

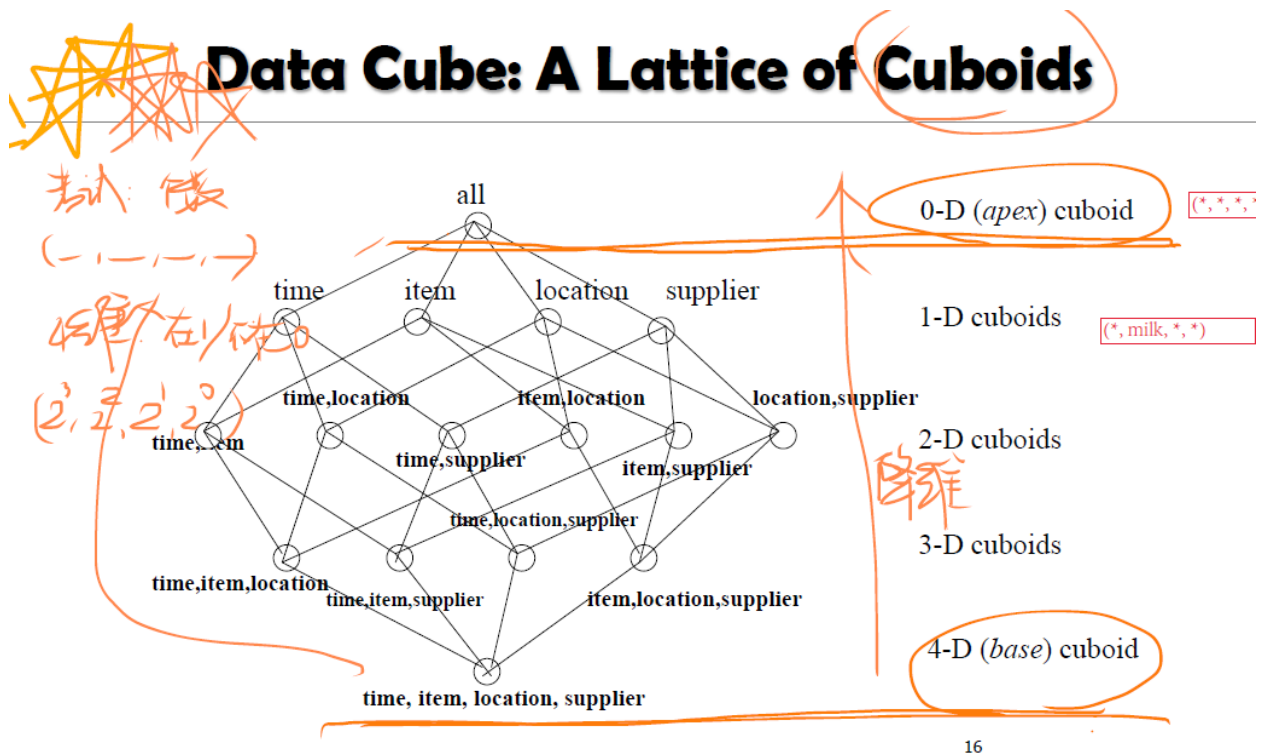
Chap 4 Data Warehousing and On-line Analytical Processing

- Table
 - Dimension Table
 - Fact Table: contains measures
- Spreadsheet
- Data Cube (P15-16)
 - Base Cuboid: n-D base cubes
 - Apex Cuboid: highest-level of summarization

1. Data Warehouse: A multi-dimensional model of a data warehouse

1) A Data Cube consists of dimensions & measures: (P15-16)

