

Data Warehousing Overview

CS245 Notes 11

Hector Garcia-Molina
Stanford University

CS 245

Notes11

1

Warehousing

- Growing industry: \$8 billion in 1998
- Range from desktop to huge:
 - ◆ Walmart: 900-CPU, 2,700 disk, 23TB Teradata system
- Lots of buzzwords, hype
 - ◆ slice & dice, rollup, MOLAP, pivot, ...

CS 245

Notes11

2

Outline

- What is a data warehouse?
- Why a warehouse?
- Models & operations
- Implementing a warehouse
- Future directions

CS 245

Notes11

3

What is a Warehouse?

- Collection of diverse data
 - ◆ subject oriented
 - ◆ aimed at executive, decision maker
 - ◆ often a copy of operational data
 - ◆ with value-added data (e.g., summaries, history)
 - ◆ integrated
 - ◆ time-varying
 - ◆ non-volatile



CS 245

Notes11

4

What is a Warehouse?

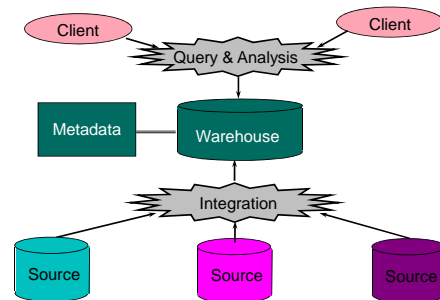
- Collection of tools
 - ◆ gathering data
 - ◆ cleansing, integrating, ...
 - ◆ querying, reporting, analysis
 - ◆ data mining
 - ◆ monitoring, administering warehouse

CS 245

Notes11

5

Warehouse Architecture



CS 245

Notes11

6

Motivating Examples

- Forecasting
- Comparing performance of units
- Monitoring, detecting fraud
- Visualization

CS 245

Notes11

7

Why a Warehouse?

- Two Approaches:
 - ◆ Query-Driven (Lazy)
 - ◆ Warehouse (Eager)

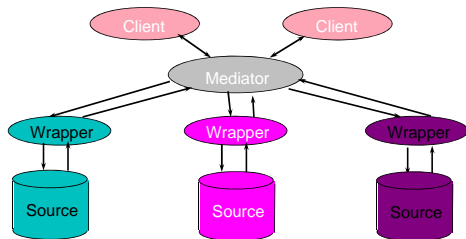


CS 245

Notes11

8

Query-Driven Approach



CS 245

Notes11

9

Advantages of Warehousing

- High query performance
- Queries not visible outside warehouse
- Local processing at sources unaffected
- Can operate when sources unavailable
- Can query data not stored in a DBMS
 - ◆ Modify, summarize (store aggregates)
 - ◆ Add historical information

CS 245

Notes11

10

Advantages of Query-Driven

- No need to copy data
 - ◆ less storage
 - ◆ no need to purchase data
- More up-to-date data
- Query needs can be unknown
- Only query interface needed at sources
- May be less draining on sources

CS 245

Notes11

11

OLTP vs. OLAP

- OLTP: On Line Transaction Processing
 - ◆ Describes processing at operational sites
- OLAP: On Line Analytical Processing
 - ◆ Describes processing at warehouse

CS 245

Notes11

12

OLTP vs. OLAP

OLTP

- Mostly updates
- Many small transactions
- Mb-Tb of data
- Raw data
- Clerical users
- Up-to-date data
- Consistency, recoverability critical

OLAP

- Mostly reads
- Queries long, complex
- Gb-Tb of data
- Summarized, consolidated data
- Decision-makers, analysts as users

CS 245

Notes11

13

Data Marts

- Smaller warehouses
- Spans part of organization
 - ◆ e.g., marketing (customers, products, sales)
- Do not require enterprise-wide consensus
 - ◆ but long term integration problems?

