

Last Lecture

- Frequent itemset mining
 - Level-wise, join-based approach: Apriori (Agrawal & Srikant@VLDB'94)
 - Direct hashing and pruning: DHP (Park, Chen, Yu@SIGMOD'95)
 - Vertical data format approach: Eclat (Zaki, Parthasarathy, Ogihara, Li@KDD'97)
 - Frequent pattern projection and growth: FPgrowth (Han, Pei, Yin @SIGMOD'00)
- Closed itemset mining
 - Pattern growth-based approach: CLOSET+ (Wang et al. @KDD'03)

Learning Goals in Last Lecture

- Describe DHP, Eclat, FP Growth, and CLOSET+
- Implement FP Growth
 - Solve the *frequent itemset mining* problem **by hand** if the database is small, say, 10 transactions
 - Solve the problem by **programming** given an arbitrary size of transaction database and minimum support

Exercise: FP-Growth

Transaction ID	Items Bought
T ₁	{Mango, Onion, Nintendo, Key-chain, Eggs, Yo-yo}
T ₂	{Doll, Onion, Nintendo, Key-chain, Eggs, Yo-yo}
T ₃	{Mango, Apple, Key-chain, Eggs}
T ₄	{Mango, Umbrella, Corn, Key-chain, Yo-yo}
T ₅	{Corn, Onion, Onion, Key-chain, Ice-cream, Eggs}

Transaction ID	Items Bought
T ₁	MONKEY
T ₂	DONKEY
T ₃	MAKE
T ₄	MUCKY
T ₅	COKIE

if min_sup = 60%



Chapter 6. Frequent Pattern Mining: Pattern Evaluation

Meng Jiang
Data Science

How to Judge if a Rule/Pattern Is Interesting?

- Pattern/association mining will generate a large set of patterns/rules
 - Not all the generated patterns/association rules are interesting

How to Judge if a Rule/Pattern Is Interesting? (cont.)

- Interestingness measures: Subjective vs. Objective
 - Subjective interestingness measures: One person's trash could be another person's treasure
 - Query-based: Relevant to a user's particular request
 - Against one's knowledge-base: unexpected, freshness, timeliness
 - Visualization tools: Multi-dimensional, interactive examination

How to Judge if a Rule/Pattern Is Interesting? (cont.)

- Interestingness measures: Subjective vs. Objective
 - Subjective interestingness measures: One person's trash could be another person's treasure
 - Query-based: Relevant to a user's particular request
 - Against one's knowledge-base: unexpected, freshness, timeliness
 - Visualization tools: Multi-dimensional, interactive examination
 - **Objective interestingness measures**
 - **Support, confidence, correlation, ...**

Judge an Association Rule

- Are s and c interesting in association rules:
 - “ $A \Rightarrow B$ ” [s, c]?
- Example: Suppose one school may have the following statistics on # of students who may play basketball and/or eat cereal:

2-way contingency table

	play-basketball	not play-basketball	sum (row)
eat-cereal	400	350	750
not eat-cereal	200	50	250
sum(col.)	600	400	1000

Limitation of the Support-Confidence Framework

- Association rule mining may generate the following:
 - play-basketball \Rightarrow eat-cereal [40%, 66.7%] (higher s & c)
- But this strong association rule is misleading: The overall % of students eating cereal is 75% $>$ 66.7%, a more telling rule:
 - \neg play-basketball \Rightarrow eat-cereal [35%, 87.5%] (high s & higher c)

2-way contingency table

	play-basketball	not play-basketball	sum (row)
eat-cereal	400	350	750
not eat-cereal	200	50	250
sum(col.)	600	400	1000

Interestingness Measure: Lift

- Measure of dependent/correlated events: **lift**

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

Symmetric?

- 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

Interestingness Measure: Lift (cont.)

- For our example,

	B	$\neg B$	Σ_{row}
C	400	350	750
$\neg C$	200	50	250
$\Sigma_{\text{col.}}$	600	400	1000

- Thus, B and C are negatively correlated since $\text{lift}(B, C) < 1$;
 - B and $\neg C$ are positively correlated since $\text{lift}(B, \neg C) > 1$
 - C and $\neg B$ are positively correlated since $\text{lift}(C, \neg B) > 1$

$$\text{lift}(B, C) = \frac{400/1000}{600/1000 \times 750/1000} = 0.89$$

$$\text{lift}(B, \neg C) = \frac{200/1000}{600/1000 \times 250/1000} = 1.33$$

$$\begin{aligned}\text{lift}(\neg B, C) &= \text{lift}(C, \neg B) \\ &= 0.35/0.4/0.75 = 1.17\end{aligned}$$

Interestingness Measure: χ^2

- Another measure to test correlated events: χ^2

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

- General rules
 - $\chi^2 = 0$: independent
 - $\chi^2 > 0$: correlated, either positive or negative, so it needs additional test

Interestingness Measure: χ^2 (cont.)

- Now,

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

- χ^2 shows B and C are *negatively correlated* since the expected value is 450 but the observed is only 400

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

Lift and χ^2 : Are They Always Good Measures?

- Null transactions: Transactions that contain neither B nor C
- Let's examine the dataset D
 - BC (100) is much rarer than B¬C (1000) and ¬BC (1000), but there are many ¬B¬C (100000)
 - Unlikely B & C will happen together!

	B	$\neg B$	Σ_{row}
C	100	1000	1100
$\neg C$	1000	100000	101000
$\Sigma_{\text{col.}}$	1100	101000	102100

null transactions

Lift and χ^2 : Are They Always Good Measures?

- But, $lift(B, C) = 8.44 \gg 1$
(Lift shows B and C are strongly positively correlated!)
- $\chi^2 = 670$: Observed(BC) \gg expected value (11.85)
- Too many null transactions may “spoil the soup”!

	B	$\neg B$	Σ_{row}
C	100	1000	1100
$\neg C$	1000	100000	101000
$\Sigma_{\text{col.}}$	1100	101000	102100

null transactions

Contingency table with expected values added

	B	$\neg B$	Σ_{row}
C	100 (11.85)	1000	1100
$\neg C$	1000 (988.15)	100000	101000
$\Sigma_{\text{col.}}$	1100	101000	102100

Interestingness Measures & Null-Invariance

- *Null invariance*: Value does not change with the # of null-transactions
- A few interestingness measures: Some are null invariant

Measure	Definition	Range	Null-Invariant
$\chi^2(A, B)$	$\sum_{i,j=0,1} \frac{(e(a_i b_j) - o(a_i b_j))^2}{e(a_i b_j)}$	$[0, \infty]$	No
$Lift(A, B)$	$\frac{s(A \cup B)}{s(A) \times s(B)}$	$[0, \infty]$	No
$AllConf(A, B)$	$\frac{s(A \cup B)}{\max\{s(A), s(B)\}}$	$[0, 1]$	Yes
$Jaccard(A, B)$	$\frac{s(A \cup B)}{s(A) + s(B) - s(A \cup B)}$	$[0, 1]$	Yes
$Cosine(A, B)$	$\frac{s(A \cup B)}{\sqrt{s(A) \times s(B)}}$	$[0, 1]$	Yes
$Kulczynski(A, B)$	$\frac{1}{2} \left(\frac{s(A \cup B)}{s(A)} + \frac{s(A \cup B)}{s(B)} \right)$	$[0, 1]$	Yes
$MaxConf(A, B)$	$\max\left\{ \frac{s(A)}{s(A \cup B)}, \frac{s(B)}{s(A \cup B)} \right\}$ $\max\{ s(A \cup B) / s(A), s(A \cup B) / s(B) \}$	$[0, 1]$	Yes

χ^2 and lift are not null-invariant

Jaccard, consine, AllConf, MaxConf, and Kulczynski are null-invariant measures

Null Invariance: An Important Property

- Why is null invariance crucial for the analysis of massive transaction data?
 - Many transactions may contain neither milk nor coffee!

milk vs. coffee contingency table

	<i>milk</i>	$\neg\text{milk}$	Σ_{row}
<i>coffee</i>	<i>mc</i>	$\neg\text{mc}$	<i>c</i>
$\neg\text{coffee}$	<i>m</i> $\neg\text{c}$	$\neg\text{m}$ $\neg\text{c}$	$\neg\text{c}$
Σ_{col}	<i>m</i>	$\neg\text{m}$	Σ

- ❑ Lift and χ^2 are not null-invariant: not good to evaluate data that contain too many or too few null transactions!
- ❑ Many measures are not null-invariant!