- Base Cuboid: n-D base cube
- Apex Cuboid: highest-level of summarization



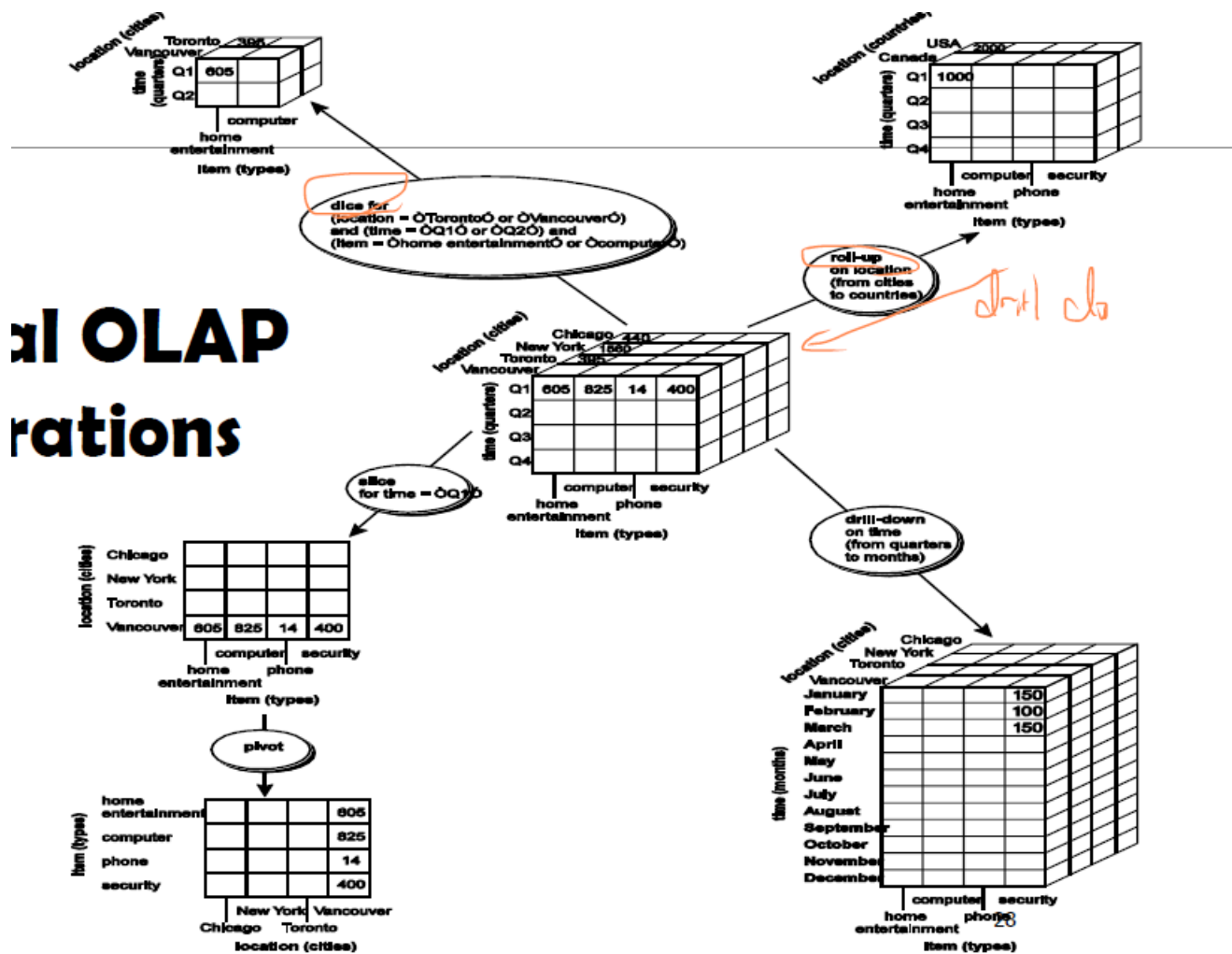
2) Star schema, snowflake schema, fact constellations (P17-20)

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

- Star Schema: 有一个中心的fact table连接着很多dimension tables
- Snowflake Schema: 连接的dimension tables有hierarchy
- Fact Constellation/ Galaxy Schema / Fact Constellation: 有多个fact table

3) OLAP operations: drilling, rolling, slicing, dicing and pivoting

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



- Roll Up/ Drill Up: summary并且降维
- Drill Down/ Roll Down: 增加维度
- Slice and Dice: slice 切片就是固定某个维度的值, dice是类似从多个维度挑选几个值
- Pivot / Rotate
- Drill Across
- Drill Through

4) Measures: Distributive, Algebraic, Holistic

Data Cube Measures: Three Categories

- **Distributive**: if the result derived by applying the function to n aggregate values is the same as that derived by applying the function on all the data without partitioning

□ E.g., count(), sum(), min(), max()

- **Algebraic**: if it can be computed by an algebraic function with M arguments (where M is a bounded integer), each of which is obtained by applying a distributive aggregate function

□ $avg(x) = sum(x) / count(x)$

□ Is min_N() an algebraic measure?

- **Holistic**: if there is no constant bound on the storage size needed to describe a subaggregate.

□ E.g., median(), mode(), rank()

2. Data Warehouse: Architecture, Design and Usage

1) Multi-tiered architecture

2) Business analysis design framework

3) Information processing, analytical processing, data mining

3. Implementation: Efficient computation of data cubes

1) Partial Materialization vs. Full Materialization vs. No Materialization

- 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?