CS 245

Notes11

14

Warehouse Models & Operators

- Data Models
 - ◆ relations
 - ◆ stars & snowflakes
 - ◆ cubes
- Operators
 - ◆ slice & dice
 - ◆ roll-up, drill down
 - ◆ pivoting
 - ◆ other

CS 245

Notes11

15

Star

product	prodid	name	price
	p1	bolt	10
	p2	nut	5

store	storeid	city
	c1	nyc
	c2	sfo
	c3	la

sale	orderid	date	custid	prodid	storeid	qty	amt
	o100	1/7/97	53	p1	c1	1	12
	o102	2/7/97	53	p2	c1	2	11
	o105	3/8/97	111	p1	c3	5	50

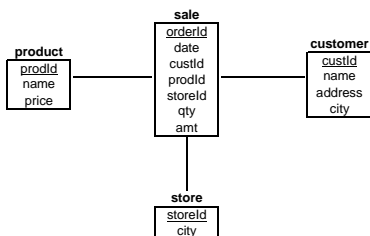
customer	custid	name	address	city
	53	joe	10 main	sfo
	81	fred	12 main	sfo
	111	sally	80 willow	la

CS 245

Notes11

16

Star Schema



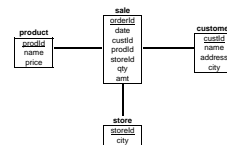
CS 245

Notes11

17

Terms

- Fact table
- Dimension tables
- Measures

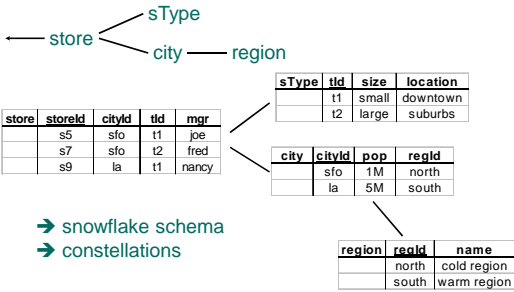


CS 245

Notes11

18

Dimension Hierarchies



CS 245

Notes11

19

Cube

Fact table view:

sale	prold	storeld	amt
p1	c1	12	
p2	c1	11	
p1	c3	50	
p2	c2	8	

Multi-dimensional cube:

	c1	c2	c3
p1	12		50
p2	11	8	

dimensions = 2

CS 245

Notes11

20

3-D Cube

Fact table view:

sale	prold	storeld	date	amt
p1	c1	1	12	
p2	c1	1	11	
p1	c3	1	50	
p2	c2	1	8	
p1	c1	2	44	
p1	c2	2	4	

Multi-dimensional cube:

	c1	c2	c3
day 2	p1 44	c2 4	c3 50
day 1	p1 12	c2 8	c3 50
	p2 11	8	

dimensions = 3

CS 245

Notes11

21

ROLAP vs. MOLAP

- ROLAP: Relational On-Line Analytical Processing
- MOLAP: Multi-Dimensional On-Line Analytical Processing

CS 245

Notes11

22

Aggregates

- Add up amounts for day 1
- In SQL: `SELECT sum(amt) FROM SALE WHERE date = 1`

sale	prold	storeld	date	amt
p1	c1	1	12	
p2	c1	1	11	
p1	c3	1	50	
p2	c2	1	8	
p1	c1	2	44	
p1	c2	2	4	



81

CS 245

Notes11

23

Aggregates

- Add up amounts by day
- In SQL: `SELECT date, sum(amt) FROM SALE GROUP BY date`

sale	prold	storeld	date	amt
p1	c1	1	12	
p2	c1	1	11	
p1	c3	1	50	
p2	c2	1	8	
p1	c1	2	44	
p1	c2	2	4	



ans	date	sum
	1	81
	2	48

CS 245

Notes11

24

Another Example

- Add up amounts by day, product
- In SQL: `SELECT date, sum(amt) FROM SALE GROUP BY date, prodId`

sale	prodId	storeId	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

sale	prodId	date	amt
	p1	1	62
	p2	1	19
	p1	2	48

→ rollup →
← drill-down ←

CS 245 Notes11 25

Aggregates

- Operators: sum, count, max, min, median, ave
- “Having” clause
- Using dimension hierarchy
 - ◆ average by region (within store)
 - ◆ maximum by month (within date)

CS 245 Notes11 26

Cube Aggregation

Example: computing sums

→ rollup →
← drill-down ←

CS 245 Notes11 27

Cube Operators

CS 245 Notes11 28

Extended Cube

CS 245 Notes11 29

Aggregation Using Hierarchies

customer
|
region
|
country

(customer c1 in Region A;
customers c2, c3 in Region B)