Null-transactions
w.r.t. m and c

Data set	<i>mc</i>	$\neg\text{mc}$	<i>m</i> $\neg\text{c}$	$\neg\text{m}$ $\neg\text{c}$	χ^2	Lift
D_1	10,000	1,000	1,000	100,000	90557	9.26
D_2	10,000	1,000	1,000	100	0	1
D_3	100	1,000	1,000	100,000	670	8.44
D_4	1,000	1,000	1,000	100,000	24740	25.75
D_5	1,000	100	10,000	100,000	8173	9.18
D_6	1,000	10	100,000	100,000	965	1.97

Comparison of Null-Invariant Measures

- Not all null-invariant measures are created equal
- Which one is better?
 - D_4 — D_6 differentiate the null-invariant measures
 - Kulc (Kulczynski 1927) holds firm and is in balance of both directional implications

2-variable contingency table

	<i>milk</i>	$\neg\text{milk}$	Σ_{row}
<i>coffee</i>	<i>mc</i>	$\neg\text{mc}$	<i>c</i>
$\neg\text{coffee}$	<i>m</i> $\neg\text{c}$	$\neg\text{m}$ $\neg\text{c}$	$\neg\text{c}$
Σ_{col}	<i>m</i>	$\neg\text{m}$	Σ

All 5 are null-invariant

Data set	<i>mc</i>	$\neg\text{mc}$	<i>m</i> $\neg\text{c}$	$\neg\text{m}$ c	AllConf	Jaccard	Cosine	Kulc	MaxConf
D_1	10,000	1,000	1,000	100,000	0.91	0.83	0.91	0.91	0.91
D_2	10,000	1,000	1,000	100	0.91	0.83	0.91	0.91	0.91
D_3	100	1,000	1,000	100,000	0.09	0.05	0.09	0.09	0.09
D_4	1,000	1,000	1,000	100,000	0.5	0.33	0.5	0.5	0.5
D_5	1,000	100	10,000	100,000	0.09	0.09	0.29	0.5	0.91
D_6	1,000	10	100,000	100,000	0.01	0.01	0.10	0.5	0.99

Subtle: They disagree on those cases

Analysis of DBLP Coauthor Relationships

- Recent DB conferences, removing balanced associations, low sup, etc.

ID	Author A	Author B	$s(A \cup B)$	$s(A)$	$s(B)$	Jaccard	Cosine	Kulc
1	Hans-Peter Kriegel	Martin Ester	28	146	54	0.163 (2)	0.315 (7)	0.355 (9)
2	Michael Carey	Miron Livny	26	104	58	0.191 (1)	0.335 (4)	0.349 (10)
3	Hans-Peter Kriegel	Joerg Sander	24	146	36	0.152 (3)	0.331 (5)	0.416 (8)
4	Christos Faloutsos	Spiros Papadimitriou	20	162	26	0.119 (7)	0.308 (10)	0.446 (7)
5	Hans-Peter Kriegel	Martin Pfeifle	18	146	18	0.123 (6)	0.351 (2)	0.562 (2)
6	Hector Garcia-Molina	Wilbert Labio	16	144	18	0.110 (9)	0.314 (8)	0.500 (4)
7	Divyakant Agrawal	Wang Hsiung	16	120	16	0.133 (5)	0.365 (1)	0.567 (1)
8	Elke Rundensteiner	Murali Mani	16	104	20	0.148 (4)	0.351 (3)	0.477 (6)
9	Divyakant Agrawal	Oliver Po	12	120	12	0.100 (10)	0.316 (6)	0.550 (3)
10	Gerhard Weikum	Martin Theobald	12	106	14	0.111 (8)	0.312 (9)	0.485 (5)

Advisor-advisee relations

- Which pairs of authors are strongly related?
 - Use Kulc to find Advisor-advisee, close collaborators

Imbalance Ratio with Kulczynski Measure

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

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

- Kulczynski and Imbalance Ratio (IR) together present a clear picture for all the three datasets D_4 through D_6
 - D_4 is neutral & balanced; D_5 is neutral but imbalanced
 - D_6 is neutral but very imbalanced

Data set	mc	$\neg mc$	$m \neg c$	$\neg m \neg c$	Jaccard	Cosine	Kulc	IR
D_1	10,000	1,000	1,000	100,000	0.83	0.91	0.91	0
D_2	10,000	1,000	1,000	100	0.83	0.91	0.91	0
D_3	100	1,000	1,000	100,000	0.05	0.09	0.09	0
D_4	1,000	1,000	1,000	100,000	0.33	0.5	0.5	0
D_5	1,000	100	10,000	100,000	0.09	0.29	0.5	0.89
D_6	1,000	10	100,000	100,000	0.01	0.10	0.5	0.99

What Measures to Choose for Effective Pattern Evaluation?

- Null value cases are predominant in many large datasets
 - Neither milk nor coffee is in most of the baskets; neither Mike nor Jim is an author in most of the papers;
- Null-invariance is an important property
- Lift and χ^2 are good measures if null transactions are not predominant
 - Otherwise, Kulczynski + Imbalance Ratio should be used to judge the interestingness of a pattern
 - AllConf, Jaccard, Cosine, MaxConf...

Summary

- Basic Concepts:
 - Frequent Patterns, Association Rules, Closed Patterns and Max-Patterns
- Frequent Itemset Mining Methods
 - The Downward Closure Property and The Apriori Algorithm
 - Extensions or Improvements of Apriori
 - Mining Frequent Patterns by Exploring Vertical Data Format
 - FP-Growth: A Frequent Pattern-Growth Approach
 - Mining Closed Patterns
- Which Patterns Are Interesting?—Pattern Evaluation Methods
 - Interestingness Measures: Lift and χ^2
 - Null-Invariant Measures
 - Comparison of Interestingness Measures

References

- R. Agrawal, T. Imielinski, and A. Swami, "Mining association rules between sets of items in large databases", in Proc. of SIGMOD'93
- R. J. Bayardo, "Efficiently mining long patterns from databases", in Proc. of SIGMOD'98
- N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal, "Discovering frequent closed itemsets for association rules", in Proc. of ICDT'99
- J. Han, H. Cheng, D. Xin, and X. Yan, "Frequent Pattern Mining: Current Status and Future Directions", Data Mining and Knowledge Discovery, 15(1): 55-86, 2007
- R. Agrawal and R. Srikant, "Fast algorithms for mining association rules", VLDB'94
- A. Savasere, E. Omiecinski, and S. Navathe, "An efficient algorithm for mining association rules in large databases", VLDB'95
- J. S. Park, M. S. Chen, and P. S. Yu, "An effective hash-based algorithm for mining association rules", SIGMOD'95
- S. Sarawagi, S. Thomas, and R. Agrawal, "Integrating association rule mining with relational database systems: Alternatives and implications", SIGMOD'98
- M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li, "Parallel algorithm for discovery of association rules", Data Mining and Knowledge Discovery, 1997
- J. Han, J. Pei, and Y. Yin, "Mining frequent patterns without candidate generation", SIGMOD'00

References (cont.)

- M. J. Zaki and Hsiao, "CHARM: An Efficient Algorithm for Closed Itemset Mining", SDM'02
- J. Wang, J. Han, and J. Pei, "CLOSET+: Searching for the Best Strategies for Mining Frequent Closed Itemsets", KDD'03
- C. C. Aggarwal, M.A., Bhuiyan, M. A. Hasan, "Frequent Pattern Mining Algorithms: A Survey", in Aggarwal and Han (eds.): Frequent Pattern Mining, Springer, 2014
- C. C. Aggarwal and P. S. Yu. A New Framework for Itemset Generation. PODS'98
- S. Brin, R. Motwani, and C. Silverstein. Beyond market basket: Generalizing association rules to correlations. SIGMOD'97
- M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A. I. Verkamo. Finding interesting rules from large sets of discovered association rules. CIKM'94
- E. Omiecinski. Alternative Interest Measures for Mining Associations. TKDE'03
- P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the Right Interestingness Measure for Association Patterns. KDD'02
- T. Wu, Y. Chen and J. Han, Re-Examination of Interestingness Measures in Pattern Mining: A Unified Framework, Data Mining and Knowledge Discovery, 21(3):371-397, 2010

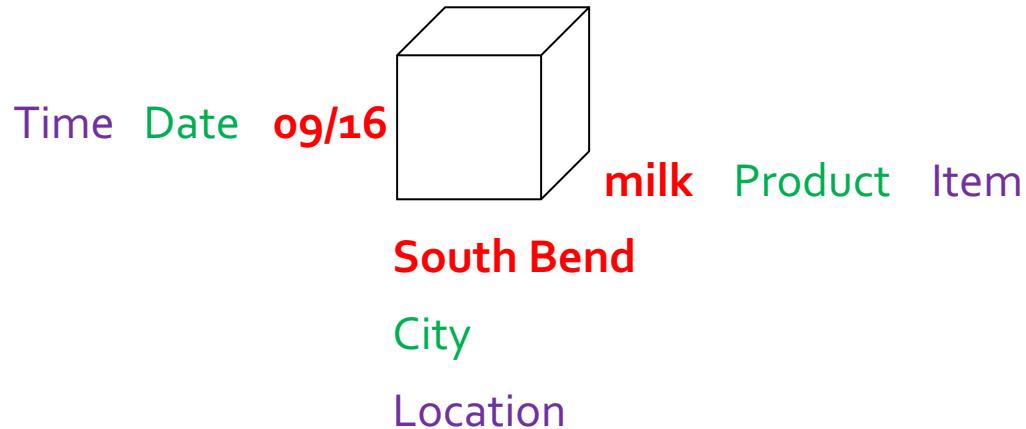


Chapter 4&5. Data Cube

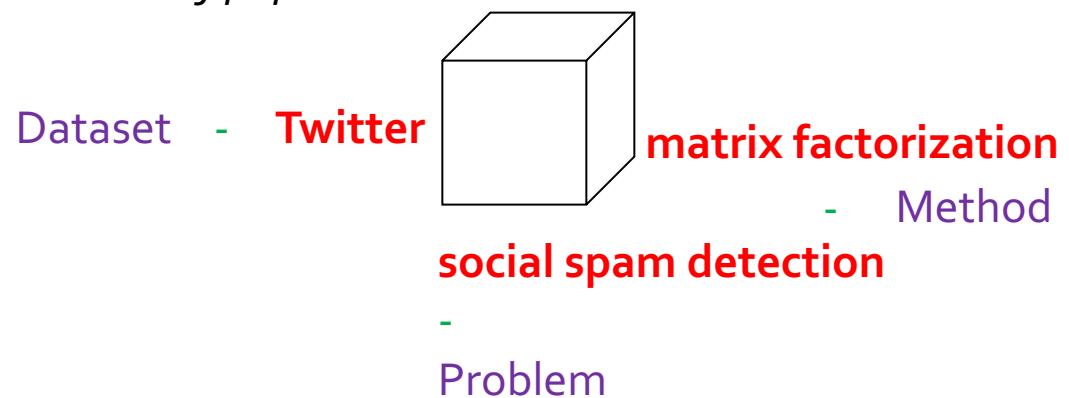
Meng Jiang
Data Science

Cells: Dimension, Dimension Level and Dimension Value

A cell of transactions:



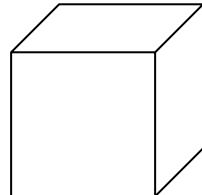
A cell of papers:



Cells: Dimension Level and Concept Hierarchy

A cell of transactions:

Time Date 09/16



milk Product Item

South Bend

City

Location

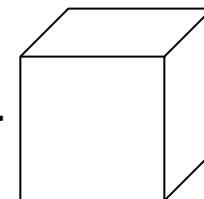
Time: Year-Quarter-Month-Week-Day

Location: Country-State-City-Street

Item: Department-Product-Model

A cell of papers:

Dataset - Twitter



matrix factorization

- Method

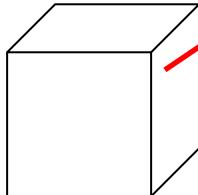
social spam detection

-

Problem

Cells: Facts or Measures

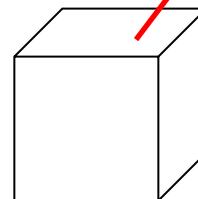
A *cell of transactions*:

Time	Date	09/16		milk	Product	Item
				South Bend		
				City		
				Location		

{TID45, TID137, TID451},
count=3,
dollars_sold=157

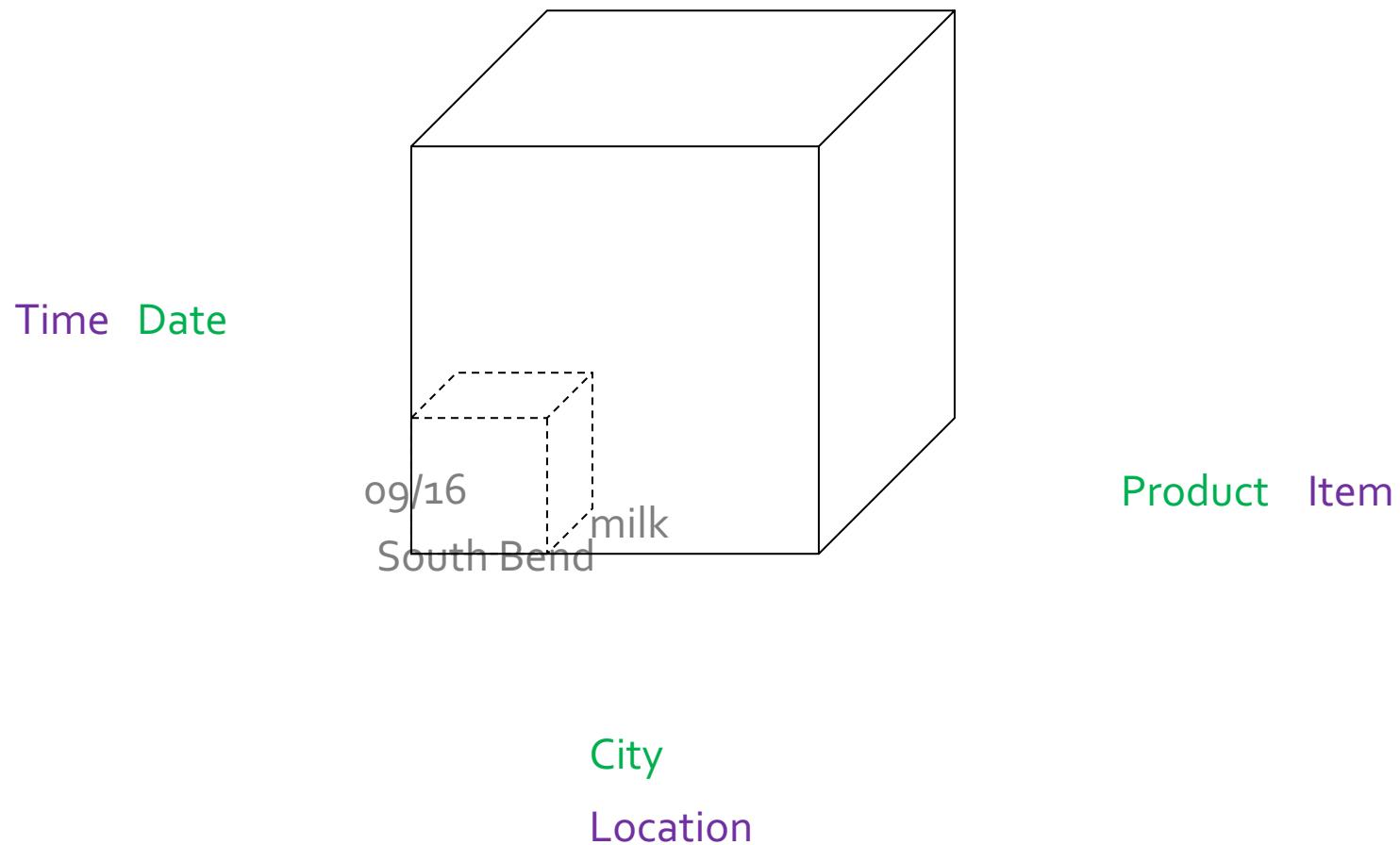
A *cell of papers*:

Dataset - Twitter


matrix factorization
- Method
social spam detection
-
Problem

{PID31, PID217},
count=2,
citations=3317

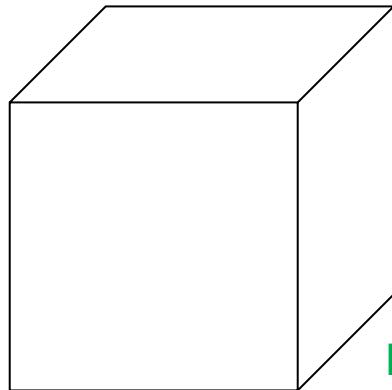
Cuboids: Dimension, Dimension Level



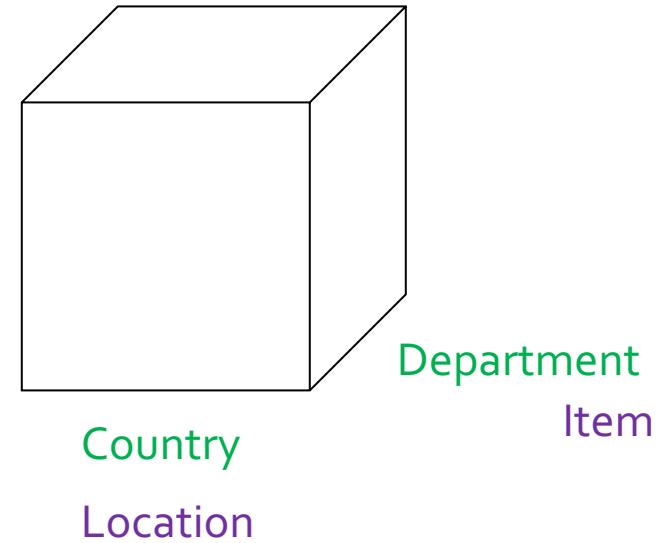
Base Cells and Aggregate Cells

- Suppose a cuboid has 3 dimensions (time, location, item) at specific dimension levels (date, city, product).
 - Base cells
 - (09/16, South Bend, milk)
 - Aggregate cells
 - (*, South Bend, milk)
 - (09/16, *, milk)
 - (09/16, South Bend, *)
 - (*, *, milk)
 - (*, South Bend, *)
 - (09/16, *, *)
 - (*, *, *), called the **Apex cell**
- parent vs child cells
ancestor vs descendant cells
sibling cell:
(09/16, Mishawaka, milk)

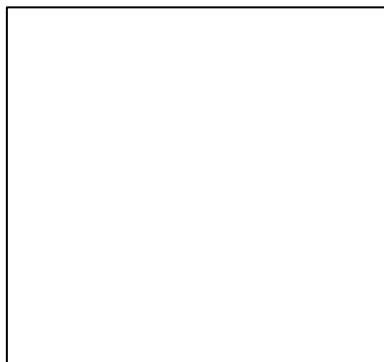
Base Cuboids and Aggregate Cuboids



(Date, City, Product)



(Month, Country, Department)

