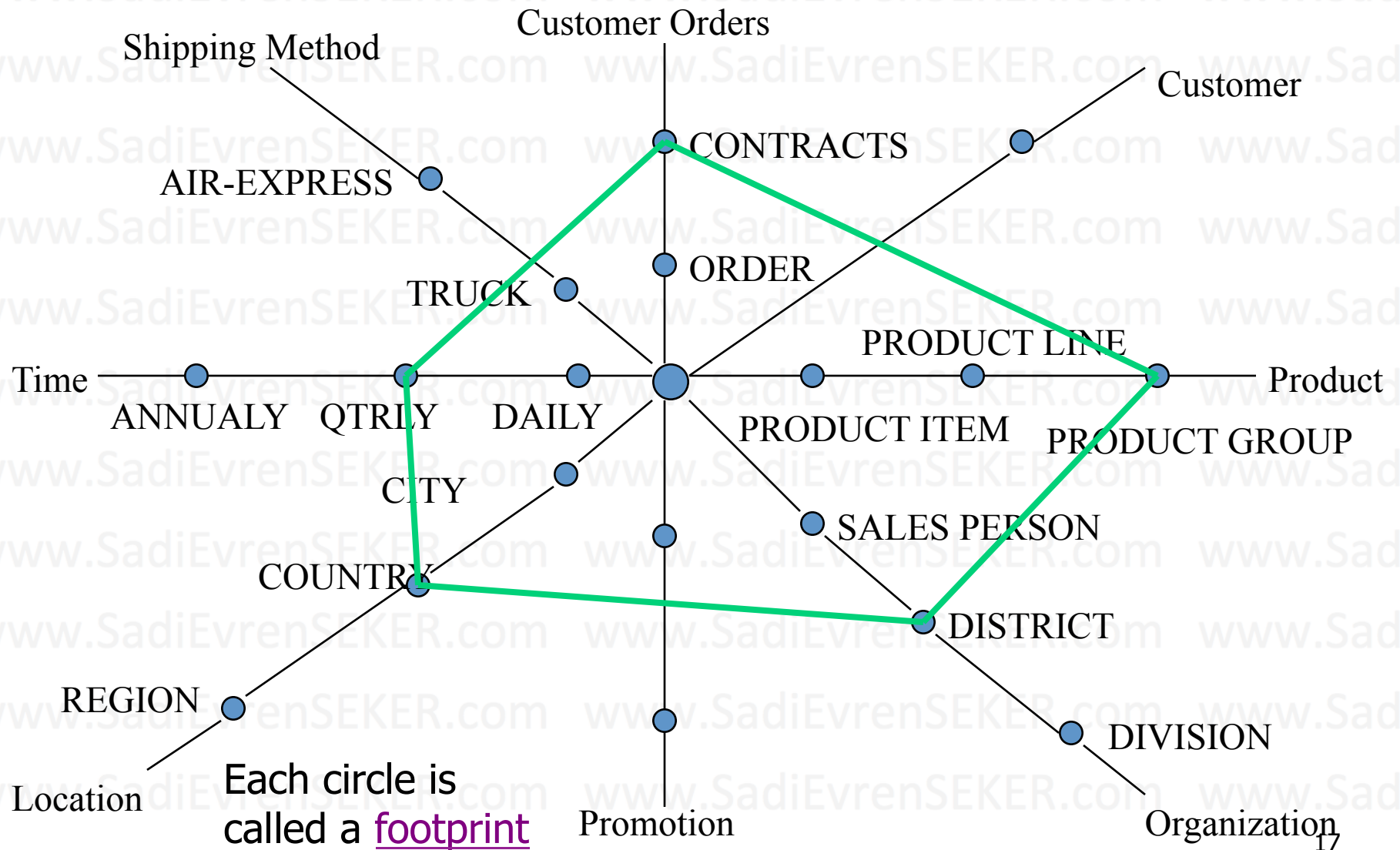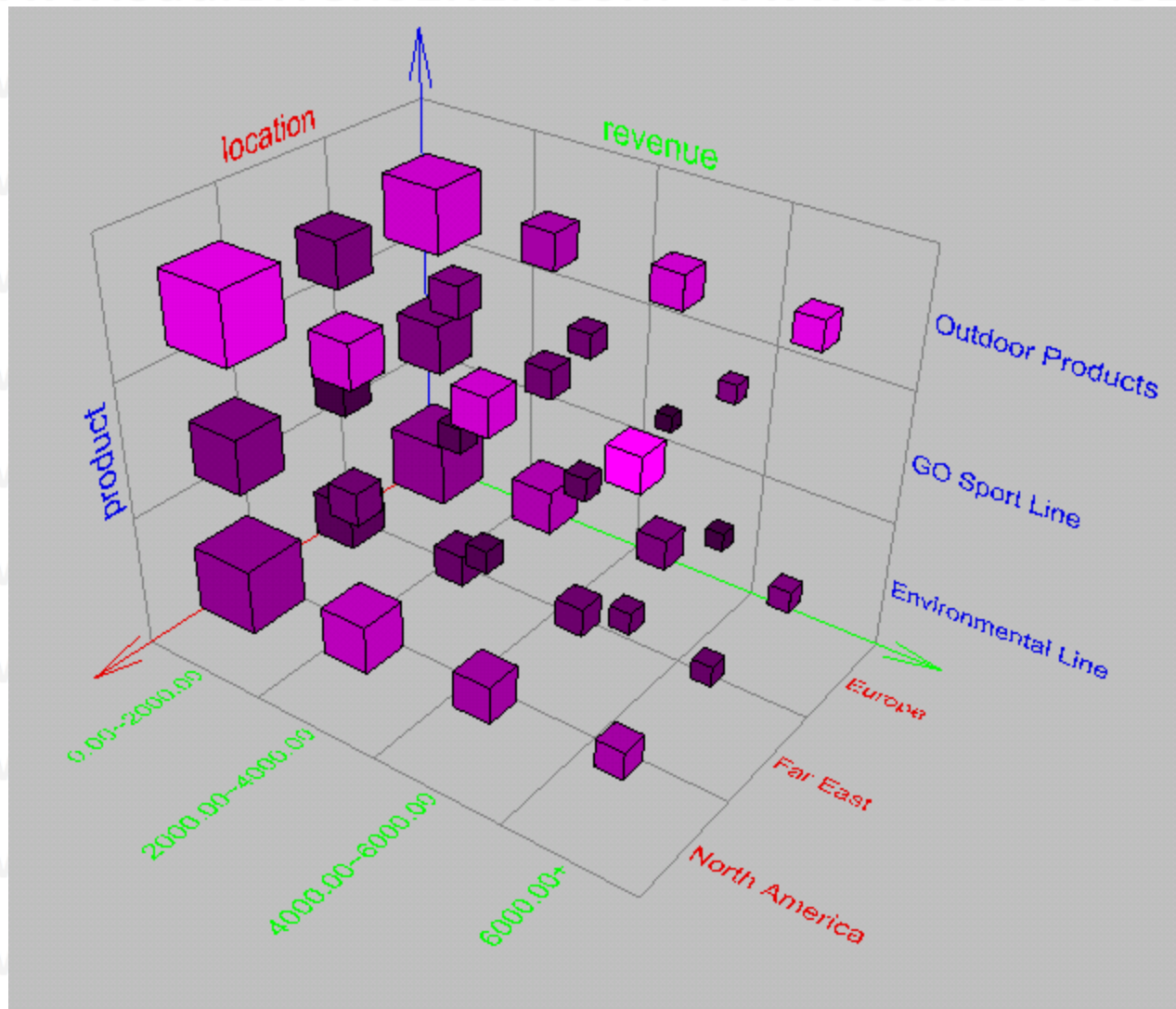
15

