

# Data Warehousing Overview

## CS245 Notes 11

Hector Garcia-Molina  
Stanford University

CS 245

Notes11

1

## Outline

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

CS 245

Notes11

2

## 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

3

## What is a Warehouse?

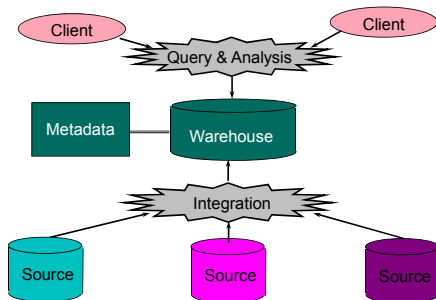
- Collection of tools
  - ◆ gathering data
  - ◆ cleansing, integrating, ...
  - ◆ querying, reporting, analysis
  - ◆ data mining
  - ◆ monitoring, administering warehouse

CS 245

Notes11

4

## Warehouse Architecture



CS 245

Notes11

5

## Motivating Examples

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

CS 245

Notes11

6

## Alternative to Warehousing

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

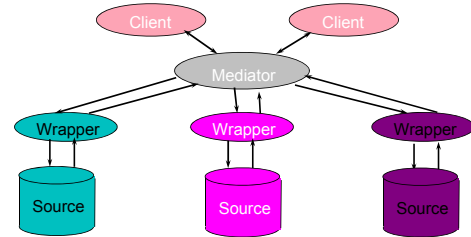


CS 245

Notes11

7

## Query-Driven Approach



CS 245

Notes11

8

## 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

9

## 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

10

## Warehouse Models & Operators

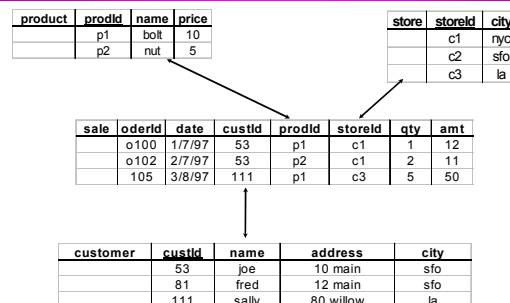
- Data Models
  - relational
  - cubes
- Operators

CS 245

Notes11

11

## Star

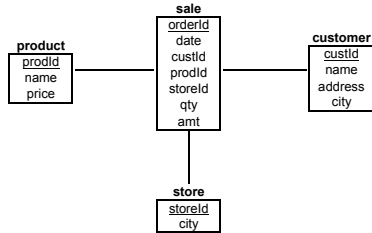


CS 245

Notes11

12

## Star Schema



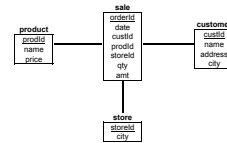
CS 245

Notes11

13

## Terms

- Fact table
- Dimension tables
- Measures

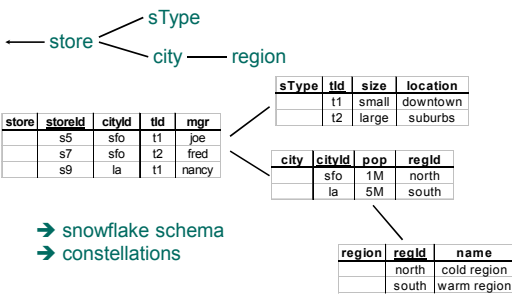


CS 245

Notes11

14

## Dimension Hierarchies



CS 245

Notes11

15

## Cube

Fact table view:

sale	prodid	storeid	amt
	p1	c1	12
	p2	c1	11
	p1	c3	50
	p2	c2	8

Multi-dimensional cube:

	c1	c2	c3
p1	12		
p2	11	8	

dimensions = 2

CS 245

Notes11

16

## 3-D Cube

Fact table view:

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

Multi-dimensional cube:

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

dimensions = 3

CS 245

Notes11

17

## Operators

- Traditional
  - ◆ selection
  - ◆ aggregation
  - ◆ ...
- Analysis
  - ◆ clean data
  - ◆ find trends
  - ◆ ...
- Relational
  - Cube

