

Data stream فصل اول ۹۳، ۹۴، ۹۵

data warehousing
کتابخانه

Chapter 3: Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- Data warehouse architecture طراحی آن
- From data warehousing to data mining

ماتریز داده‌ها و استخراج داده‌ها از آن

December 1, 2014

Data Mining: Concepts and Techniques

1

What is 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

December 1, 2014

Data Mining: Concepts and Techniques

2

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

December 1, 2014

Data Mining: Concepts and Techniques

3

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.

December 1, 2014

Data Mining: Concepts and Techniques

4

→ Random Sampling تصویف هدف آنست از داده ها بخش کوچکی را به طریقی تصادفی برگزینیم

داده ها چند است. با فرض اینکه داده ها به هم وابسته اند و به هم وابسته اند.

Reservoir Sampling به این احوال است.

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"

December 1, 2014

Data Mining: Concepts and Techniques

5

Data Warehouse—Nonvolatile

- 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*

December 1, 2014

Data Mining: Concepts and Techniques

6

حافظه استفاده شده متفاوت است و جواب دهنده متن متغیر است

کوئی حاد (ارزاد است) - متن در آن کاشی به انجام داد

challenge (چالنج) یعنی آنکه برای هر داده باید روشی بکار برده می شود
تفاوت مابین این دو است - مثل ترافیک شبکه که مقدار زیادی داده

باید در طول زمان بر روی سیستم که ساده است

روش کار استفاده شده تغییر است - وقت باقی است

روش کار استفاده شده برابر است با Random Sampling است

داده بیست است - disjoint یعنی در بین غرض برای گرفتن حجم داده کم شود

حسب نام دیگر روش داده که برای آن تعداد و زمان را حدس بزنم

ساده داده را از یک نمونه رد کنم و جابجایی آنها - threshold - کنیم - بخوانم داده را

و نوری احف کنیم - برابر سرع با کار الگوریتم Randomize (تغییر داده)

Data Warehouse vs. Heterogeneous DBMS

- Traditional heterogeneous DB integration: A query driven approach
 - Build wrappers/mediators on top of heterogeneous databases
 - 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

December 1, 2014

Data Mining: Concepts and Techniques

7

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

December 1, 2014

Data Mining: Concepts and Techniques

8

معاری : DSMS

یکدیگر داده داریم که در طول زمان تولید می شوند مثل این باید چندین استرس کنیم
چون که استرس است و بر آورد است و بخش بندی و بخش بندی و غیره
در این چندین که داریم می توانیم اطلاعات را به یک بخش دیگر که داریم انتقال دهیم
خبره کنیم به آنرا (مثل مصوفین)

چندین challenge برای Stream Data Processing داریم :
داده های غیر منتهی (کار نیست به شود) داده های بی نهایت می آیند و می بینیم
مندی باید کار کنیم به این دلیل که داده ها به سرعت می آیند و می بینیم
کوئری ها می توانند اسفند و می توانند اسفند و کار ما محدود می شود
تقسیم کردن از قبل تعریف شده یا (ad-hoc) باشد یعنی همان دسی (میکروس)
حفظ خبری را بر آورد می کنیم و یک بر آورد خوب به دست می آوریم و به دست می آوریم

OLTP vs. OLAP

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

December 1, 2014

Data Mining: Concepts and Techniques

9

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
- Note: There are more and more systems which perform OLAP analysis directly on relational databases

December 1, 2014

Data Mining: Concepts and Techniques

10

1
تفاوت DBMS با DSMS (Data stream management system)
کار با DBMS سخت تر است. کار کردن با داده های جاری به زبان ساده
است در برابر آن.