Fig. 3.10 Typical OLAP Operations

**Top left cube (dice result):**
- location (cities): Toronto 395, Vancouver
- time (quarters): Q1 605, Q2
- item (types): home entertainment, computer

**dice for**
(location = "Toronto" or "Vancouver")
and (time = "Q1" or "Q2") and
(item = "home entertainment" or "computer")

**Top right cube (roll-up result):**
- location (countries): USA 2000, Canada
- time (quarters): Q1 1000, Q2, Q3, Q4
- item (types): home entertainment, computer, phone, security

**roll-up**
on location
(from cities
to countries)

**Center cube:**
- location (cities): Chicago 440, New York 1560, Toronto 395, Vancouver
- time (quarters): Q1 605 825 14 400, Q2, Q3, Q4
- item (types): home entertainment, computer, phone, security

**slice**
for time = "Q1"

**Slice result table:**

| location (cities) | home entertainment | computer | phone | security |
|---|---|---|---|---|
| Chicago | | | | |
| New York | | | | |
| Toronto | | | | |
| Vancouver | 605 | 825 | 14 | 400 |

item (types)

**pivot**

**Pivot result table:**

| item (types) | Chicago | New York | Toronto | Vancouver |
|---|---|---|---|---|
| home entertainment | | | | 605 |
| computer | | | | 825 |
| phone | | | | 14 |
| security | | | | 400 |

location (cities)

**drill-down**
on time
(from quarters
to months)

**Drill-down cube:**
- location (cities): Chicago, New York, Toronto, Vancouver
- time (months): January 150, February 100, March 150, April, May, June, July, August, September, October, November, December
- item (types): home entertainment, computer, phone, security