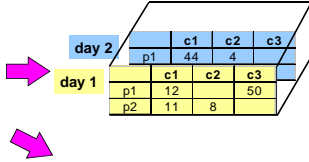
CS 245 Notes11 30

Pivoting

Fact table view:

sale	prold	storeld	date	amt
p1	c1	1	1	12
p2	c1	1	1	11
p1	c3	1	1	50
p2	c2	1	1	8
p1	c1	2	2	44
p1	c2	2	2	4

Multi-dimensional cube:



	c1	c2	c3
day 1	12	8	50
day 2	44	4	

CS 245

Notes11

31

Query & Analysis Tools

- Query Building
- Report Writers (comparisons, growth, graphs,...)
- Spreadsheet Systems
- Web Interfaces
- Data Mining

CS 245

Notes11

32

Other Operations

- Time functions
 - ◆ e.g., time average
- Computed Attributes
 - ◆ e.g., commission = sales * rate
- Text Queries
 - ◆ e.g., find documents with words X AND B
 - ◆ e.g., rank documents by frequency of words X, Y, Z

CS 245

Notes11

33

Data Mining

- Decision Trees
- Clustering
- Association Rules

CS 245

Notes11

34

Decision Trees

Example:

- Conducted survey to see what customers were interested in new model car
- Want to select customers for advertising campaign

sale	custld	car	age	city	newCar
	c1	taurus	27	sf	yes
	c2	van	35	la	yes
	c3	van	40	sf	yes
	c4	taurus	22	sf	yes
	c5	merc	50	la	no
	c6	taurus	25	la	no

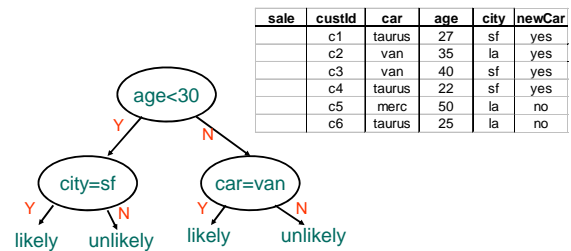
training set

CS 245

Notes11

35

One Possibility

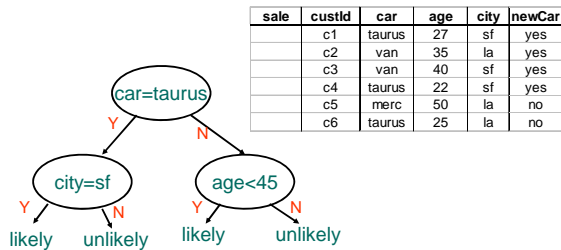


CS 245

Notes11

36

Another Possibility



CS 245

Notes11

37

Issues

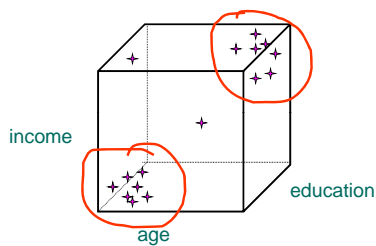
- Decision tree cannot be “too deep”
 - would not have statistically significant amounts of data for lower decisions
- Need to select tree that most reliably predicts outcomes

CS 245

Notes11

38

Clustering



CS 245

Notes11

39

Another Example: Text

- Each document is a vector
 - ◆ e.g., <100110...> contains words 1,4,5,...
- Clusters contain “similar” documents
- Useful for understanding, searching documents



CS 245

Notes11

40

Issues

- Given desired number of clusters?
- Finding “best” clusters
- Are clusters semantically meaningful?
 - ◆ e.g., “yuppies” cluster?
- Using clusters for disk storage

CS 245

Notes11

41

Association Rule Mining

transaction id	customer id	products bought
tran1	cust33	p2, p5, p8
tran2	cust45	p5, p8, p11
tran3	cust12	p1, p9
tran4	cust40	p5, p8, p11
tran5	cust12	p2, p9
tran6	cust12	p9

sales records:

market-basket data

- Trend: Products p5, p8 often bought together
- Trend: Customer 12 likes product p9