Random access: داده ها را می توان به صورت تصادفی جمع آوری کرد و به راحتی دسترسی
آنها ممکن است بر روی سیستم ها با DSMS چون به زبان ساده است و به راحتی
در دسترس است برای دسترسی به آنها به راحتی می توان

در DBMS با داده های جاری کار می کنند اما در DSMS داده ها به صورت جریان می آید

دلیل در DBMS نیاز به Real time بودن ندارد اما در DSMS real time

است. update کردن در DBMS به راحتی می توان اما در DSMS به راحتی نمی توان
درست کردن به راحتی می توان. در DBMS داده ها را می توان به راحتی در DSMS

در DBMS داده ها را می توان به راحتی در DSMS به راحتی می توان

Chapter 3: Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- Data warehouse architecture
- From data warehousing to data mining

December 1, 2014

Data Mining: Concepts and Techniques

11

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

December 1, 2014

Data Mining: Concepts and Techniques

12

DBMS نسبی داده‌ها نیست و ضبط شده

تفاوت اصلی جریان داده در DBMS هم است

مسئله‌های جریان داده به سه دسته است: مثل مصرف‌کننده، کپی، کپی برداشتن

مصرف‌کننده: مصرف‌کننده real time دیده می‌شوند یا برای داده‌های ذخیره‌سازی
در واحدهای مختلف است.

عملیات این امور برای مرتب‌سازی داده‌ها در نظر می‌گیریم

مسئله 5: داده در مورد میزان برآورد مصرف‌کننده، اگر می‌بینیم و خود داده‌ها را

در زیر می‌بینیم و تغییرات می‌شود مثل متوسط آنها را ذخیره می‌کنیم

مثال: داده‌ها را با سرعت 1 تغییر می‌دهند multiple چندین

باز می‌بینیم. مثل میزان برآورد مصرف‌کننده، مثلاً فاصله از زمان

کردیم که هیچ‌کس مصرف‌کننده را نمی‌بیند (خود می‌بیند)

چرا استفاده می‌کنیم و در مورد زمان تفاوت مصرف‌کننده، برآورد

وقت تعیین، مصرف‌کننده، ماشین‌ها در دسترس برای مصرف‌کننده

برای مصرف‌کننده، حجم داده‌ها در زمان داده‌ها می‌بینیم

Data Warehouse Design Process

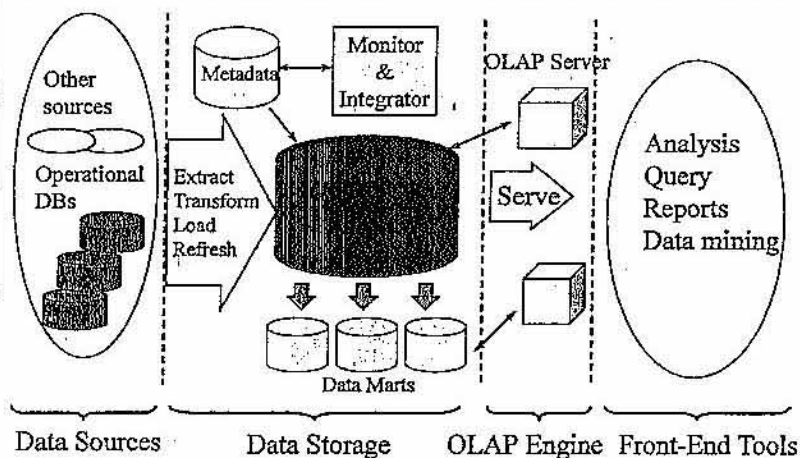
- Top-down, bottom-up approaches or a combination of both
 - Top-down: Starts with overall design and planning (mature)
 - Bottom-up: Starts with experiments and prototypes (rapid)
- From software engineering point of view
 - Waterfall: structured and systematic analysis at each step before proceeding to the next
 - Spiral: 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

December 1, 2014

Data Mining: Concepts and Techniques

13

Data Warehouse: A Multi-Tiered Architecture



December 1, 2014

Data Mining: Concepts and Techniques

14

(ساخته شده)

Three Data Warehouse Models

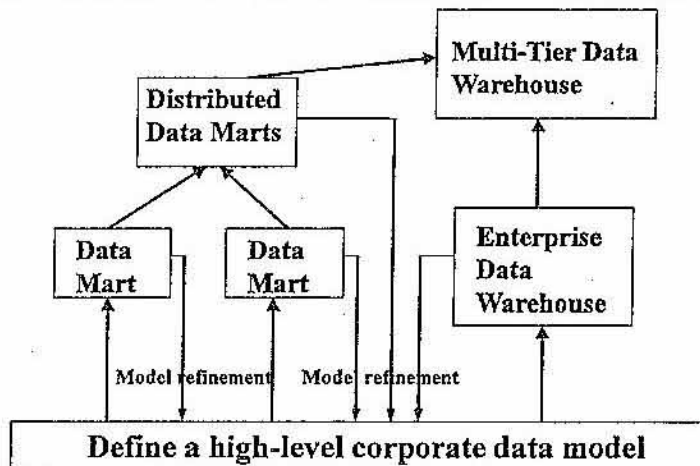
- Enterprise warehouse
 - collects all of the information about subjects spanning the entire organization
- Data Mart
 - a subset of corporate-wide data that is of value to a specific groups of users. Its scope is confined to specific, selected groups, such as marketing data mart
 - Independent vs. dependent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