16

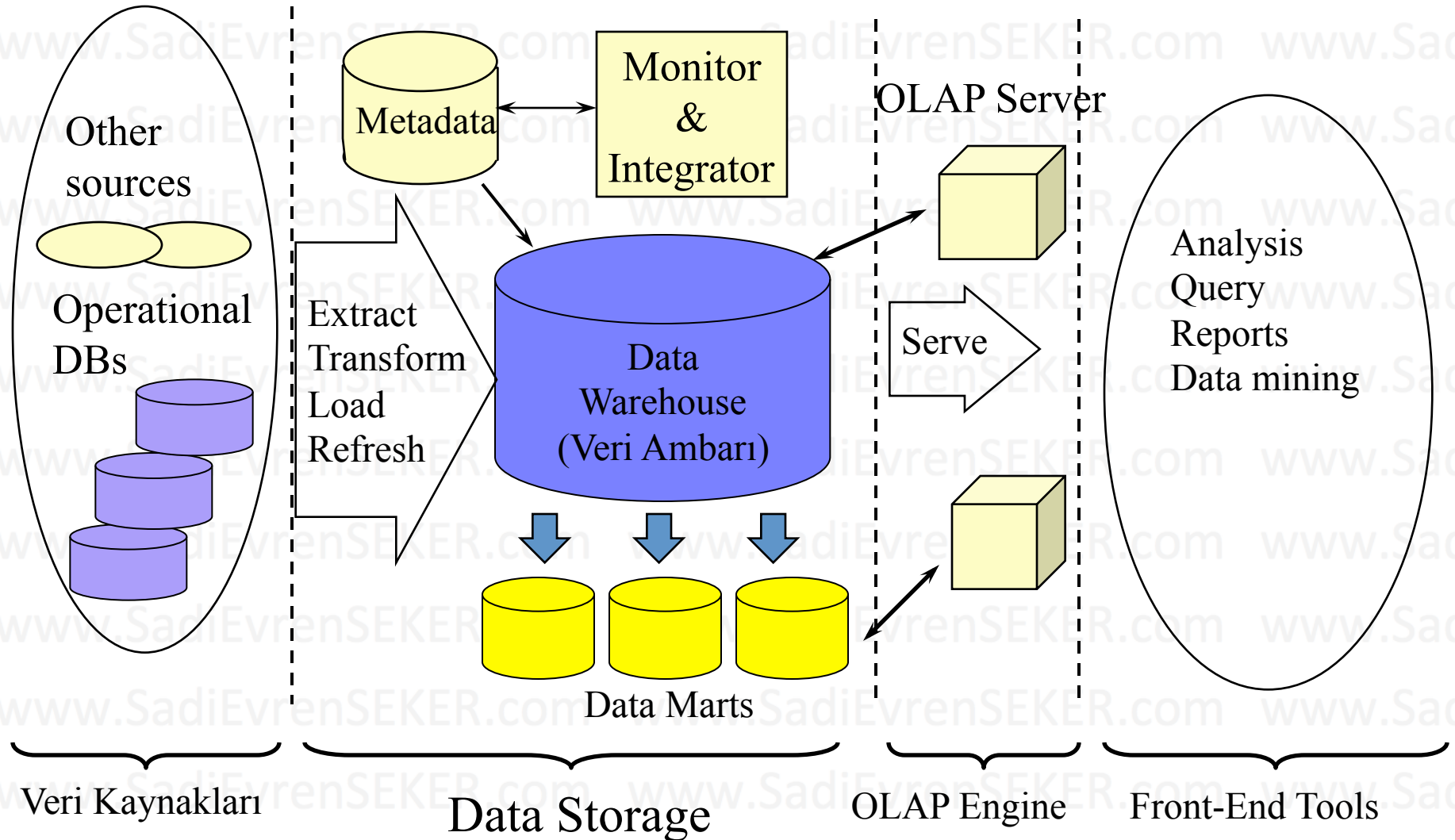# A Star-Net Query Model

# Browsing a Data Cube



ation
apabilities
tive manipulation

18

# Chapter 4: Data Warehousing and On-line Analytical Processing

- Data Warehouse: Basic Concepts