CS 245

Notes11

18

## Aggregates

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

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



81

CS 245

Notes11

19

## Aggregates

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

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



ans	date	sum
	1	81
	2	48

CS 245

Notes11

20

## 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

21

## Aggregates

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

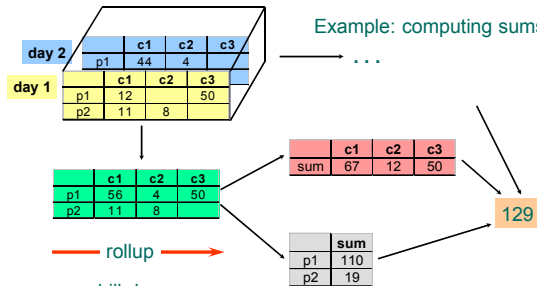
CS 245

Notes11

22

## Cube Aggregation

Example: computing sums

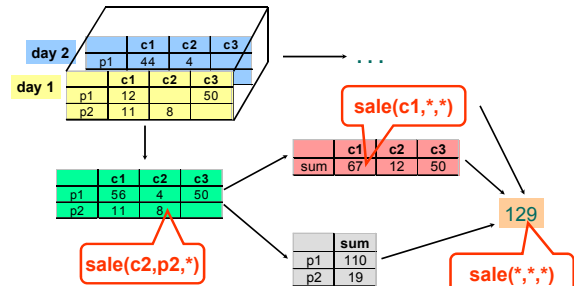


CS 245

Notes11

23

## Cube Operators

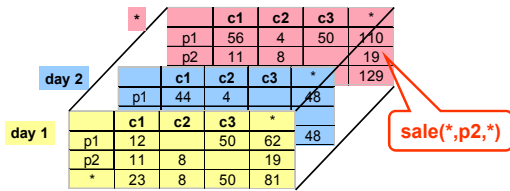


CS 245

Notes11

24

## Extended Cube

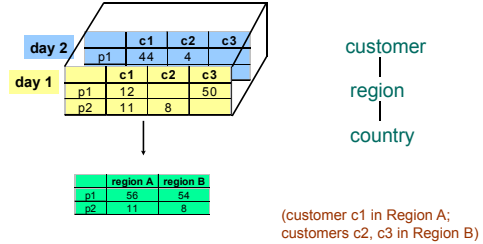


CS 245

Notes11

25

## Aggregation Using Hierarchies



CS 245

Notes11

26

## Data Analysis

- Decision Trees
- Clustering
- Association Rules

CS 245

Notes11

27

## Decision Trees

Example:

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

sale	custid	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

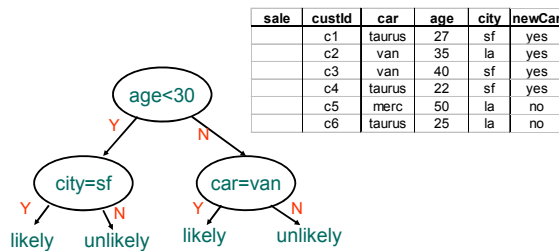
training set

CS 245

Notes11

28

## One Possibility

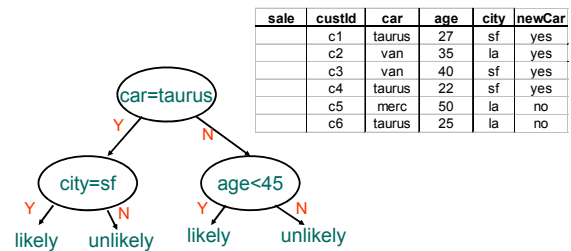


CS 245

Notes11

29

## Another Possibility



CS 245

Notes11

30

## Issues

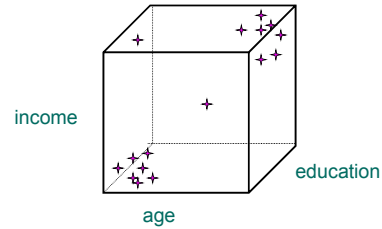
- Decision tree cannot be “too deep”
  - would not have statistically significant amounts of data for lower decisions
- Need to select tree that most reliably predicts outcomes

CS 245

Notes11

31

## Clustering

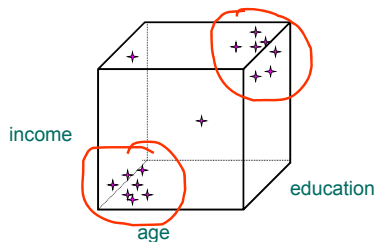


CS 245

Notes11

32

## Clustering



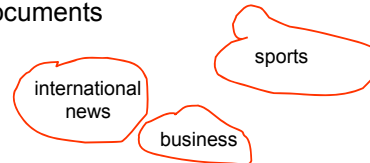
CS 245

Notes11

33

## 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

34

## 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

35

## Association Rule Mining

sales records:

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

market-basket data

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

CS 245

Notes11

36

## 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

37

## Implementation Issues

- ETL (Extraction, transformation, loading)
  - ◆ Getting data to the warehouse
  - ◆ Entity Resolution
- What to materialize?
- Efficient Analysis
  - ◆ Association rule mining
  - ◆ ...

CS 245

Notes11

38

## ETL: 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

39

## ETL: Data Cleaning

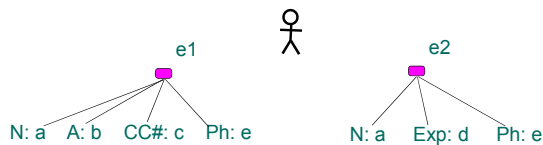
- Migration (e.g., yen  $\Rightarrow$  dollars)
  - Scrubbing: use domain-specific knowledge (e.g., social security numbers)
  - Fusion (e.g., mail list, customer merging)
- billing DB  $\longrightarrow$  customer1(Joe)  $\longrightarrow$  merged\_customer(Joe)  
 service DB  $\longrightarrow$  customer2(Joe)  $\longrightarrow$  merged\_customer(Joe)
- Auditing: discover rules & relationships (like data mining)

CS 245

Notes11

40

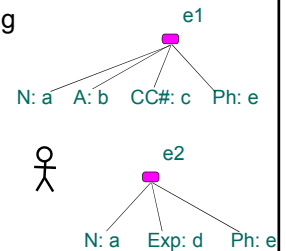
## More details: Entity Resolution



41

## Applications

- comparison shopping
- mailing lists
- classified ads
- customer files
- counter-terrorism



42

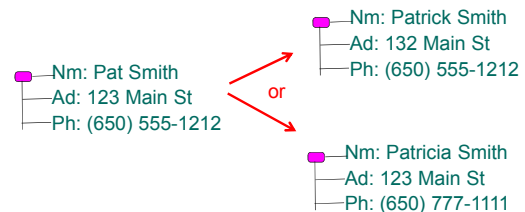
## Why is ER Challenging?

- Huge data sets
- No unique identifiers
- Lots of uncertainty
- Many ways to skin the cat

43

## Taxonomy: Pairwise vs Global

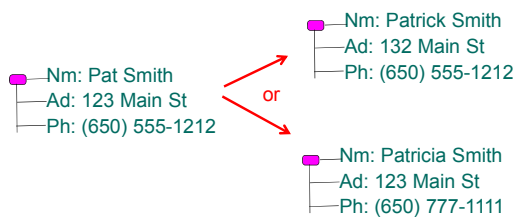
- Decide if  $r, s$  match only by looking at  $r, s$ ?
- Or need to consider more (all) records?



44

## Taxonomy: Pairwise vs Global

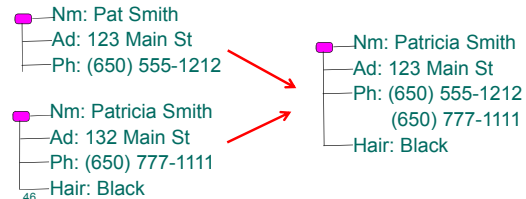
- Global matching complicates things a lot!
  - ◆ e.g., change decision as new records arrive



45

## Taxonomy: Outcome

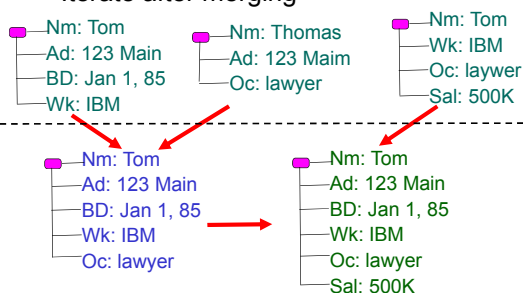
- Partition of records
  - ◆ e.g., comparison shopping
- Merged records



46

## Taxonomy: Outcome

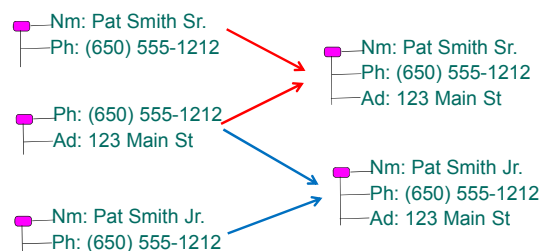
- Iterate after merging



47

## Taxonomy: Record Reuse

- One record related to multiple entities?