December 1, 2014

Data Mining: Concepts and Techniques

15

Data Warehouse Development: A Recommended Approach



December 1, 2014

Data Mining: Concepts and Techniques

16

صبرتان داده: مثل محاسبه نرخ تغییرات کار، محاسبه جریان آن (طبیعتاً) مقدار به کار

کشور در روز ماه و هفته می دهد

چنانچه داده می تواند بصورت پیوسته تولید شود و صرفاً بر یک خط به عنوان آنرا اندازه گیری کرد

در داده های stream دنبال change point detection می بینیم. (مقاطع)

داده های stream همیشه به یکدیگر وابسته اند و نیاز است

در صورت مشاهده تغییر در مقدار

سیستم تولید می شوند

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

December 1, 2014

Data Mining: Concepts and Techniques

17

Metadata Repository

- Meta data is the data defining warehouse objects. It stores:
- Description of the structure of the data 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
 - warehouse schema, view and derived data definitions
- Business data
 - business terms and definitions, ownership of data, charging policies

December 1, 2014

Data Mining: Concepts and Techniques

18

داده های طولی Stream (است) داده های درختی
جمع آوری می شوند مثل فشار خون یک بیمار بصورت روزانه. شبکه بارش داده های طولی
داده های Stream سری های زمانی اینجاست که اینجای داده مستقل از هم نیستند
ارتباط فضایی است که در فضا است مثل داده های درختی که درختان در یک منطقه هستند
و این قوی تر است. داده های درختی که در طول زمان هستند و به هم وابسته هستند
داده های زمانی برای هم نیست Sequential Data می باشد.

اما در هم حرارت شب و روز با هم متفاوت است و این هم نیست
فصل به فصل از یک فصل به فصل تر باشد مثلاً که با فصل کار می کنیم بعد فصل بعدی باشد
مثل درختان که یک درخت یک ساله یک ساله ۱۰۰۰۰ (یک میلیون بار) است
اگر در یک زمان مختلف داده ها را اندازه گیری شود به آن Spatial-Temporal می گویند
که این حالت را به عنوان یک سری زمانی با فضا می نامند و این حالت را به عنوان یک سری زمانی
یک فصل زمانی در یک فصل است و این حالت را به عنوان یک سری زمانی با فضا می نامند
می شود. در این حالت با یک سری زمانی با فضا می نامند و این حالت را به عنوان یک سری زمانی

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 SQL Server)
 - Flexibility, e.g., low level: relational, high-level: array
- Specialized SQL servers (e.g., Redbricks)
 - Specialized support for SQL queries over star/snowflake schemas

December 1, 2014

Data Mining: Concepts and Techniques

19

Chapter 3: Data Warehousing and OLAP Technology: An Overview

- What is a data warehouse?
- Data warehouse architecture
- From data warehousing to data mining