- Data Warehouse Modeling: Data Cube and OLAP

- Data Warehouse Design and Usage

- Data Warehouse Implementation

- Data Generalization by Attribute-Oriented Induction

- Summary
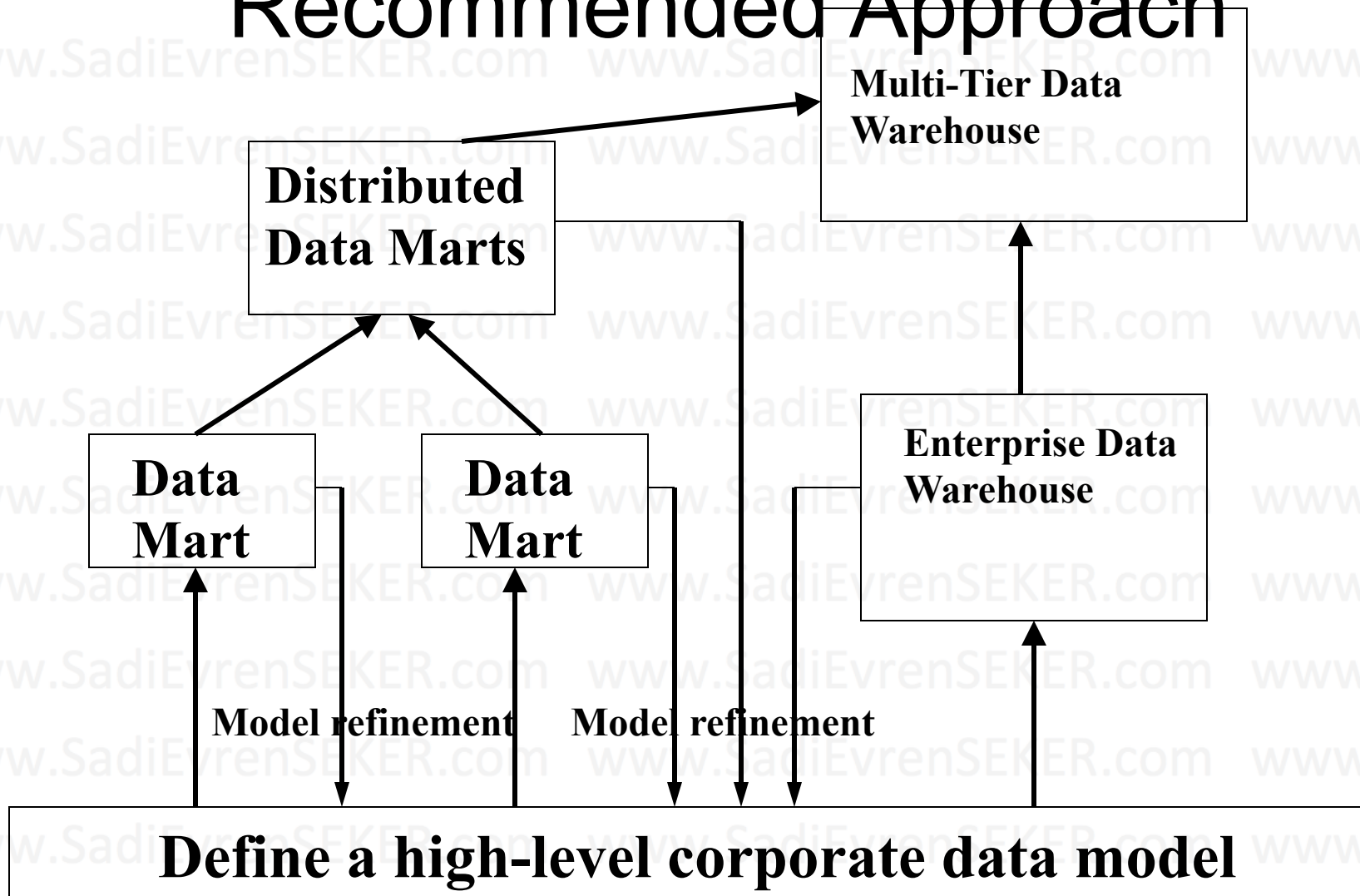
# Data Warehouse: A Multi-Tiered Architecture