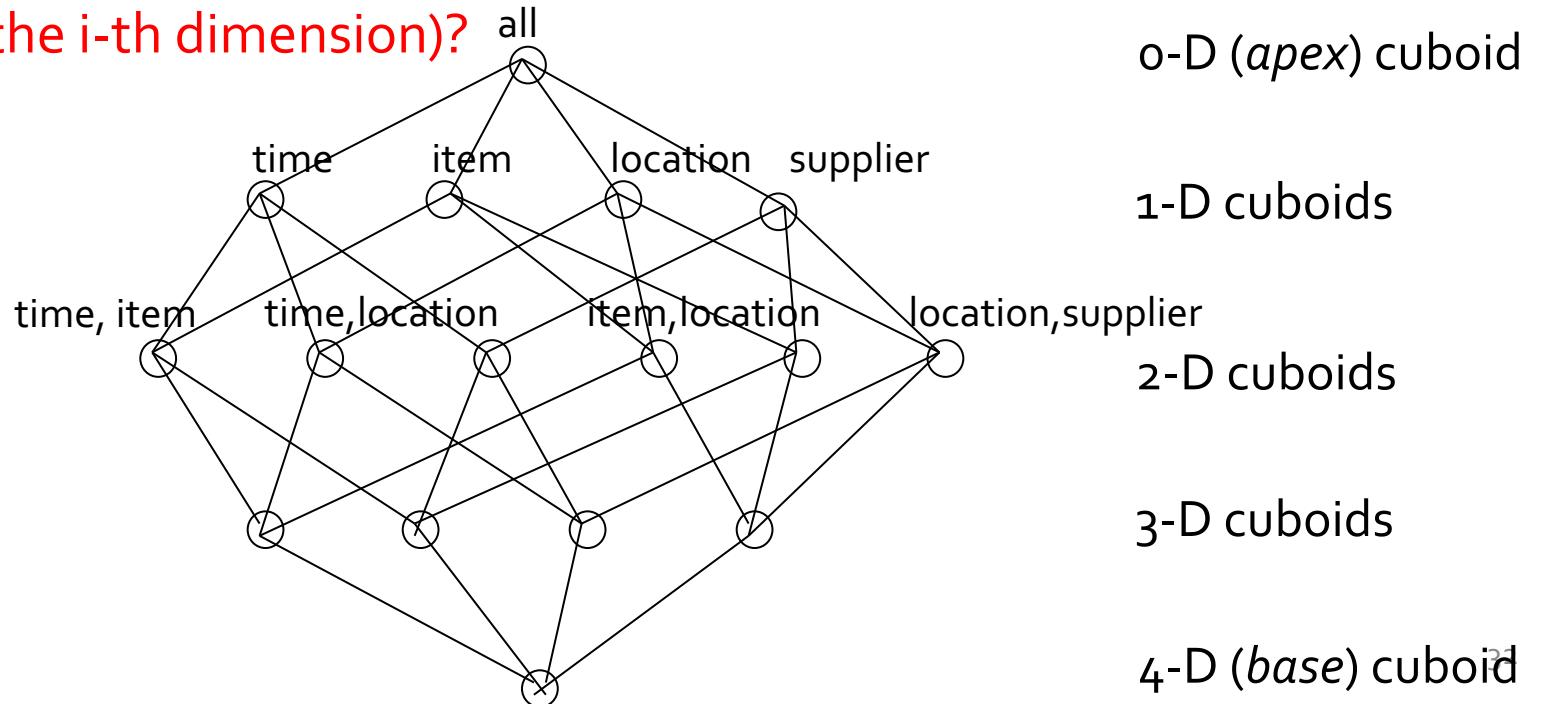


(Date, City, *)
Location

Apex cuboid: (*, *, *)

(N-Dimensional) Data Cube

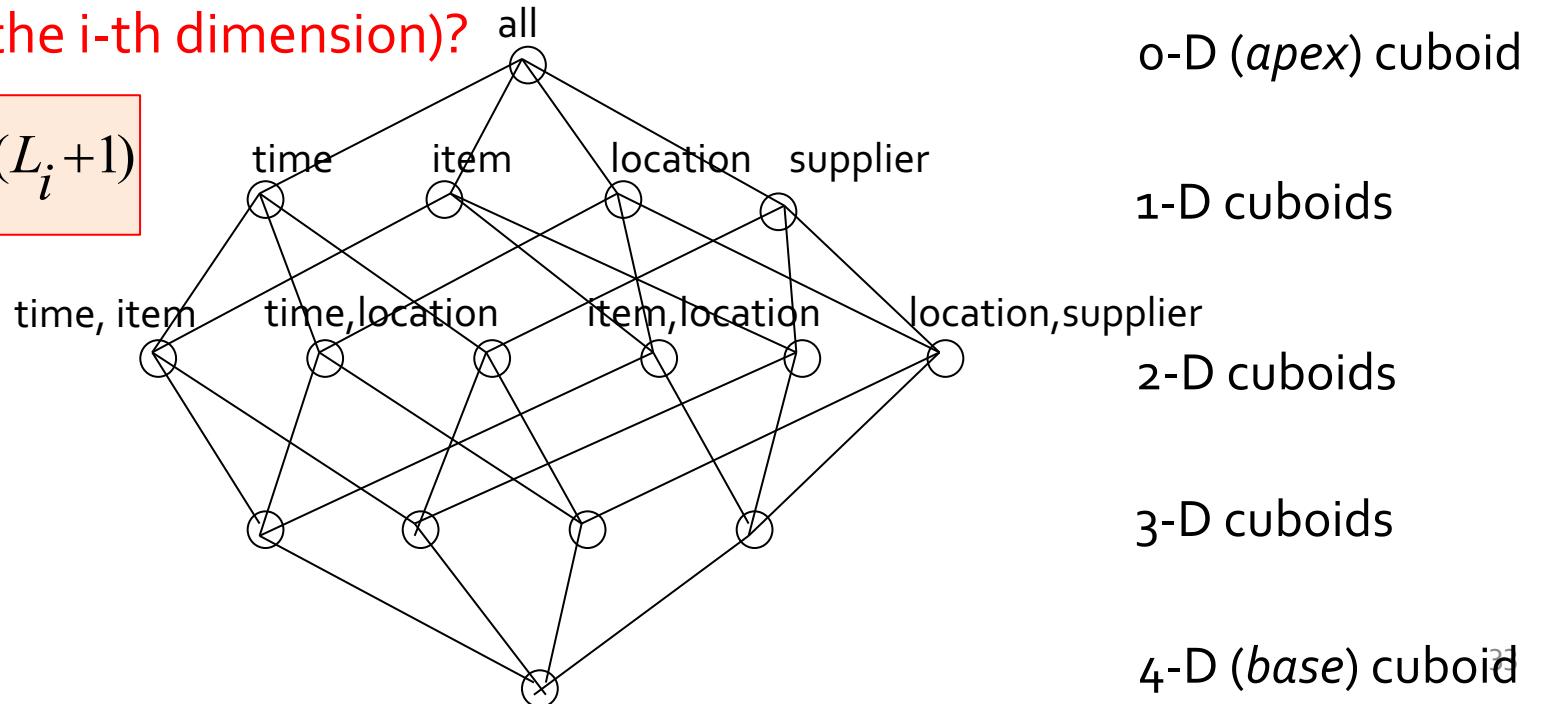
- Data cube can be viewed as a lattice of cuboids
 - The bottom-most cuboid is the base cuboid
 - The top-most cuboid (apex) contains only one cell
 - How many cuboids in an n-dimensional cube with L_i levels (at the i-th dimension)?



(N-Dimensional) Data Cube

- 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_i levels (at the i-th dimension)?