- Materialization of data cube

□ **Full materialization**: Materialize every (cuboid)

□ **No materialization**: Materialize none (cuboid)

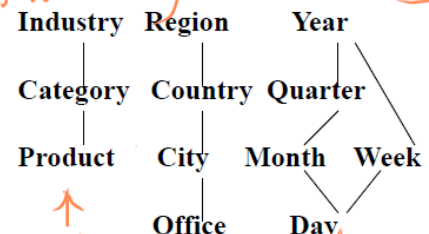
□ **Partial materialization**: Materialize some cuboids

□ Which cuboids to materialize?

□ Selection based on size, sharing, access frequency, etc.

Why this formula?

$$T = \prod_{i=1}^n (L_i + 1)$$



2) Indexing OLAP data: Bitmap index and join index

3) OLAP query processing

Chap 5 Data Cube Technology

1. Data Cube Computation: Preliminary Concepts

- Base Cell vs Aggregate Cell
- Ancestor Cell vs Descendant Cell
- Parent Cell vs Child Cell
- Full Cube vs Iceberg Cube
- **Close Cube & Close Cell**
- Cube Shell

2. Data Cube Computation Methods

- MultiWay Array Aggregation — small number of dimensions

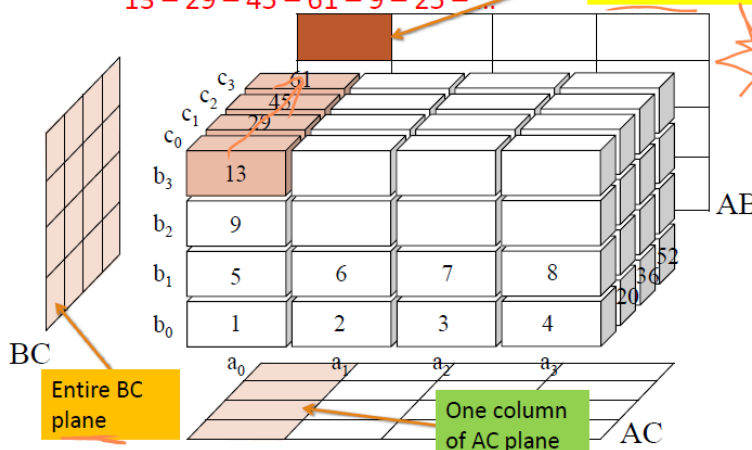
Cube Computation: Multi-Way Array Aggregation (MOLAP)

- Reducing memory and I/O

- Suppose we scan using order:

13 – 29 – 45 – 61 – 9 – 25 – ...

One chunk of AB plane



Example:
A: 4000, B: 400, C: 40
Chunk:
1000 x 100 x 10

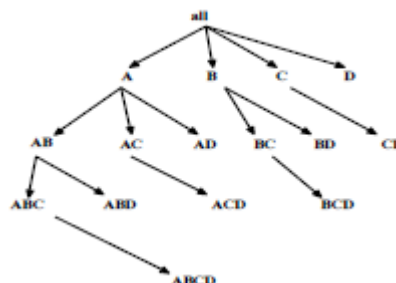
memory: *Minimizing*

- 40×400 (BC) + 40×1000 (AC) + 100×1000 (AB) = 156,000 units

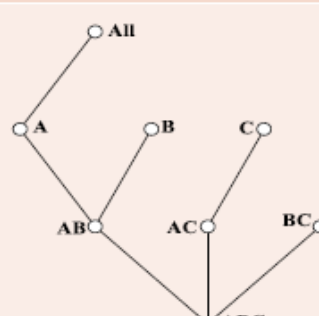
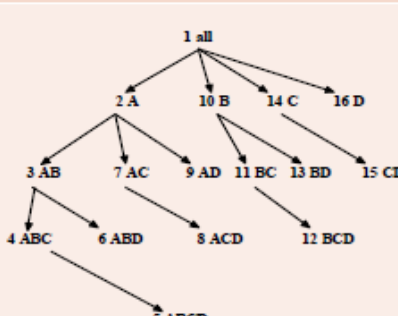
- Keep the smallest plane in main memory, fetch and compute only one chunk at a time for the largest plane

- The planes should be sorted and computed according to their size in ascending order