Other sources

Metadata

Monitor & Integrator

OLAP Server

Operational DBs

Extract
Transform
Load
Refresh

Data Warehouse (Veri Ambarı)

Serve

Analysis
Query
Reports
Data mining

Data Marts

Veri Kaynakları

Data Storage

OLAP Engine

Front-End Tools

# Design of Data Warehouse: A Business Analysis Framework

- Four views regarding the design of a data warehouse
  - Top-down view
    - allows selection of the relevant information necessary for the data warehouse
  - Data source view
    - exposes the information being captured, stored, and managed by operational systems
  - Data warehouse view
    - consists of fact tables and dimension tables
  - Business query view
    - sees the perspectives of data in the warehouse from the view of end-user

# Data Warehouse Design Process

- **Top-down, bottom-up approaches or a combination** of both
  - <u>Top-down</u>: Starts with overall design and planning (mature)
  - <u>Bottom-up</u>: Starts with experiments and prototypes (rapid)
- **From software engineering point of view**
  - <u>Waterfal</u>l: structured and systematic analysis at each step before proceeding to the next
  - <u>Spiral</u>:  rapid generation of increasingly functional systems, short turn around time, quick turn around
- **Typical data warehouse design process**
  - Choose a business process to model, e.g., orders, invoices, etc.
  - Choose the *grain* (*atomic level of data*) of the business process
  - Choose the dimensions that will apply to each fact table record
  - Choose the measure that will populate each fact table record

# Data Warehouse Development: A Recommended Approach

Multi-Tier Data Warehouse

Distributed Data Marts

Data Mart

Data Mart

Enterprise Data Warehouse

Model refinement    Model refinement

**Define a high-level corporate data model**

# 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

# From On-Line Analytical Processing (OLAP) to On Line 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

# Chapter 4: Data Warehousing and On-line Analytical Processing

- Data Warehouse: Basic Concepts

- Data Warehouse Modeling: Data Cube and OLAP

- Data Warehouse Design and Usage

- Data Warehouse Implementation

- Data Generalization by Attribute-Oriented Induction

- Summary

# 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 <u>every</u> (cuboid) (**full materialization**), <u>none</u> (**no materialization**), or <u>some</u> (**partial materialization**)
  - Selection of which cuboids to materialize
    - Based on size, sharing, access frequency, etc.

# The "Compute Cube" Operator

- Cube definition and computation in DMQL

  define cube sales [item, city, year]: sum (sales_in_dollars)

  compute cube sales

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

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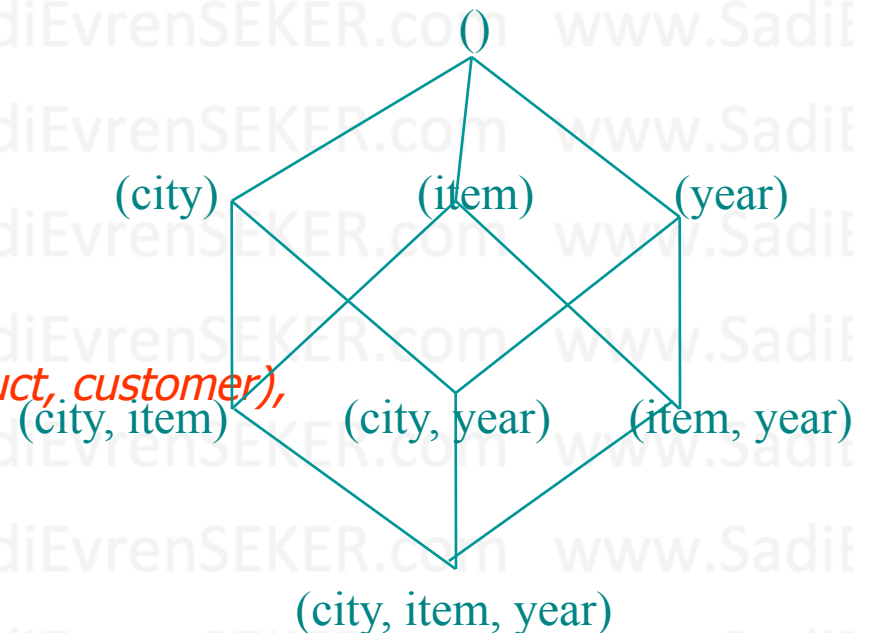
  FROM SALES