December 1, 2014

Data Mining: Concepts and Techniques

20

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

December 1, 2014

Data Mining: Concepts and Techniques

21

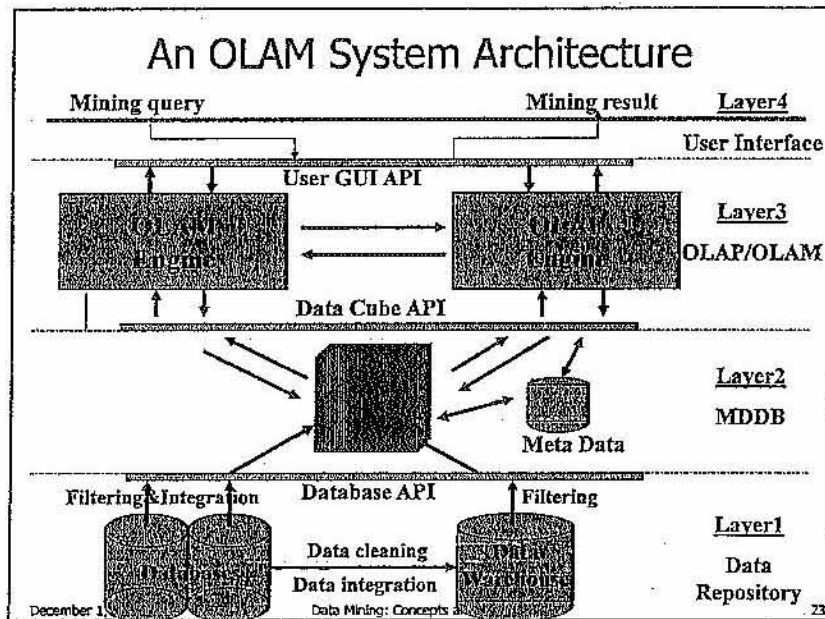
From On-Line Analytical Processing (OLAP) to On Line Analytical Mining (OLAM)

- 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

December 1, 2014

Data Mining: Concepts and Techniques

22



Chapter 3: Data Warehousing and OLAP Technology: An Overview

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

Summary: Data Warehouse and OLAP Technology

- Why data warehousing?
- Data warehouse architecture
- From OLAP to OLAM (on-line analytical mining)

December 1, 2014

Data Mining: Concepts and Techniques

25

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*

December 1, 2014

Data Mining: Concepts and Techniques

26

References (II)

- C. Imhoff, N. Gallempo, 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 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
- A. Shoshani. OLAP and statistical databases: Similarities and differences. PODS'00.
- S. Sarawagi and M. Stonebraker. Efficient organization of large multidimensional arrays. ICDE'94
- OLAP council. MDAPI specification version 2.0. In <http://www.olapcouncil.org/research/aply.htm>, 1998
- E. Thomsen. OLAP Solutions: Building Multidimensional Information Systems. John Wiley, 1997
- P. Valduriez. Join indices. ACM Trans. Database Systems, 12:218-246, 1987.
- J. Widom. Research problems in data warehousing. CIKM'95.

December 1, 2014

Data Mining: Concepts and Techniques

27

Data Mining: Concepts and Techniques

— Chapter 8 —

8.1. Mining data streams

Jiawei Han and Micheline Kamber
Department of Computer Science
University of Illinois at Urbana-Champaign

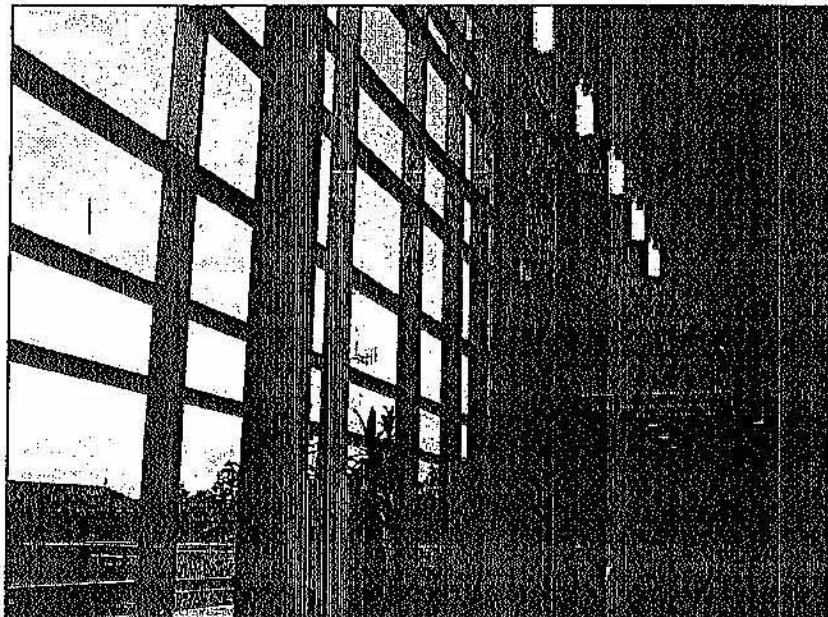
www.cs.uiuc.edu/~hanj

©2006 Jiawei Han and Micheline Kamber. All rights reserved.