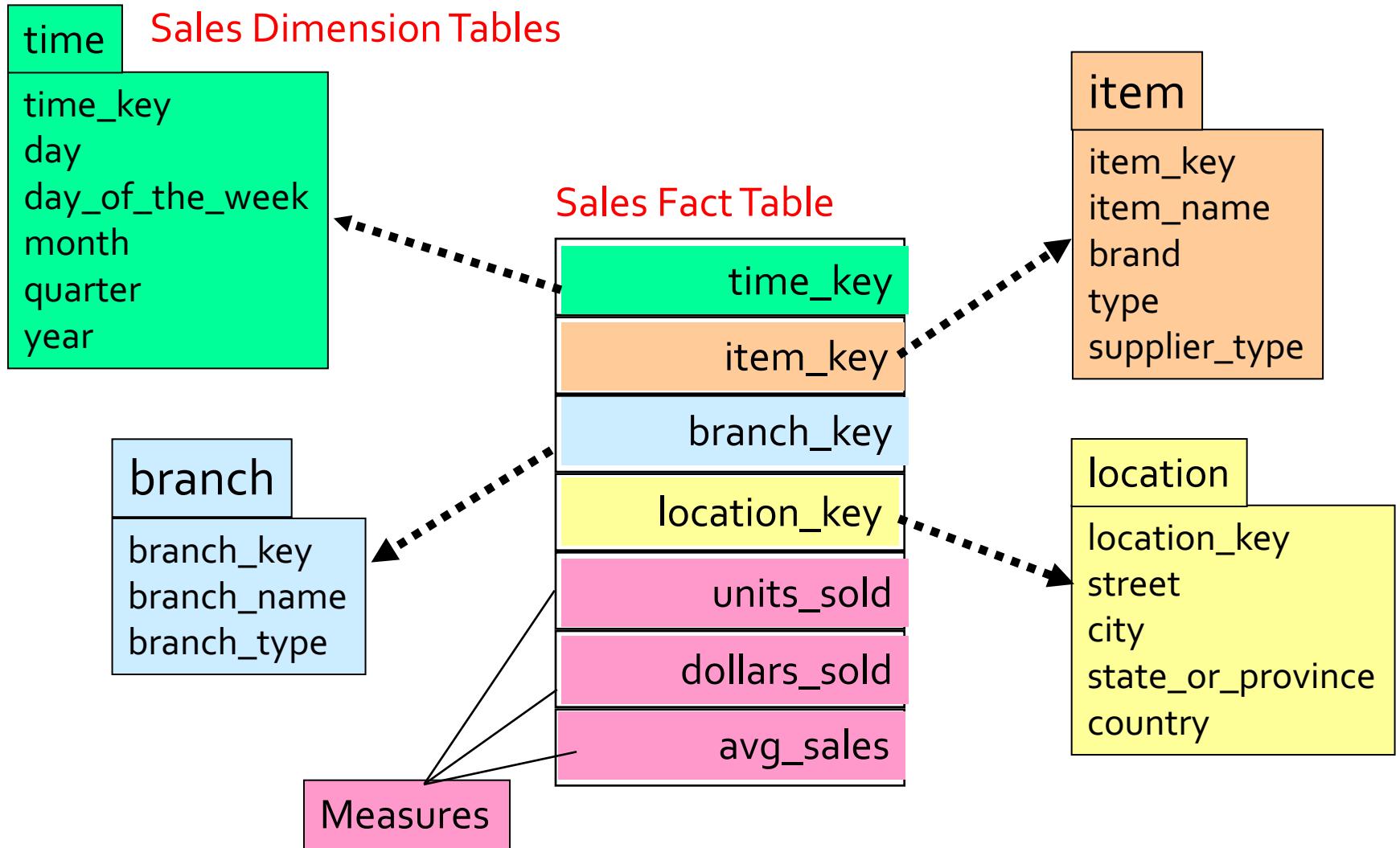
$$T = \prod_{i=1}^n (L_i + 1)$$



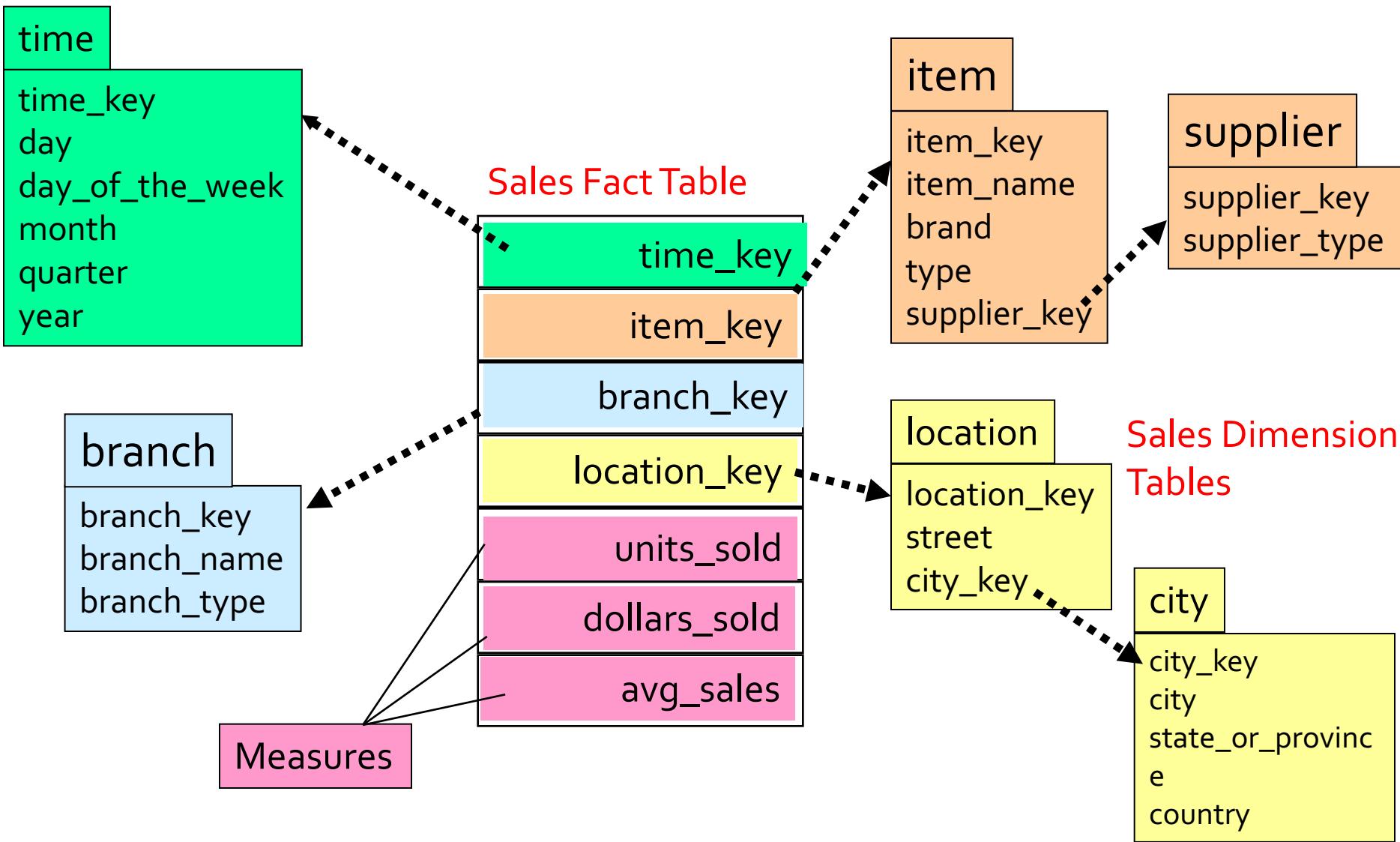
Data Cube: Definition

- **Data cube:** A lattice of cuboids
 - In data warehousing literature, an **n-D base cube** is called a **base cuboid**
 - The top most **0-D cuboid**, which holds the highest-level of summarization, is called the **apex cuboid**
 - The lattice of cuboids forms a **data cube**
- A data cube, such as sales, allows data to be modeled and viewed in multiple dimensions
 - **Dimension tables**, such as item (item_name, brand, type), or time (day, week, month, quarter, year)
 - **Fact table** contains **measures** (such as dollars_sold) and keys to each of the related dimension tables
 - **Schemas:** Dimension tables and Fact tables