CS 245

Notes11

42

Association Rule

- Rule: $\{p_1, p_3, p_8\}$
- Support: number of baskets where these products appear
- High-support set: support \geq threshold s
- Problem: find all high support sets

CS 245

Notes11

43

Finding High-Support Pairs

- Baskets(basket, item)
- `SELECT I.item, J.item, COUNT(I.basket)`
`FROM Baskets I, Baskets J`
`WHERE I.basket = J.basket AND`
`I.item < J.item` WHY?
`GROUP BY I.item, J.item`
`HAVING COUNT(I.basket) >= s;`

CS 245

Notes11

44

Example

basket	item
t1	p2
t1	p5
t1	p8
t2	p5
t2	p8
t2	p11
...	...

→

basket	item1	item2
t1	p2	p5
t1	p2	p8
t1	p5	p8
t2	p5	p8
t2	p5	p11
t2	p8	p11
...

check if count \geq s

CS 245

Notes11

45

Issues

- Performance for size 2 rules

basket	item
t1	p2
t1	p5
t1	p8
t2	p5
t2	p8
t2	p11
...	...

big ↓

basket	item1	item2
t1	p2	p5
t1	p2	p8
t1	p5	p8
t2	p5	p8
t2	p8	p11
...

even bigger! ↓

- Performance for size k rules

CS 245

Notes11

46

Implementing a Warehouse

- *Monitoring*: Sending data from sources
- *Integrating*: Loading, cleansing,...
- *Processing*: Query processing, indexing, ...
- *Managing*: Metadata, Design, ...

CS 245

Notes11

47

Monitoring

- Source Types: relational, flat file, IMS, VSAM, IDMS, WWW, news-wire, ...
- Incremental vs. Refresh

customer	id	name	address	city
	53	joe	10 main	sfo
	81	fred	12 main	sfo
	111	sally	80 willow	la

← new

CS 245

Notes11

48

Monitoring Techniques

- Periodic snapshots
- Database triggers
- Log shipping
- Data shipping (replication service)
- Transaction shipping
- Polling (queries to source)
- Screen scraping
- Application level monitoring

Advantages & Disadvantages!!

CS 245

Notes11

49

Monitoring Issues

- Frequency
 - ◆ periodic: daily, weekly, ...
 - ◆ triggered: on "big" change, lots of changes, ...
- Data transformation
 - ◆ convert data to uniform format
 - ◆ remove & add fields (e.g., add date to get history)
- Standards (e.g., ODBC)
- Gateways

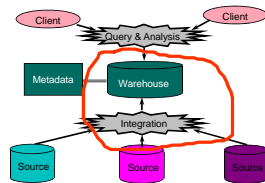
CS 245

Notes11

50

Integration

- Data Cleaning
- Data Loading
- Derived Data



CS 245

Notes11

51

Data Cleaning

- Migration (e.g., yen ⇔ dollars)
- Scrubbing: use domain-specific knowledge (e.g., social security numbers)
- Fusion (e.g., mail list, customer merging)


```

            billing DB → customer1(Joe)
            service DB → customer2(Joe)
            customer1(Joe) → merged_customer(Joe)
            customer2(Joe) → merged_customer(Joe)
            
```
- Auditing: discover rules & relationships (like data mining)

CS 245

Notes11

52

Loading Data

- Incremental vs. refresh
- Off-line vs. on-line
- Frequency of loading
 - ◆ At night, 1x a week/month, continuously
- Parallel/Partitioned load

CS 245

Notes11

53

Derived Data

- Derived Warehouse Data
 - ◆ indexes
 - ◆ aggregates
 - ◆ materialized views (next slide)
- When to update derived data?
- Incremental vs. refresh

CS 245

Notes11

54

Materialized Views

- Define new warehouse relations using SQL expressions

sale	prodid	storeid	date	amt	product	id	name	price
	p1	c1	1	12		p1	bolt	10
	p2	c1	1	11				
	p1	c3	1	50				
	p2	c2	1	8				
	p1	c1	2	44				
	p1	c2	2	4				

joinTb	prodid	name	price	storeid	date	amt
	p1	bolt	10	c1	1	12
	p2	nut	5	c1	1	11
	p1	bolt	10	c3	1	50
	p2	nut	5	c2	1	8
	p1	bolt	10	c1	2	44
	p1	bolt	10	c2	2	4

does not exist at any source

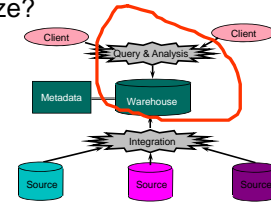
CS 245

Notes11

55

Processing

- ROLAP servers vs. MOLAP servers
- Index Structures
- What to Materialize?
- Algorithms