December 1, 2014

Data Mining: Concepts and Techniques

1



Data and Information Systems (DAIS:) Course Structures at CS/UIUC

- Three streams: Database, data mining and text information systems
- Database Systems:
 - Database mgmt systems (CS411: Fall and Spring)
 - Advanced database systems (CS511: Fall)
 - Web information systems (Kevin Chang)
 - Information integration (An-Hai Doan)
- Data mining
 - Intro. to data mining (CS412: Han—Fall)
 - Data mining: Principles and algorithms (CS512: Han—Spring)
 - Seminar: Advanced Topics in Data mining (CS591 Han—Fall and Spring)
- Text information systems and Bioinformatics
 - Text information system (CS410 Zhai)
 - Introduction to Bioinformatics (CS598 Sinha, CS498 Zhai)

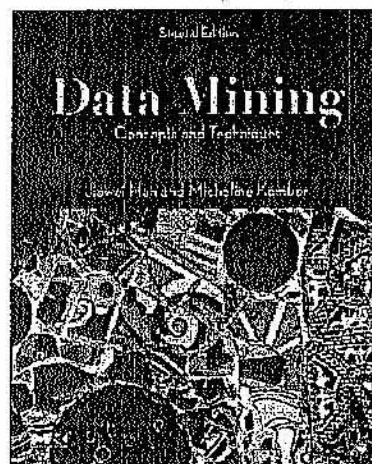
December 1, 2014

Data Mining: Concepts and Techniques

3

Data Mining: Concepts and Techniques, 2ed. 2006

- Seven chapters (Chapters 1-7) are covered in the Fall semester
- Four chapters (Chapters 8-11) are covered in the Spring semester



December 1, 2014

Data Mining: Concepts and Techniques

4

Coverage of CS412@UIUC (Intro. to Data Warehousing and Data Mining)

1. Introduction
2. Data Preprocessing
3. Data Warehouse and OLAP Technology: An Introduction
4. Advanced Data Cube Technology and Data Generalization
5. Mining Frequent Patterns, Association and Correlations
6. Classification and Prediction
7. Cluster Analysis

December 1, 2014

Data Mining: Concepts and Techniques

5

Coverage of CS512@UIUC (Data Mining: Principles and Algorithms)

- | | |
|---|--|
| <ol style="list-style-type: none"> 8. Mining stream, time-series, and sequence data <ul style="list-style-type: none"> ■ Mining data streams ■ Mining time-series data ■ Mining sequence patterns in transactional databases ■ Mining sequence patterns in biological data 9. Graph mining, social network analysis, and multi-relational data mining <ul style="list-style-type: none"> ■ Graph mining ■ Social network analysis ■ Multi-relational data mining | <ol style="list-style-type: none"> 10. Mining Object, Spatial, Multimedia, Text and Web data <ul style="list-style-type: none"> ■ Mining object data ■ Spatial and spatiotemporal data mining ■ Multimedia data mining ■ Text mining ■ Web mining 11. Applications and trends of data mining <ul style="list-style-type: none"> ■ Data mining applications ■ Data mining products and research prototypes ■ Additional themes on data mining ■ Social impacts of data mining ■ Trends in data mining |
|---|--|

December 1, 2014

Data Mining: Concepts and Techniques

6

Chapter 8. Mining Stream, Time-Series, and Sequence Data

- Mining data streams
- Mining time-series data
- Mining sequence patterns in transactional databases
- Mining sequence patterns in biological data

December 1, 2014

Data Mining: Concepts and Techniques

7

Mining Data Streams

- What is stream data? Why Stream Data Systems?
- Stream data management systems: Issues and solutions
- Stream data cube and multidimensional OLAP analysis
- Stream frequent pattern analysis
- Stream classification
- Stream cluster analysis
- Research issues