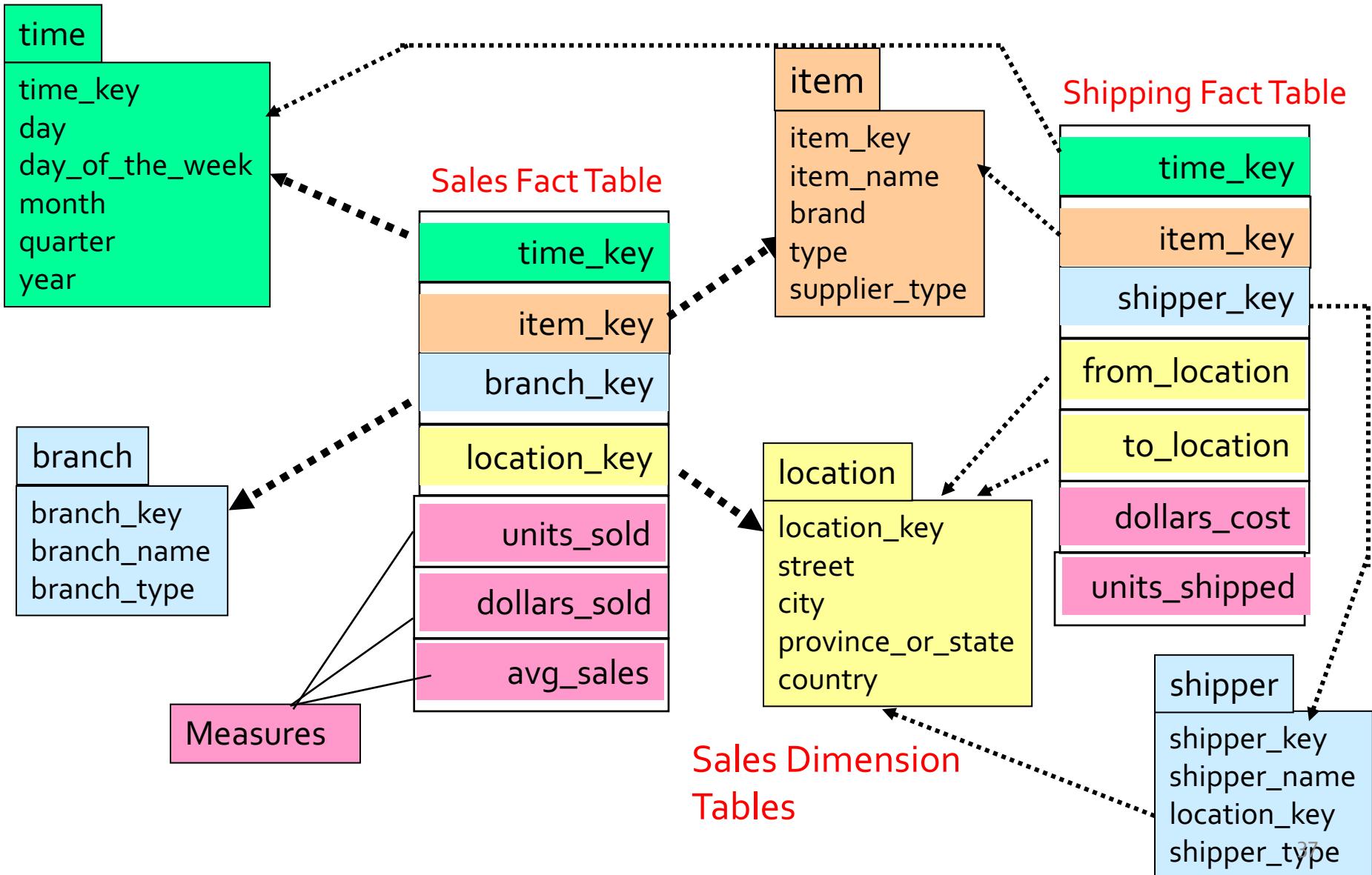
Star Schema



Snowflake Schema



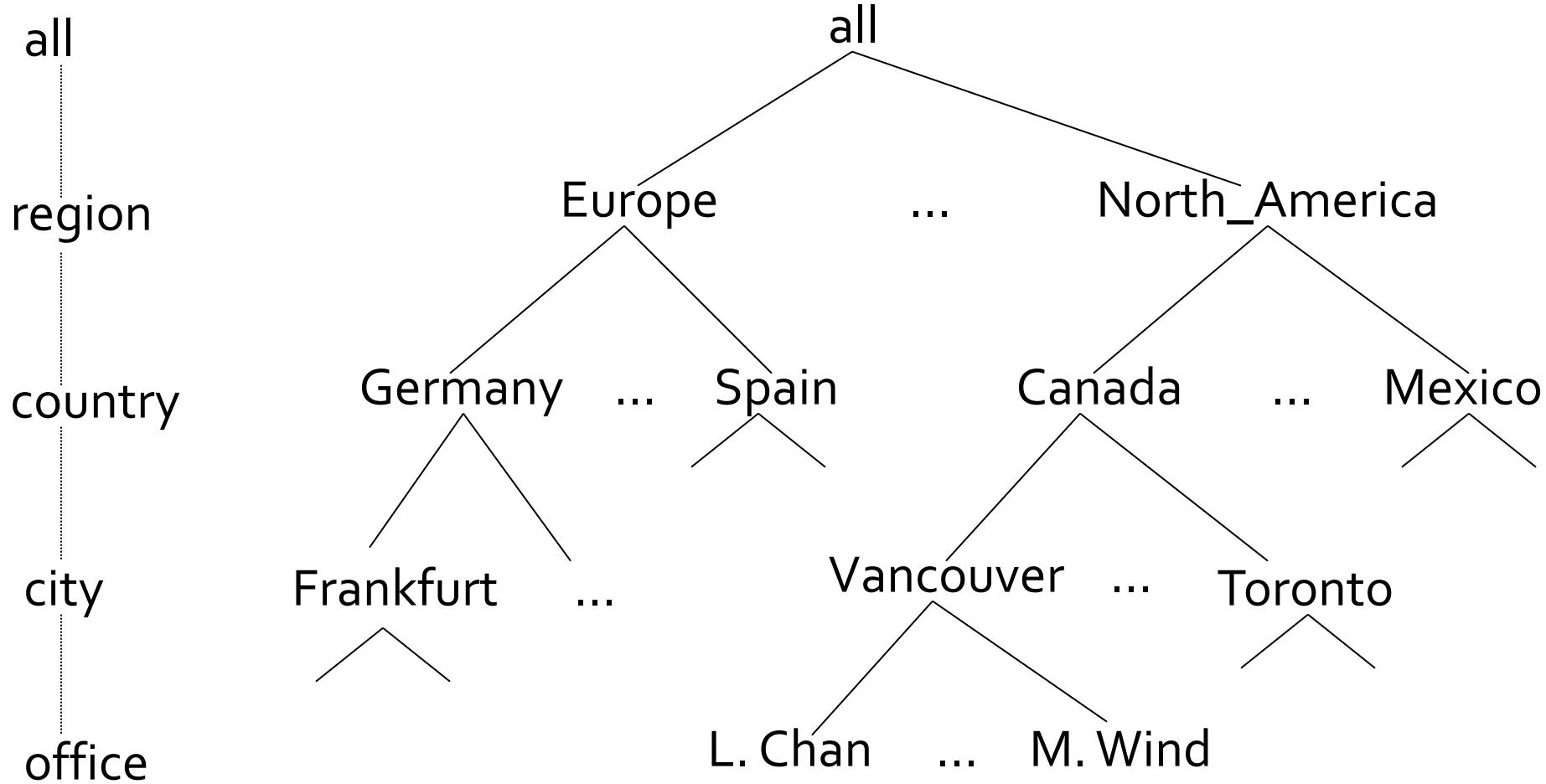
Fact Constellation



Modeling of Data Cubes

- Modeling data cubes: dimensions & measures
 - **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