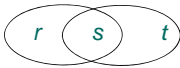


48

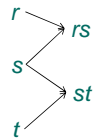


## Taxonomy: Record Reuse

### • Partitions



### • Merges



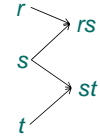
49

## Taxonomy: Record Reuse

### • Partitions



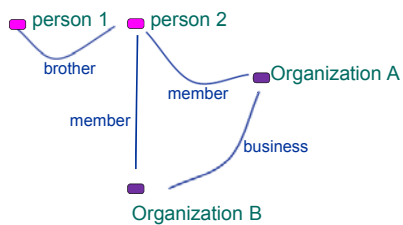
### • Merges



• Record reuse → complex and expensive!

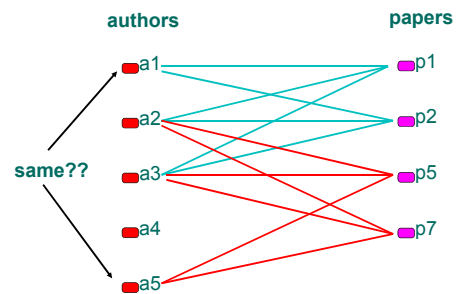
50

## Taxonomy: Multiple Entity Types



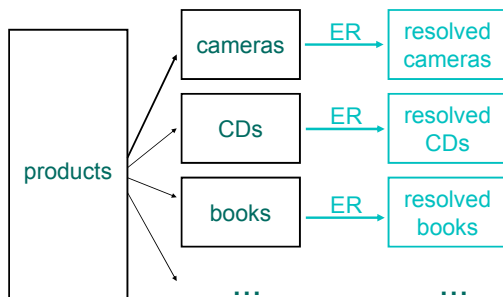
51

## Taxonomy: Multiple Entity Types



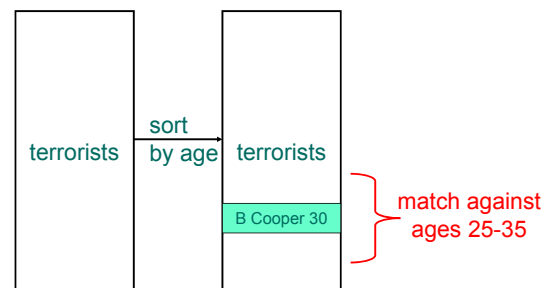
52

## Taxonomy: Exact vs Approximate




53

## Taxonomy: Exact vs Approximate



54

## Implementation Issues

- ETL (Extraction, transformation, loading)
  - ◆ Getting data to the warehouse
  - ◆ Entity Resolution
- What to materialize? 
- Efficient Analysis
  - ◆ Association rule mining
  - ◆ ...

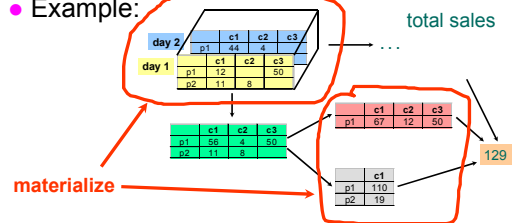
CS 245

Notes11

55

## What to Materialize?

- Store in warehouse results useful for common queries
- Example:



CS 245

Notes11

56

## Materialization Factors

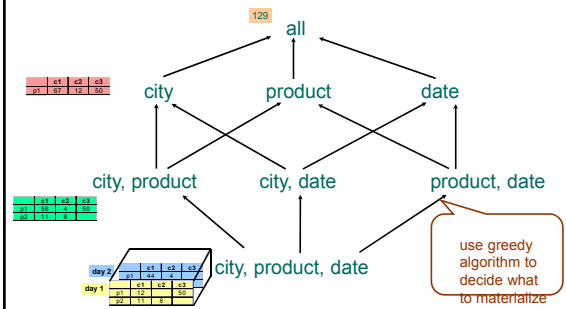
- Type/frequency of queries
- Query response time
- Storage cost
- Update cost