  CUBE BY item, city, year

- Need compute the following Group-Bys
  (date, product, customer),
  (date,product),(date, customer), (product, customer),
  (date), (product), (customer)
  ()

()

(city)          (item)          (year)

(city, item)    (city, year)    (item, year)

(city, item, year)

# Indexing OLAP Data: **Bitmap Index**

- Index on a particular column
- Each value in the column has a bit vector: bit-op is fast
- The length of the bit vector: # of records in the base table
- The $i$-th bit is set if the $i$-th row of the base table has the value for the indexed column
- not suitable for high cardinality domains
- A recent bit compression technique, Word-Aligned Hybrid (WAH), makes it work for high cardinality domain as well [Wu, et al. TODS' 06]

### Base table

| Cust | Region | Type |
|------|--------|--------|
| C1 | Asia | Retail |
| C2 | Europe | Dealer |
| C3 | Asia | Dealer |
| C4 | America | Retail |
| C5 | Europe | Dealer |

### Index on Region

| RecID | Asia | Europe | America |
|-------|------|--------|---------|
| 1 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 |
| 3 | 1 | 0 | 0 |
| 4 | 0 | 0 | 1 |
| 5 | 0 | 1 | 0 |

### Index on Type

| RecID | Retail | Dealer |
|-------|--------|--------|
| 1 | 1 | 0 |
| 2 | 0 | 1 |
| 3 | 0 | 1 |
| 4 | 1 | 0 |
| 5 | 0 | 1 |

29

# Indexing OLAP Data: **Join Indices**

- Join index: JI(R-id, S-id) where R (R-id, …) ⊳⊲ S (S-id, …)
- Traditional indices map the values to a list of record ids
  - It materializes relational join in JI file and speeds up relational join
- In data warehouses, join index relates the values of the dimensions of a start schema to rows in the fact table.
  - E.g. fact table: *Sales* and two dimensions *city* and *product*
    - A join index on *city* maintains for each distinct city a list of R-IDs of the tuples recording the Sales in the city
  - Join indices can span multiple dimensions

location

Main Street

sales

T57

T238

T459

T884

item

Sony-TV

# Efficient Processing 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 which materialized cuboid(s)** should be selected for OLAP op.
  - Let the query to be processed be on {*brand, province_or_state*} with the condition "*year = 2004*", and there are 4 materialized cuboids available:

    1) {*year, item_name, city*}

    2) {*year, brand, country*}

    3) {*year, brand, province_or_state*}

    4) {*item_name, province_or_state*}  where *year = 2004*

    Which should be selected to process the query?
- Explore indexing structures and compressed vs. dense array structs in MOLAP

# OLAP Server Architectures

- **Relational OLAP (ROLAP)**
  - Use relational or extended-relational DBMS to store and manage warehouse data and OLAP middle ware
  - Include optimization of DBMS backend, implementation of aggregation navigation logic, and additional tools and services
  - Greater scalability
- **Multidimensional OLAP (MOLAP)**
  - Sparse array-based multidimensional storage engine
  - Fast indexing to pre-computed summarized data
- **Hybrid OLAP (HOLAP)** (e.g., Microsoft SQLServer)
  - Flexibility, e.g., low level: relational, high-level: array
- **Specialized SQL servers** (e.g., Redbricks)
  - Specialized support for SQL queries over star/snowflake schemas

# Chapter 4: Data Warehousing and On-line Analytical Processing

- Data Warehouse: Basic Concepts

- Data Warehouse Modeling: Data Cube and OLAP

- Data Warehouse Design and Usage

- Data Warehouse Implementation

- Data Generalization by Attribute-Oriented Induction

- Summary

# Attribute-Oriented Induction

- Proposed in 1989 (KDD '89 workshop)
- Not confined to categorical data nor particular measures
- How it is done?
  - Collect the task-relevant data (*initial relation*) using a relational database query
  - Perform generalization by <u>attribute removal</u> or <u>attribute generalization</u>
  - Apply aggregation by merging identical, generalized tuples and accumulating their respective counts
  - Interaction with users for knowledge presentation