December 1, 2014

Data Mining: Concepts and Techniques

8

Characteristics of Data Streams

- Data Streams
 - Data streams—continuous, ordered, changing, fast, huge amount
 - Traditional DBMS—data stored in finite, persistent data sets
- Characteristics
 - Huge volumes of continuous data, possibly infinite
 - Fast changing and requires fast, real-time response
 - Data stream captures nicely our data processing needs of today
 - Random access is expensive—single scan algorithm (*can only have one look*)
 - Store only the summary of the data seen thus far
 - Most stream data are at pretty low-level or multi-dimensional in nature, needs multi-level and multi-dimensional processing

December 1, 2014

Data Mining: Concepts and Techniques

9

Stream Data Applications

- Telecommunication calling records
- Business: credit card transaction flows
- Network monitoring and traffic engineering
- Financial market: stock exchange
- Engineering & industrial processes: power supply & manufacturing
- Sensor, monitoring & surveillance: video streams, RFIDs
- Security monitoring
- Web logs and Web page click streams
- Massive data sets (even saved but random access is too expensive)

December 1, 2014

Data Mining: Concepts and Techniques

10

به این ارقام تعاریف Data Warehouse و یکای مقیاس به تفصیل در

DBMS versus DSMS

- | | |
|---|---|
| ▪ Persistent relations | ▪ Transient streams |
| ▪ One-time queries | ▪ Continuous queries |
| ▪ Random access | ▪ Sequential access |
| ▪ "Unbounded" disk store | ▪ Bounded main memory |
| ▪ Only current state matters | ▪ Historical data is important |
| ▪ No real-time services | ▪ Real-time requirements |
| ▪ Relatively low update rate | ▪ Possibly multi-GB arrival rate |
| ▪ Data at any granularity | ▪ Data at fine granularity |
| ▪ Assume precise data | ▪ Data stale/imprecise |
| ▪ Access plan determined by query processor, physical DB design | ▪ Unpredictable/variable data arrival and characteristics |

Ack. From Motwani's PODS tutorial slides

December 1, 2014

Data Mining: Concepts and Techniques

11

Mining Data Streams

- What is stream data? Why Stream Data Systems?
- Stream data management systems: Issues and solutions
- Stream data cube and multidimensional OLAP analysis
- Stream frequent pattern analysis
- Stream classification
- Stream cluster analysis
- Research issues

December 1, 2014

Data Mining: Concepts and Techniques

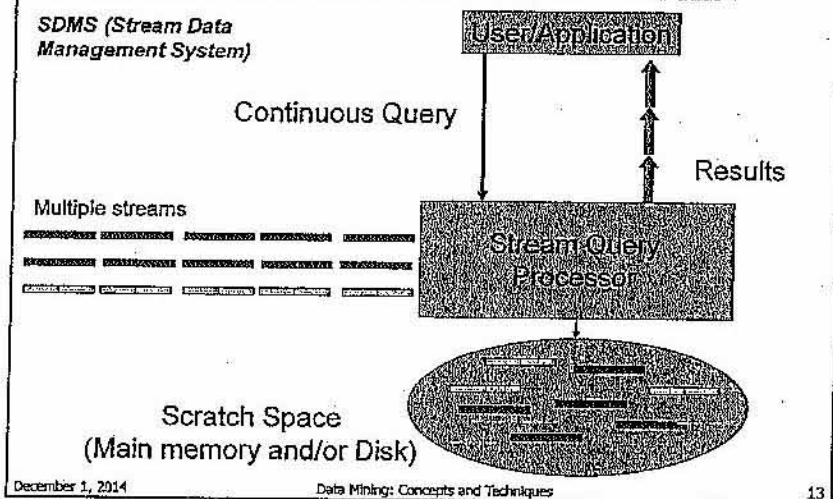
12

OLTP برای متری است OLAP برای مارت است
OLAP داده در جریان است با جزئیات OLAP برای مارت است
(ست)

برای هر کدام از سیستم‌ها یک کامپوننت با نام سرخا هم برای دسترسی اندیس مارتی، برچسب‌گذاری
کنترل یا ریاضی معاری مارت در دسترس باشد مثل معایم ماری OLAP یا

Google
ماتریس مارتی داده‌های مختلف به هم می‌زنند برای تحلیل (Data warehouse)
missing data داده‌های گم شده هستند مثل پاسخ‌دهی به یک سوال در پرسش نامه
با EM الگوریتم می‌توانم داده‌های گم شده را infill کنم
کیفیت داده‌ها مهم است اینجاست جایی که جمع‌آوری شده و در دسترس است و یا با نام مارت
تغییر یافته یا نسبت به وقت جمع‌آوری شده