CS 245

Notes11

57

## Cube Aggregates Lattice



CS 245

Notes11

58

## Dimension Hierarchies

all  
|  
state  
|  
city

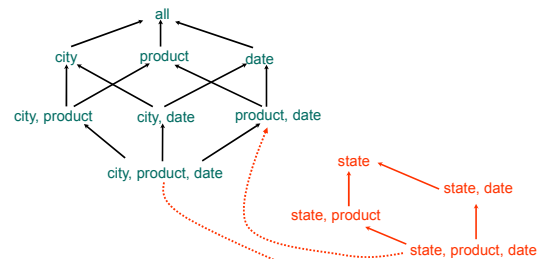
cities	city	state
	c1	CA
	c2	NY

CS 245

Notes11

59

## Dimension Hierarchies



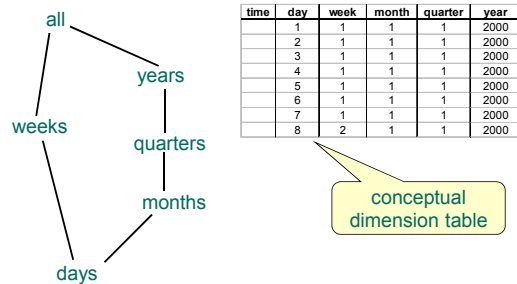
not all arcs shown...

CS 245

Notes11

60

## Interesting Hierarchy



CS 245

Notes11

61

## Implementation Issues

- ETL (Extraction, transformation, loading)
  - Getting data to the warehouse
  - Entity Resolution
- What to materialize?
- Efficient Analysis
  - Association rule mining
  - ...

CS 245

Notes11

62

## 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
 GROUP BY I.item, J.item
 HAVING COUNT(I.basket) >= s;

CS 245

Notes11

63

## 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
 GROUP BY I.item, J.item
 HAVING COUNT(I.basket) >= s;

WHY?

CS 245

Notes11

64

## 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
...	...	...

CS 245

Notes11

65

## 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 ≥ s

CS 245

Notes11

66

## 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	p5	p11
t2	p8	p11
...	...	...

even bigger!

- Performance for size  $k$  rules

CS 245

Notes11

67

## Association Rules

- How do we perform rule mining efficiently?

CS 245

Notes11

68

## 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  $t$

CS 245

Notes11

69

## 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

70

## 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

71

## Algorithm for 2-Sets

- INSERT INTO okBaskets(basket, item)
 

```
SELECT basket, item
FROM Baskets
GROUP BY item
HAVING COUNT(basket) >= s;
```

CS 245

Notes11

72

## 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

73

## Counting Efficiently

- One way:

threshold = 3

basket	I.item	J.item
t1	p5	p8
t2	p5	p8
t2	p8	p11
t3	p2	p3
t3	p5	p8
t3	p2	p8
...	...	...

CS 245

Notes11

74

## Counting Efficiently

- One way:

threshold = 3

basket	I.item	J.item
t1	p5	p8
t2	p5	p8
t2	p8	p11
t3	p2	p3
t3	p5	p8
t3	p2	p8
...	...	...

sort

basket	I.item	J.item
t3	p2	p3
t3	p2	p8
t1	p5	p8
t2	p5	p8
t3	p5	p8
t2	p8	p11
...	...	...

CS 245

Notes11

75

## Counting Efficiently

- One way:

threshold = 3

basket	I.item	J.item
t1	p5	p8
t2	p5	p8
t2	p8	p11
t3	p2	p3
t3	p5	p8
t3	p2	p8
...	...	...

sort

basket	I.item	J.item
t3	p2	p3
t3	p2	p8
t1	p5	p8
t2	p5	p8
t3	p5	p8
t2	p8	p11
...	...	...

count & remove

count	I.item	J.item
3	p5	p8
5	p12	p18
...	...	...

CS 245

Notes11

76

## Counting Efficiently

- Another way:

threshold = 3

basket	I.item	J.item
t1	p5	p8
t2	p5	p8
t2	p8	p11
t3	p2	p3
t3	p5	p8
t3	p2	p8
...	...	...