# Attribute-Oriented Induction: An Example

Example:  Describe general characteristics of graduate students in the University database

- Step 1. Fetch relevant set of data using an SQL statement, e.g.,

   **Select** * (i.e., name, gender, major, birth_place, birth_date, residence, phone#, gpa)

   **from** student

   **where**  student_status in {"Msc", "MBA", "PhD" }

- Step 2. Perform attribute-oriented induction
- Step 3. Present results in generalized relation, cross-tab, or rule forms

# Class Characterization: An Example

**Initial Relation**

| Name | Gender | Major | Birth-Place | Birth_date | Residence | Phone # | GPA |
|---|---|---|---|---|---|---|---|
| Jim Woodman | M | CS | Vancouver,BC, Canada | 8-12-76 | 3511 Main St., Richmond | 687-4598 | 3.67 |
| Scott Lachance | M | CS | Montreal, Que, Canada | 28-7-75 | 345 1st Ave., Richmond | 253-9106 | 3.70 |
| Laura Lee | F | Physics | Seattle, WA, USA | 25-8-70 | 125 Austin Ave., Burnaby | 420-5232 | 3.83 |
| … | … | … | … | … | … | … | … |
| **Removed** | **Retained** | **Sci,Eng, Bus** | **Country** | **Age range** | **City** | **Removed** | **Excl, VG,..** |

**Prime Generalized Relation**

| Gender | Major | Birth_region | Age_range | Residence | GPA | Count |
|---|---|---|---|---|---|---|
| M | Science | Canada | 20-25 | Richmond | Very-good | 16 |
| F | Science | Foreign | 25-30 | Burnaby | Excellent | 22 |
| … | … | … | … | … | … | … |

| Gender \ Birth_Region | Canada | Foreign | Total |
|---|---|---|---|
| M | 16 | 14 | 30 |
| F | 10 | 22 | 32 |
| Total | 26 | 36 | 62 |

36

# Basic Principles of Attribute-Oriented Induction

- Data focusing: task-relevant data, including dimensions, and the result is the *initial relation*

- Attribute-removal: remove attribute *A* if there is a large set of distinct values for *A* but (1) there is no generalization operator on *A*, or (2) *A*'s higher level concepts are expressed in terms of other attributes

- Attribute-generalization: If there is a large set of distinct values for *A*, and there exists a set of generalization operators on *A*, then select an operator and generalize *A*

- Attribute-threshold control: typical 2-8, specified/default

- Generalized relation threshold control: control the final relation/ rule size

37

# Attribute-Oriented Induction: Basic Algorithm

- **InitialRel**: Query processing of task-relevant data, deriving the *initial relation*.

- **PreGen**: Based on the analysis of the number of distinct values in each attribute, determine generalization plan for each attribute: removal? or how high to generalize?

- **PrimeGen**: Based on the PreGen plan, perform generalization to the right level to derive a "prime generalized relation", accumulating the counts.

- **Presentation**: User interaction: (1) adjust levels by drilling, (2) pivoting, (3) mapping into rules, cross tabs, visualization presentations.

# Presentation of Generalized Results

- <u>Generalized relation</u>:
  - Relations where some or all attributes are generalized, with counts or other aggregation values accumulated.

- <u>Cross tabulation</u>:
  - Mapping results into cross tabulation form (similar to contingency tables).

  - <u>Visualization techniques</u>:

  - Pie charts, bar charts, curves, cubes, and other visual forms.

- <u>Quantitative characteristic rules</u>:

  - Mapping generalized result into characteristic rules with quantitative information associated with it, e.g.

$$grad(x) \wedge male(x) \Rightarrow$$
$$birth\_region(x) = "Canada"[t:53\%] \vee birth\_region(x) = "foreign"[t:47\%].$$

# Mining Class Comparisons

- <u>Comparison:</u> Comparing two or more classes
- <u>Method:</u>
  - Partition the set of relevant data into the target class and the contrasting class(es)
  - Generalize both classes to the same high level concepts
  - Compare tuples with the same high level descriptions
  - Present for every tuple its description and two measures
    - support - distribution within single class
    - comparison - distribution between classes
  - Highlight the tuples with strong discriminant features
- <u>Relevance Analysis:</u>
  - Find attributes (features) which best distinguish different classes