Concept Hierarchy: Dimension Level and Dimension Value

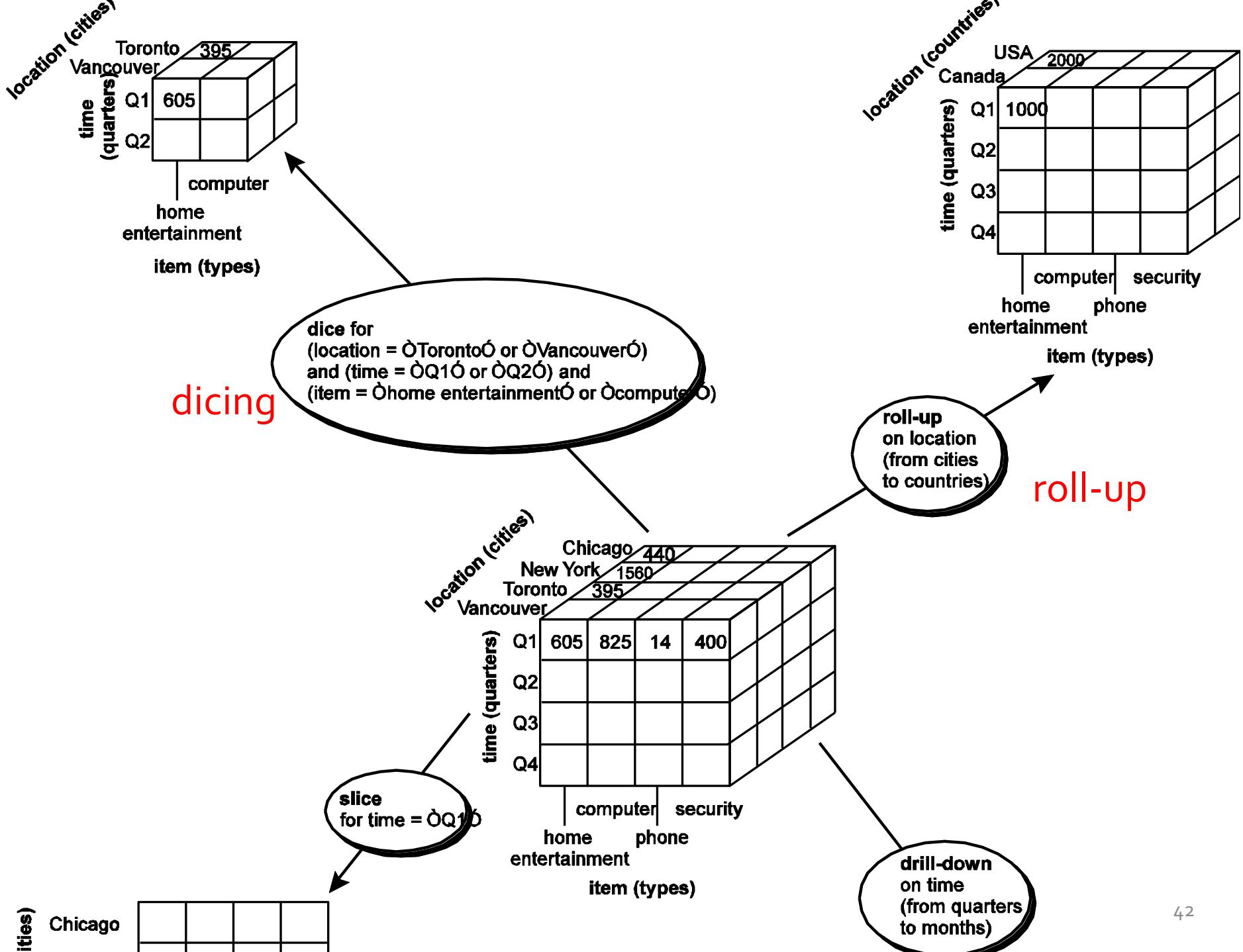


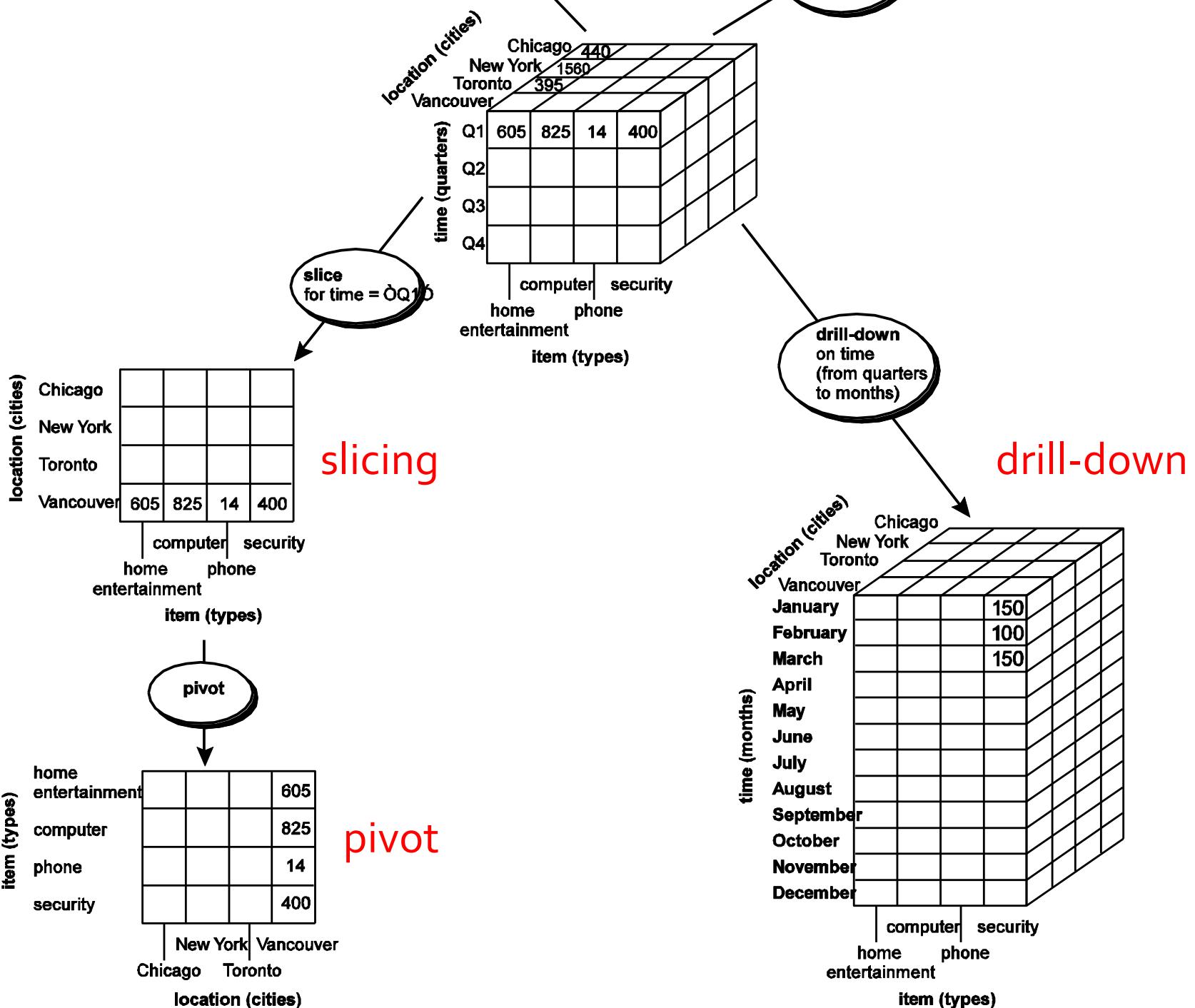
Data Cube Measures: Three Categories

- **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**
 - $\text{avg}(x) = \text{sum}(x) / \text{count}(x)$
- **Holistic**: if there is no constant bound on the storage size needed to describe a sub-aggregate.
 - E.g., `median()`, `mode()`, `rank()`
- Q: How about `standard_deviation()`, `Q1()`, `Q3()`?

Typical Data Cube Operations

- 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*



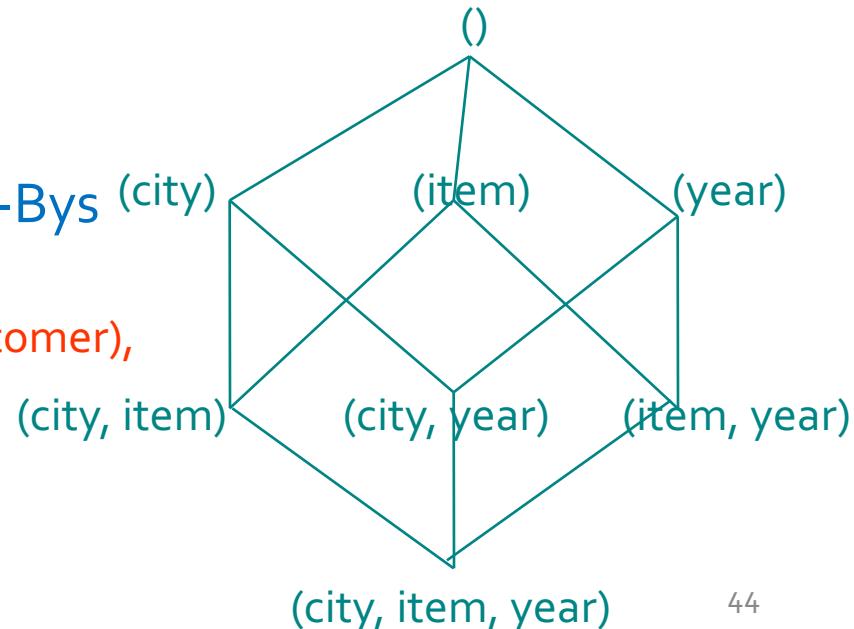


The “Compute Cube” Operator

- Cube definition and computation

```
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.'97**)

```
SELECT item, city, year, SUM (amount)  
FROM SALES  
CUBE BY item, city, year
```
- Need compute the following **Group-Bys**
 (city)
 $(\text{year, product, customer})$,
 (year, product) , (year, customer) , $(\text{product, customer})$,
 (year) , (product) , (customer)
 $()$



Data Cube History

Data cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals

2981

1997

J Gray, S Chaudhuri, A Bosworth, A Layman, D Reichart, M Venkatrao, ...

Data Mining and Knowledge Discovery 1 (1), 29-53

Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tab, and Sub-Totals

Jim Gray
Surajit Chaudhuri
Adam Bosworth
Andrew Layman
Don Reichart
Murali Venkatrao
Frank Pellow
Hamid Pirahesh¹

May 1997

Technical Report
MSR-TR-97-32

Microsoft Research
Microsoft Corporation
One Microsoft Way
Redmond, WA 98052

Surajit Chaudhuri is a computer scientist best known for his contributions to database management systems. He is currently a distinguished scientist at Microsoft Research, where he leads the Data Management, Exploration and Mining group.

Adam Bosworth is a former Vice President of Product Management at Google Inc. from 2004–2007; prior to that, he was senior VP Engineering and Chief Software Architect at BEA Systems responsible for ...

Hamid Pirahesh, Ph.D., is an IBM fellow, ACM Fellow and a senior manager responsible for the exploratory database department at IBM Research - Almaden in San Jose, California. Dr. Hamid Pirahesh is the senior manager at IBM Almaden Research Center in San Jose, California.

Jim Gray Summary Home Page

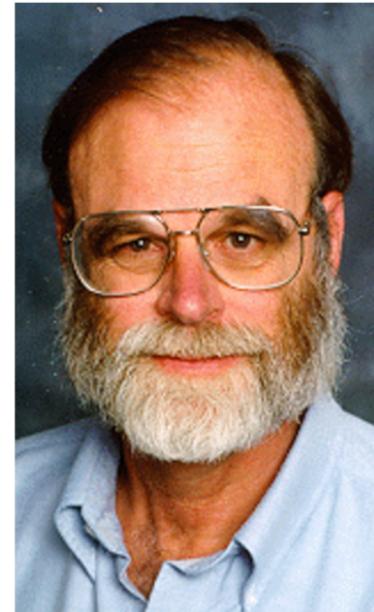
[Microsoft eScience Group](#)

As you may be aware, Jim Gray has [gone missing](#).

We (his colleagues in Microsoft Research) have heard from many of his collaborators about projects and collaborations that he had underway with them and who are unsure how to proceed. If you find yourself in this situation, please email grayproj@microsoft.com and we will follow up with you to find the best way forward.

Jim Gray is a researcher and manager of Microsoft Research's [eScience Group](#). His primary research interests are in databases and transaction processing systems -- with particular focus on using computers to make scientists more productive. He and his group are working in the areas of astronomy, geography, hydrology, oceanography, biology, and health care. He continues a long-standing interest on building supercomputers with commodity components, thereby reducing the cost of storage, processing, and networking by factors of 10x to 1000x over low-volume solutions. This includes work on building fast networks, on building huge web servers with *CyberBricks*, and building very inexpensive and very high-performance storage servers.

Jim also is working with the astronomy community to build the [world-wide telescope](#) and has been active in building online databases like <http://terraService.Net> and <http://skyserver.sdss.org>. When the entire world's astronomy data is on the Internet and is accessible as a single distributed database, the Internet will be the world's best telescope. This is part of the larger agenda of getting all information online and easily accessible (digital libraries, digital government, online science ...). More generally, he is working with the science community (Oceanography, Hydrology, environmental monitoring, ..) to build the world-wide digital library that integrates all the world's scientific literature and the data in one easily-accessible collection. He is active in the research community, is an ACM, NAE, NAS, and AAAS Fellow, and received the ACM Turing Award for his work on transaction processing. He also edits of a series of books on data management.



[https://en.wikipedia.org/wiki/Jim_Gray_\(computer_scientist\)](https://en.wikipedia.org/wiki/Jim_Gray_(computer_scientist))

James Nicholas "Jim" Gray (born January 12, 1944; presumed lost at sea January 28, 2007; declared deceased May 16, 2012^[4]) was an American computer scientist who received the Turing Award^[5] in 1998 "for seminal contributions to database and transaction processing research and technical leadership in system implementation."

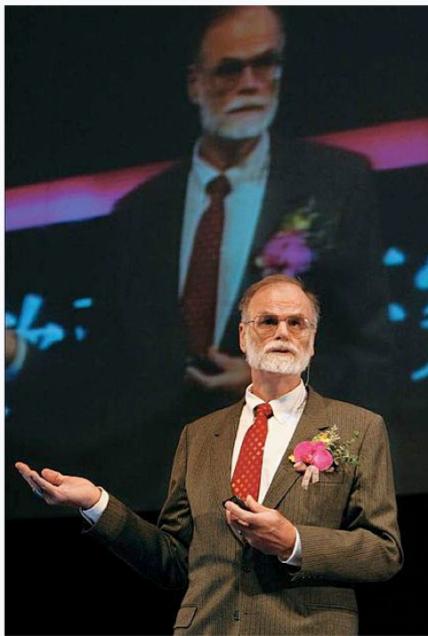
Contents [hide]

- 1 Early years
- 2 Research
- 3 Disappearance
- 4 Personal life
- 5 Jim Gray eScience Award
- 6 References
- 7 External links

Early years [edit]

Gray was born in San Francisco, California, the second child of a mother who was a teacher and a father in the U.S. Army; the family moved to Rome where Gray spent most of the first three years of his life, learning to speak Italian before English.^[2] The family then moved to Virginia, spending about four years there, until Gray's parents divorced, after which he returned to San Francisco with

Jim Gray



Gray in 2006

Born	James Nicholas Gray January 12, 1944 ^[1] San Francisco, California ^[2]
Disappeared	January 28, 2007 (aged 63) Waters near San Francisco
Status	Dead in absentia, May 16, 2012 (aged 68)
Nationality	American
Alma mater	University of California, Berkeley (Ph.D)
Occupation	Computer scientist
Employer	IBM Tandem Computers DEC Microsoft

On Sunday, January 28, 2007, during a short solo sailing trip to the Farallon Islands near San Francisco to scatter his mother's ashes, Gray and his 40-foot yacht, *Tenacious*, were reported missing by his wife, Donna Carnes. The Coast Guard searched for four days using a C-130 plane, helicopters, and patrol boats but found no sign of the vessel.^{[21][22][23][24]}

Gray's boat was equipped with an automatically deployable EPIRB (Emergency Position-Indicating Radio Beacon), which should have deployed and begun transmitting the instant his vessel sank. The area around the Farallon Islands where Gray was sailing is well north of the East-West ship channel used by freighters entering and leaving San Francisco Bay. The weather was clear that day and no ships reported striking his boat, nor were any distress radio transmissions reported.

On February 1, 2007, the DigitalGlobe satellite did a scan of the area, generating thousands of images.^[25] The images were posted to Amazon Mechanical Turk in order to distribute the work of searching through them, in hopes of spotting his boat.

In the immediate aftermath of the disappearance, many theories were put forward on how Gray disappeared.^[26]

After being missing for five years, Gray was legally assumed to have died at sea on January 28, 2012.^{[4][33]}



Jim Gray on the *Tenacious* in January 2006

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_i 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.

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

Cube Materialization:

Q: What do they hate the most?



Iceberg



Cube Materialization: Full Cube vs. Iceberg Cube

- Full cube vs. iceberg cube