CS 245

Notes11

77

## Counting Efficiently

- Another way:

threshold = 3

basket	I.item	J.item
t1	p5	p8
t2	p5	p8
t2	p8	p11
t3	p2	p3
t3	p5	p8
t3	p2	p8
...	...	...

scan & count

count	I.item	J.item
1	p2	p3
2	p2	p8
3	p5	p8
5	p12	p18
1	p21	p22
2	p21	p23
...	...	...

keep counter array in memory

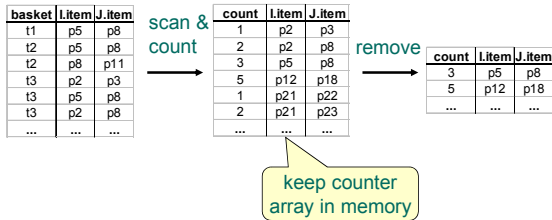
CS 245

Notes11

78

## Counting Efficiently

- Another way: threshold = 3



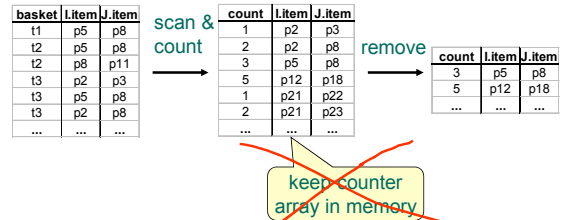
CS 245

Notes11

79

## Counting Efficiently

- Another way: threshold = 3

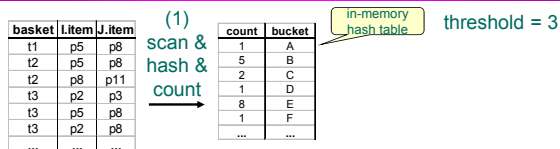


CS 245

Notes11

80

## Yet Another Way

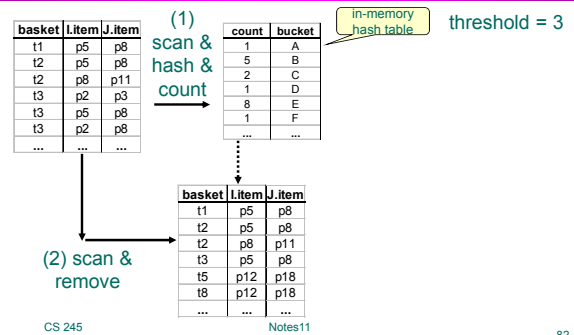


CS 245

Notes11

81

## Yet Another Way

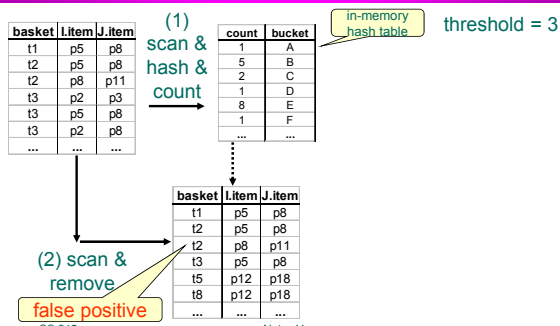


CS 245

Notes11

82

## Yet Another Way

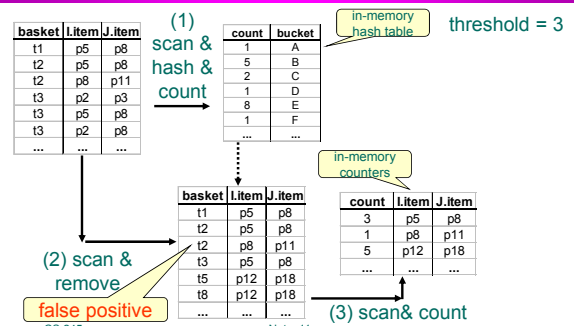


CS 245

Notes11

83

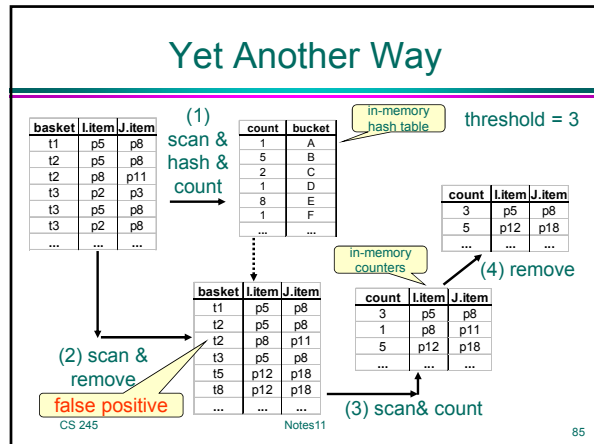
## Yet Another Way



CS 245

Notes11

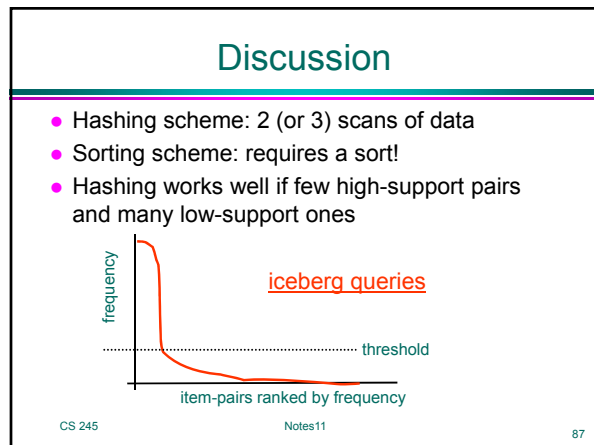
84



### 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

86



### Implementation Issues

- ETL (Extraction, transformation, loading)
  - ◆ Getting data to the warehouse
  - ◆ Entity Resolution
- What to materialize?
- Efficient Analysis
  - ◆ Association rule mining
  - ◆ ...

88

### Extra: Data Mining in the InfoLab

#### Recommendations in CourseRank

user	quarters			