- BUC — large number of dimensions



MultiWay VS BUC

	multiway	BUC
Input format	Multi-dimensional array	Relational database
Good for	Full cube	Iceberg cube
Key idea	Simultaneously Aggregation	Partition and sort
Calculation direction		

- High-Dimensional OLAP with Shell-Fragments

Shell Fragment Cubes: Ideas

- Generalize the **1-D** inverted indices to **multi-dimensional** ones in the data cube sense
- Compute all cuboids for data cubes ABC and DE while retaining the inverted indices
 - Ex. shell fragment cube ABC contains 7 cuboids:
 - A, B, C; AB, AC, BC; ABC
- This completes the offline computation

ID_Measure Table

- If measures other than count are present, store in *ID_measure* table separate from the shell fragments

tid	count	sum
1	5	70
2	3	10
3	8	20
4	5	40
5	2	30

Shell-fragment AB

Attribute Value	TID List	List Size
a1	1 2 3	3
a2	4 5	2
b1	1 4 5	3
b2	2 3	2
c1	1 2 3 4 5	5
d1	1 3 4 5	4
d2	2	1
e1	1 2	2
e2	3 4	2
e3	5	1

Cell	Intersection	TID List	List Size
a1 b1	1 2 3 \cap 1 4 5	1	1
a1 b2	1 2 3 \cap 2 3	2 3	2
a2 b1	4 5 \cap 1 4 5	4 5	2
a2 b2	4 5 \cap 2 3	ϕ	0

3. Multidimensional Data Analysis in Cube Space

- Multi-feature Cubes
- Discover-Driven Exploration of Data Cubes

Chap 6 Mining Frequent Patterns, Association and Correlations

1. Basic Concepts

1) Pattern Discovery

2) Basic Concepts: Frequent Patterns and Association Rules

3) Compressed Representation: Closed Patterns and Max-Patterns

2. Efficient Pattern Mining Methods:

1) The Downward Closure Property of Frequent Patterns

2) The Apriori Algorithm

C_k : Candidate itemset of size k

F_k : Frequent itemset of size k

$K := 1$;

$F_k := \{\text{frequent items}\}$; // frequent 1-itemset

While ($F_k \neq \emptyset$) **do** { // when F_k is non-empty

$C_{k+1} := \text{candidates generated from } F_k$; // candidate generation

 Derive F_{k+1} by counting candidates in C_{k+1} with respect to TDB at minsup;