```
compute cube sales iceberg as
select date, product, city, department, count(*)
from salesInfo
cube by date, product, city
having count(*) >= min support
```
- Compute *only* the **cells** whose **measure** satisfies the **iceberg condition**
- Only a small portion of cells may be “above the water” in a **sparse cube**
- Ex.: Show only those cells whose **count** is no less than 100



Why Iceberg Cube?

- Advantages of computing iceberg cubes
 - No need to save nor show those cells whose value is below the threshold (iceberg condition)
 - Efficient methods may even avoid computing the un-needed, intermediate cells
 - Avoid explosive growth
- Example: A cube with 100 dimensions
 - Suppose it contains only 2 base cells and the count of each cell is 1:
 - $\{(a_1, a_2, a_3, \dots, a_{100}) : 1, (a_1, a_2, b_3, \dots, b_{100}) : 1\}$
 - How many **aggregate cells** if “having count ≥ 1 ” (**non-empty**)?
 - What are the **iceberg cells** with condition “having count ≥ 2 ”?

Suppose it contains only 2 base cells:

$$\{(a_1, a_2, a_3, \dots, a_{100}), (a_1, a_2, b_3, \dots, b_{100})\}$$

How many non-empty aggregate cells?

For $\{(a_1, a_2, a_3, \dots, a_{100}), (a_1, a_2, b_3, \dots, b_{100})\}$, the total # of non-base cells should be $2 * (2^{100} - 1) - 4$.

This is calculated as follows:

- $(a_1, a_2, a_3, \dots, a_{100})$ will generate $2^{100} - 1$ non-base cells
- $(a_1, a_2, b_3, \dots, b_{100})$ will generate $2^{100} - 1$ non-base cells

Among these, 4 cells are overlapped and thus minus 4 so we get:

$$2 * 2^{100} - 2 - 4$$

These 4 cells are:

- $(a_1, a_2, *, \dots, *)$: 2
- $(a_1, *, *, \dots, *)$: 2
- $(*, a_2, *, \dots, *)$: 2
- $(*, *, *, \dots, *)$: 2

Is Iceberg Cube Good Enough?

Closed Cube & Cube Shell

- Let cube P have only 2 base cells: $\{(a_1, a_2, a_3 \dots, a_{100}):10, (a_1, a_2, b_3, \dots, b_{100}):10\}$
 - How many cells will the iceberg cube contain if “having count(*) ≥ 10 ”?
 - Answer: $2^{101} - 4$ (base+aggregate; still too big!)
- **Close cube:**
 - A cell c is **closed** if there **exists no cell d**, such that d is a **descendant** of c, and d has the **same measure** value as c
 - Ex. The same cube P has only 3 closed cells:
 - $\{(a_1, a_2, *, \dots, *): 20, (a_1, a_2, a_3 \dots, a_{100}): 10, (a_1, a_2, b_3, \dots, b_{100}): 10\}$
 - A **closed cube** is a cube consisting of only closed cells

References

- 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
- K. Beyer and R. Ramakrishnan. Bottom-Up Computation of Sparse and Iceberg CUBEs.. SIGMOD'99
- J. Han, J. Pei, G. Dong, K. Wang. Efficient Computation of Iceberg Cubes With Complex Measures. SIGMOD'01
- L. V. S. Lakshmanan, J. Pei, and J. Han, Quotient Cube: How to Summarize the Semantics of a Data Cube, VLDB'02
- X. Li, J. Han, and H. Gonzalez, High-Dimensional OLAP: A Minimal Cubing Approach, VLDB'04
- X. Li, J. Han, Z. Yin, J.-G. Lee, Y. Sun, "Sampling Cube: A Framework for Statistical OLAP over Sampling Data", SIGMOD'08
- K. Ross and D. Srivastava. Fast computation of sparse datacubes. VLDB'97
- D. Xin, J. Han, X. Li, B. W. Wah, Star-Cubing: Computing Iceberg Cubes by Top-Down and Bottom-Up Integration, VLDB'03
- Y. Zhao, P. M. Deshpande, and J. F. Naughton. An array-based algorithm for simultaneous multidimensional aggregates. SIGMOD'97
- D. Burdick, P. Deshpande, T. S. Jayram, R. Ramakrishnan, and S. Vaithyanathan. OLAP over uncertain and imprecise data. VLDB'05

References (cont.)

- R. Agrawal, A. Gupta, and S. Sarawagi. Modeling multidimensional databases. ICDE'97
- B.-C. Chen, L. Chen, Y. Lin, and R. Ramakrishnan. Prediction cubes. VLDB'05
- B.-C. Chen, R. Ramakrishnan, J.W. Shavlik, and P. Tamma. Bellwether analysis: Predicting global aggregates from local regions. VLDB'06
- Y. Chen, G. Dong, J. Han, B. W. Wah, and J. Wang, Multi-Dimensional Regression Analysis of Time-Series Data Streams, VLDB'02
- R. Fagin, R. V. Guha, R. Kumar, J. Novak, D. Sivakumar, and A. Tomkins. Multi-structural databases. PODS'05
- J. Han. Towards on-line analytical mining in large databases. SIGMOD Record, 27:97–107, 1998
- T. Imielinski, L. Khachiyan, and A. Abdulghani. Cubegrades: Generalizing association rules. Data Mining & Knowledge Discovery, 6:219–258, 2002.
- R. Ramakrishnan and B.-C. Chen. Exploratory mining in cube space. Data Mining and Knowledge Discovery, 15:29–54, 2007.
- K. A. Ross, D. Srivastava, and D. Chatziantoniou. Complex aggregation at multiple granularities. EDBT'98
- S. Sarawagi, R. Agrawal, and N. Megiddo. Discovery-driven exploration of OLAP data cubes. EDBT'98
- G. Sathe and S. Sarawagi. Intelligent Rollups in Multidimensional OLAP Data. VLDB'01