CS 245

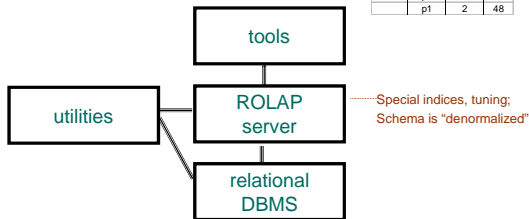
Notes11

56

ROLAP Server

- Relational OLAP Server

sale	prodid	date	sum
	p1	1	82
	p2	1	19
	p1	2	48



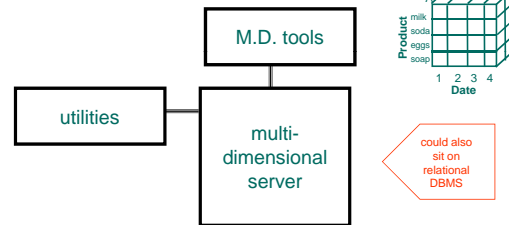
CS 245

Notes11

57

MOLAP Server

- Multi-Dimensional OLAP Server



CS 245

Notes11

58

Index Structures

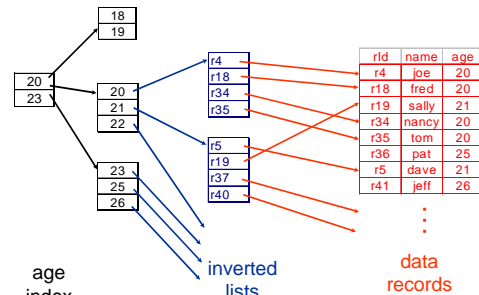
- Traditional Access Methods
 - B-trees, hash tables, R-trees, grids, ...
- Popular in Warehouses
 - inverted lists
 - bit map indexes
 - join indexes
 - text indexes

CS 245

Notes11

59

Inverted Lists



CS 245

Notes11

60

Using Inverted Lists

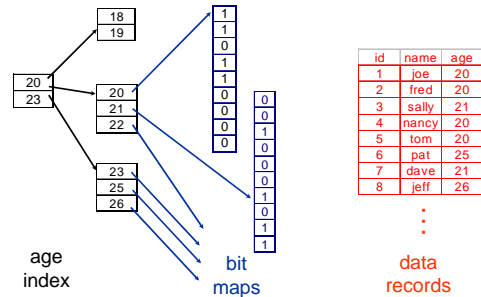
- Query:
 - ◆ Get people with age = 20 and name = "fred"
- List for age = 20: r4, r18, r34, r35
- List for name = "fred": r18, r52
- Answer is intersection: r18

CS 245

Notes11

61

Bit Maps



CS 245

Notes11

62

Using Bit Maps

- Query:
 - ◆ Get people with age = 20 and name = "fred"
- List for age = 20: 1101100000
- List for name = "fred": 0100000001
- Answer is intersection: 01000000000
- Good if domain cardinality small
- Bit vectors can be compressed

CS 245

Notes11

63

Join

- "Combine" SALE, PRODUCT relations
- In SQL: SELECT * FROM SALE, PRODUCT

sale	prodlid	storeid	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

product	id	name	price
	p1	bolt	10
	p2	nut	5

joinTb	prodlid	name	price	storeid	date	amt
	p1	bolt	10	c1	1	12
	p2	nut	5	c1	1	11
	p1	bolt	10	c3	1	50
	p2	nut	5	c2	1	8
	p1	bolt	10	c1	2	44
	p1	bolt	10	c2	2	4