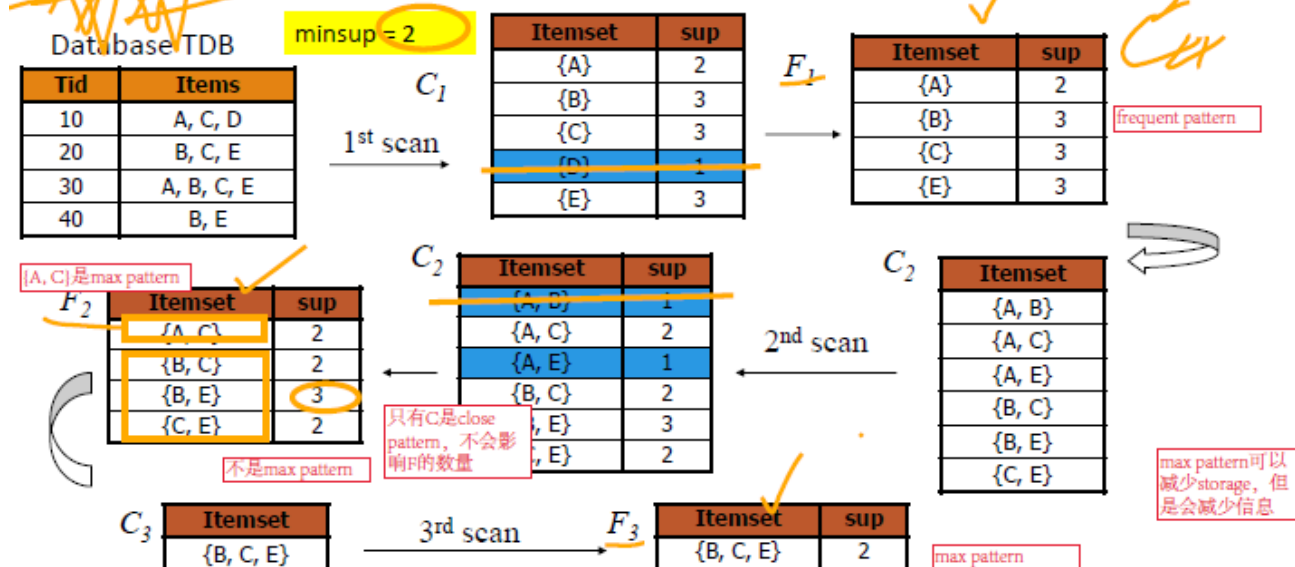
$k := k + 1$

}

return $\cup_k F_k$ // return F_k generated at each level

The Apriori Algorithm—An Example

从小的开始，然后不断增加set



3) Extensions or Improvements of Apriori

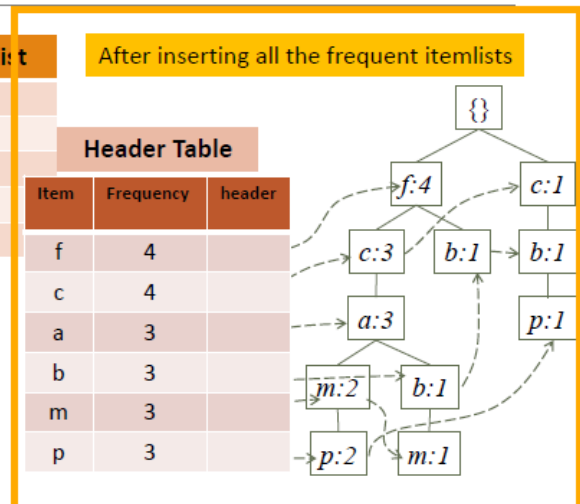
4) Mining Frequent Patterns by Exploring Vertical Data Format

5) FP Growth A Frequent Pattern-Growth Approach

Example: Construct FP-tree from a Transaction DB

TID	Items in the Transaction	Ordered, frequent itemlist
100	{f, a, c, d, g, i, m, p}	f, c, a, m, p
200	{a, b, c, f, l, m, o}	f, c, a, b, m
300	{b, f, h, j, o, w}	f, b
400	{b, c, k, s, p}	c, b, p
500	{a, f, c, e, l, p, m, n}	f, c, a, m, p

4. For each transaction, insert the ordered frequent itemlist into an FP-tree, with shared sub-branches merged, counts accumulated



memory会减少

6) Mining Closed Patterns

3. Pattern Evaluation

1) Interestingness Measures in Pattern Mining: Lift

- Measure of dependent/correlated events: lift

$$(400/1000) / (600/1000 * 750/1000)$$

Lift is more telling than s & c

$$lift(B, C) = \frac{c(B \rightarrow C)}{s(C)} = \frac{s(B \cup C)}{s(B) \times s(C)}$$

- Lift(B, C) may tell how B and C are correlated

- Lift(B, C) = 1: B and C are independent
- > 1: positively correlated
- < 1: negatively correlated

- For our example,
- $$lift(B, C) = \frac{400/1000}{600/1000 \times 750/1000} = 0.89$$
- $$lift(B, \neg C) = \frac{200/1000}{600/1000 \times 250/1000} = 1.33$$

- Thus, B and C are negatively correlated since lift(B, C) < 1;
- B and ¬C are positively correlated since lift(B, ¬C) > 1

	B	¬B	Σ _{row}
C	400	350	750
¬C	200	50	250
Σ _{col}	600	400	1000

- Another measure to test correlated events: χ^2

$$\chi^2 = \sum \frac{(Observed - Expected)^2}{Expected}$$

- For the table on the right,

	B	$\neg B$	Σ_{row}
C	400 (450)	350 (300)	750
$\neg C$	200 (150)	50 (100)	250
Σ_{col}	600	400	1000

$$\chi^2 = \frac{(400 - 450)^2}{450} + \frac{(350 - 300)^2}{300} + \frac{(200 - 150)^2}{150} + \frac{(50 - 100)^2}{100} = 55.56$$

Expected value

Observed value

- By consulting a table of critical values of the χ^2 distribution, one can conclude that the chance for B and C to be independent is very low (< 0.01)
- χ^2 -test shows B and C are negatively correlated since the expected value is 450 but the observed is only 400
- Thus, χ^2 is also more telling than the support-confidence framework

2) Interestingness Measures: Lift and χ^2

3) Null-Invariant Measures

Imbalance Ratio with Kulczynski Measure

- IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications:

$$IR(A, B) = \frac{|s(A) - s(B)|}{s(A) + s(B) - s(A \cup B)}$$

4) Comparison of Interestingness Measures