Architecture: Stream Query Processing



Challenges of Stream Data Processing

- Multiple, continuous, rapid, time-varying, ordered streams
- Main memory computations
- Queries are often continuous
 - Evaluated continuously as stream data arrives
 - Answer updated over time
- Queries are often complex
 - Beyond element-at-a-time processing
 - Beyond stream-at-a-time processing
 - Beyond relational queries (scientific, data mining, OLAP)
- Multi-level/multi-dimensional processing and data mining
 - Most stream data are at low-level or multi-dimensional in nature

فرق Datawarehouse با DBMS :

OLTP جابجایی استفاده می شود با Transaction کاربرد مثل عملیات بانکی، خرید و ... یک ارتباطی با DBMS دارد

OLAP : (برای داده های گسترده و برای تحلیل و گزارش) در این نوع سیستم داده ها به صورت گسترده و گسترده

معمولاً Datawarehouse ← OLAP است

OLTP

Processing Stream Queries

- Query types
 - One-time query vs. continuous query (being evaluated continuously as stream continues to arrive)
 - Predefined query vs. ad-hoc query (issued on-line)
- Unbounded memory requirements
 - For real-time response, main memory algorithm should be used
 - Memory requirement is unbounded if one will join future tuples
- Approximate query answering
 - With bounded memory, it is not always possible to produce exact answers
 - High-quality approximate answers are desired
 - Data reduction and synopsis construction methods
 - Sketches, random sampling, histograms, wavelets, etc.

December 1, 2014

Data Mining: Concepts and Techniques

15

Methodologies for Stream Data Processing

- Major challenges
 - Keep track of a large universe, e.g., pairs of IP address, not ages
- Methodology
 - Synopses (trade-off between accuracy and storage)
 - Use *synopsis data structure*, much smaller ($O(\log^k N)$ space) than their base data set ($O(N)$ space)
 - Compute an *approximate answer* within a *small error range* (factor ϵ of the actual answer)
- Major methods
 - Random sampling
 - Histograms
 - Sliding windows
 - Multi-resolution model
 - Sketches
 - Randomized algorithms

December 1, 2014

Data Mining: Concepts and Techniques

16

Stream Data Processing Methods (1)

- Random sampling (but without knowing the total length in advance)
 - Reservoir sampling: maintain a set of s candidates in the reservoir, which form a true random sample of the element seen so far in the stream. As the data stream flow, every new element has a certain probability (s/N) of replacing an old element in the reservoir.
- Sliding windows
 - Make decisions based only on *recent data* of sliding window size w
 - An element arriving at time t expires at time $t + w$
- Histograms
 - Approximate the frequency distribution of element values in a stream
 - Partition data into a set of contiguous buckets
 - Equal-width (equal value range for buckets) vs. V-optimal (minimizing frequency variance within each bucket)
- Multi-resolution models
 - Popular models: balanced binary trees, micro-clusters, and wavelets

December 1, 2014

Data Mining: Concepts and Techniques

17

Stream Data Processing Methods (2)

- Sketches
 - Histograms and wavelets require multi-passes over the data but sketches can operate in a single pass
 - Frequency moments of a stream $A = \{a_1, \dots, a_N\}$, $F_k = \sum_{i=1}^N m_i^k$
 where v : the universe or domain size, m_i the frequency of i in the sequence
 - Given N elts and v values, sketches can approximate F_0, F_1, F_2 in $O(\log v + \log N)$ space
- Randomized algorithms
 - Monte Carlo algorithm: bound on running time but may not return correct result
 - Chebyshev's inequality: $P(|X - \mu| > k) \leq \frac{\sigma^2}{k^2}$
 - Let X be a random variable with mean μ and standard deviation σ
 - Chernoff bound: $P[X < (1 - \delta)\mu] < e^{-\mu\delta^2/4}$
 - Let X be the sum of independent Poisson trials X_1, \dots, X_n , $\delta \in (0, 1]$
 - The probability decreases exponentially as we move from the mean

December 1, 2014

Data Mining: Concepts and Techniques

18