	q1	q2	q3	q4
u1	a: 5	b: 5	d: 5	
u2	a: 1	e: 2	d: 4	f: 3
u3	g: 4	h: 2	e: 3	f: 3
u4	b: 2	g: 4	h: 4	e: 4
u	a: 5	g: 4	e: 4	

89

### Extra: Data Mining in the InfoLab

#### Recommendations in CourseRank

user	quarters			
	q1	q2	q3	q4
u1	a: 5	b: 5	d: 5	
u2	a: 1	e: 2	d: 4	f: 3
u3	g: 4	h: 2	e: 3	f: 3
u4	b: 2	g: 4	h: 4	e: 4
u	a: 5	g: 4	e: 4	

u3 and u4 are similar to u

Recommend h

90

## Extra: Data Mining in the InfoLab

### Recommendations in CourseRank

user	quarters			
	q1	q2	q3	q4
u1	a: 5	b: 5	d: 5	
u2	a: 1	e: 2	d: 4	f: 3
u3	g: 4	h: 2	e: 3	f: 3
u4	b: 2	g: 4	h: 4	e: 4
u	a: 5	g: 4	e: 4	

Recommend d (and f, h)

CS 245

Notes11

91

## Sequence Mining

- Given a set of transcripts, use  $Pr[x|a]$  to predict if x is a good recommendation given user has taken a.
- Two issues...

CS 245

Notes11

92

## $Pr[x|a]$ Not Quite Right

transcript	containing
1	-
2	a
3	x
4	a → x
5	x → a

target user's transcript:  
[ ... a ... || unknown ]

recommend x?

$$Pr[x|a] = 2/3$$

$$Pr[x|a \sim x] = 1/2$$

CS 245

Notes11

93

## User Has Taken $\geq 1$ Course

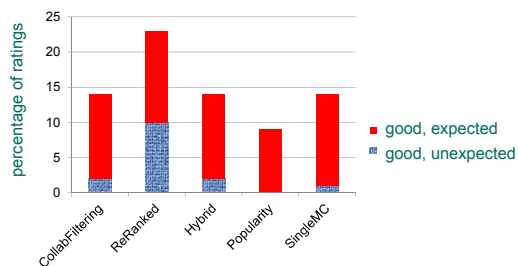
- User has taken  $T = \{a, b, c\}$
- Need  $Pr[x|T \sim x]$
- Approximate as  $Pr[x|a \sim x \wedge b \sim x \wedge c \sim x]$
- Expensive to compute, so...

CS 245

Notes11

94

## CourseRank User Study



CS 245

Notes11

95