CS 245

Notes11

64

Join Indexes

join index

product	id	name	price	index
	p1	bolt	10	r1,r3,r5,r6
	p2	nut	5	r2,r4

sale	rlid	prodlid	storeid	date	amt
	r1	p1	c1	1	12
	r2	p2	c1	1	11
	r3	p1	c3	1	50
	r4	p2	c2	1	8
	r5	p1	c1	2	44
	r6	p1	c2	2	4

CS 245

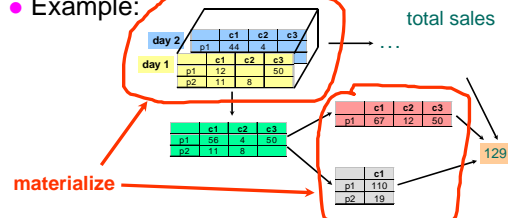
Notes11

65

What to Materialize?

- Store in warehouse results useful for common queries

- Example:



CS 245

Notes11

66

Materialization Factors

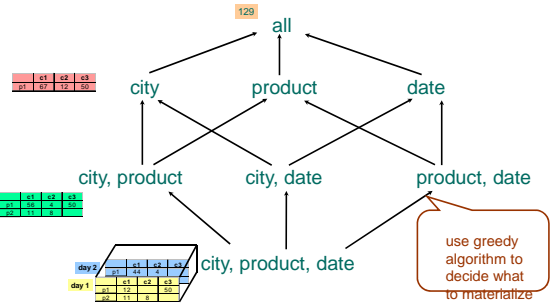
- Type/frequency of queries
- Query response time
- Storage cost
- Update cost

CS 245

Notes11

67

Cube Aggregates Lattice



CS 245

Notes11

68

Dimension Hierarchies



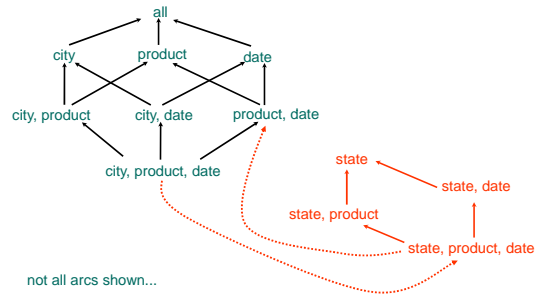
cities	city	state
	c1	CA
	c2	NY

CS 245

Notes11

69

Dimension Hierarchies



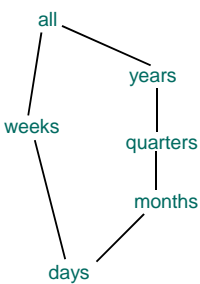
not all arcs shown...

CS 245

Notes11

70

Interesting Hierarchy



time	day	week	month	quarter	year
1	1	1	1	1	2000
2	1	1	1	1	2000
3	1	1	1	1	2000
4	1	1	1	1	2000
5	1	1	1	1	2000
6	1	1	1	1	2000
7	1	1	1	1	2000
8	2	1	1	1	2000

conceptual dimension table

CS 245

Notes11

71

Algorithms

- Query Optimization
- Parallel Processing
- Data Mining

CS 245

Notes11

72

Example: Association Rules

- How do we perform rule mining efficiently?
- Observation: If set X has support t , then each X subset must have at least support t
- For 2-sets:
 - if we need support s for $\{i, j\}$
 - then each i, j must appear in at least s baskets

CS 245

Notes11

73

Algorithm for 2-Sets

- Find OK products
 - those appearing in s or more baskets
- Find high-support pairs using only OK products

CS 245

Notes11

74

Algorithm for 2-Sets

- INSERT INTO okBaskets(basket, item)


```
SELECT basket, item
FROM Baskets
GROUP BY item
HAVING COUNT(basket) >= s;
```
- Perform mining on okBaskets


```
SELECT I.item, J.item, COUNT(I.basket)
FROM okBaskets I, okBaskets J
WHERE I.basket = J.basket AND
I.item < J.item
GROUP BY I.item, J.item
HAVING COUNT(I.basket) >= s;
```