# Concept Description vs. Cube-Based OLAP

- **Similarity**:
  - Data generalization
  - Presentation of data summarization at multiple levels of abstraction
  - Interactive drilling, pivoting, slicing and dicing
- **Differences**:
  - OLAP has systematic preprocessing, query independent, and can drill down to rather low level
  - AOI has automated desired level allocation, and may perform dimension relevance analysis/ranking when there are many relevant dimensions
  - AOI works on the data which are not in relational forms

# Chapter 4: Data Warehousing and On-line Analytical Processing

- Data Warehouse: Basic Concepts

- Data Warehouse Modeling: Data Cube and OLAP

- Data Warehouse Design and Usage

- Data Warehouse Implementation

- Data Generalization by Attribute-Oriented Induction

- Summary

# Summary

- Data warehousing: A multi-dimensional model of a data warehouse
  - A data cube consists of *dimensions* & *measures*
  - Star schema, snowflake schema, fact constellations
  - OLAP operations: drilling, rolling, slicing, dicing and pivoting
- Data Warehouse Architecture, Design, and Usage
  - Multi-tiered architecture
  - Business analysis design framework
  - Information processing, analytical processing, data mining, OLAM (Online Analytical Mining)
- Implementation: Efficient computation of data cubes
  - Partial vs. full vs. no materialization
  - Indexing OALP data: Bitmap index and join index
  - OLAP query processing
  - OLAP servers: ROLAP, MOLAP, HOLAP
- Data generalization: Attribute-oriented induction

# References (I)

- S. Agarwal, R. Agrawal, P. M. Deshpande, A. Gupta, J. F. Naughton, R. Ramakrishnan, and S. Sarawagi.   On the computation of multidimensional aggregates.  VLDB'96

- D. Agrawal, A. E. Abbadi, A. Singh, and T. Yurek.  Efficient view maintenance in data warehouses. SIGMOD'97

- R. Agrawal, A. Gupta, and S. Sarawagi. Modeling multidimensional databases.  ICDE'97

- S. Chaudhuri and U. Dayal.  An overview of data warehousing and OLAP technology.  ACM SIGMOD Record, 26:65-74, 1997

- E. F. Codd, S. B. Codd, and C. T. Salley. Beyond decision support. Computer World, 27, July 1993.

- J. Gray, et al.  Data cube: A relational aggregation operator generalizing group-by, cross-tab and sub-totals.  Data Mining and Knowledge Discovery, 1:29-54, 1997.

- A. Gupta and I. S. Mumick.  Materialized Views: Techniques, Implementations, and Applications. MIT Press, 1999.

- J. Han.  Towards on-line analytical mining in large databases. *ACM SIGMOD Record*, 27:97-107, 1998.

- V. Harinarayan, A. Rajaraman, and J. D. Ullman. Implementing data cubes efficiently. SIGMOD'96

- J. Hellerstein, P. Haas, and H. Wang. Online aggregation. SIGMOD'97

# References (II)

- C. Imhoff, N. Galemmo, and J. G. Geiger. Mastering Data Warehouse Design: Relational and Dimensional Techniques. John Wiley, 2003

- W. H. Inmon. Building the Data Warehouse. John Wiley, 1996

- R. Kimball and M. Ross.  The Data Warehouse Toolkit: The Complete Guide to Dimensional Modeling. 2ed. John Wiley, 2002

- P. O'Neil and G. Graefe. Multi-table joins through bitmapped join indices. *SIGMOD Record*, 24:8–11, Sept. 1995.

- P. O'Neil and D. Quass. Improved query performance with variant indexes. SIGMOD'97

- Microsoft. OLEDB for OLAP programmer's reference version 1.0. In http://www.microsoft.com/data/oledb/olap, 1998

- S. Sarawagi and M. Stonebraker.  Efficient organization of large multidimensional arrays. ICDE'94

- A. Shoshani.  OLAP and statistical databases: Similarities and differences. PODS'00.

- D. Srivastava, S. Dar, H. V. Jagadish, and A. V. Levy. Answering queries with aggregation using views. *VLDB'96*

- P. Valduriez. Join indices.  ACM Trans. Database Systems, 12:218-246, 1987.

- J. Widom.  Research problems in data warehousing.  CIKM'95

- K. Wu, E. Otoo, and A. Shoshani, Optimal Bitmap Indices with Efficient Compression, ACM Trans. on Database Systems (TODS), 31(1): 1-38, 2006