CS 245

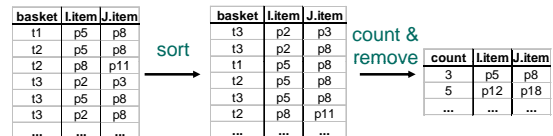
Notes11

75

Counting Efficiently

- One way:

threshold = 3



CS 245

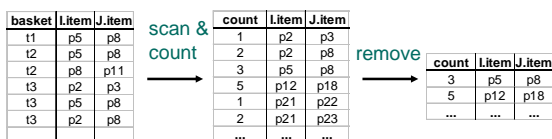
Notes11

76

Counting Efficiently

- Another way:

threshold = 3

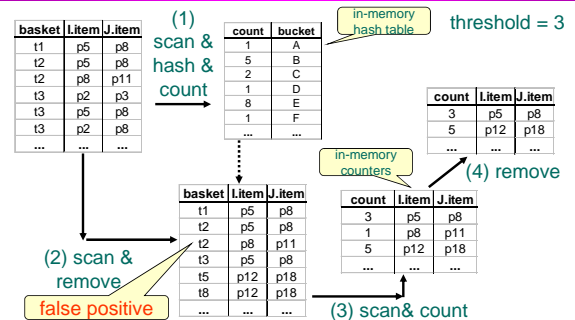


CS 245

Notes11

77

Yet Another Way



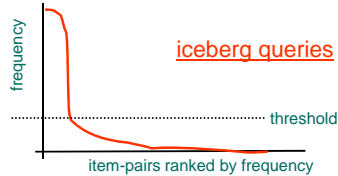
CS 245

Notes11

78

Discussion

- Hashing scheme: 2 (or 3) scans of data
- Sorting scheme: requires a sort!
- Hashing works well if few high-support pairs and many low-support ones



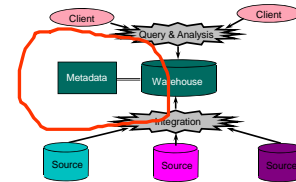
CS 245

Notes11

79

Managing

- Metadata
- Warehouse Design
- Tools



CS 245

Notes11

80

Metadata

- Administrative
 - ◆ definition of sources, tools, ...
 - ◆ schemas, dimension hierarchies, ...
 - ◆ rules for extraction, cleaning, ...
 - ◆ refresh, purging policies
 - ◆ user profiles, access control, ...

CS 245

Notes11

81

Metadata

- Business
 - ◆ business terms & definition
 - ◆ data ownership, charging
- Operational
 - ◆ data lineage
 - ◆ data currency (e.g., active, archived, purged)
 - ◆ use stats, error reports, audit trails

CS 245

Notes11

82

Design

- What data is needed?
- Where does it come from?
- How to clean data?
- How to represent in warehouse (schema)?
- What to summarize?
- What to materialize?
- What to index?

CS 245

Notes11

83

Tools

- Development
 - ◆ design & edit: schemas, views, scripts, rules, queries, reports
- Planning & Analysis
 - ◆ what-if scenarios (schema changes, refresh rates), capacity planning
- Warehouse Management
 - ◆ performance monitoring, usage patterns, exception reporting
- System & Network Management
 - ◆ measure traffic (sources, warehouse, clients)
- Workflow Management
 - ◆ "reliable scripts" for cleaning & analyzing data

CS 245

Notes11

84

Current State of Industry

- Extraction and integration done off-line
 - ◆ Usually in large, time-consuming, batches
- Everything copied at warehouse
 - ◆ Not selective about what is stored
 - ◆ Query benefit vs storage & update cost
- Query optimization aimed at OLTP
 - ◆ High throughput instead of fast response
 - ◆ Process whole query before displaying anything

CS 245

Notes11

85

Future Directions

- Better performance
- Larger warehouses
- Easier to use
- What are companies & research labs working on?

CS 245

Notes11

86

Research (1)

- Incremental Maintenance
- Data Consistency
- Data Expiration
- Recovery
- Data Quality
- Error Handling (Back Flush)

CS 245

Notes11

87

Research (2)

- Rapid Monitor Construction
- Temporal Warehouses
- Materialization & Index Selection
- Data Fusion
- Data Mining
- Integration of Text & Relational Data

CS 245

Notes11

88

Conclusions

- Massive amounts of data and complexity of queries will push limits of current warehouses
- Need better systems:
 - ◆ easier to use
 - ◆ provide quality information

CS 245

Notes11

89