## Where is the Science in the Software Sciences?

#### MAPi 2017

#### FEUP, 15th Sep 2017



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで

L'enfant terrible

yped LA

References



## Informatics $\simeq$ Computer "Science"?





Vinton Cerf (1943-) ACM President

#### DOI:10.1145/2347736.2347737

# Where is the Science in Computer Science?

"... we have a responsibility to pursue the science in computer science. We must develop better tools and much deeper understanding of the systems we invent and a far greater ability to make predictions about the behavior of these complex, connected, and interacting systems.".

> (Vinton G. Cerf, Letter from the ACM President, CACM 55(10), Oct. 2012)







▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ





#### ( "L'enfant terrible" is born )



1st NATO Conference on Software Engineering, Darmstadt, October 1968







**Hardware** and other "traditional" industrial **products** fabricated according to the laws of **physics**.

Software is not governed by the laws of physics:

- it does not weight / does not smell
- it does not warm up / cool down
- it is chemically neutral ...

Anthony Oettinger (ACM President, 1967):

"(...) the scientific, rigorous component of computing, is more like **mathematics** than it is like **physics**".





Can one *pretend* that **software production** is not affected by its **special nature** and move on?

Many people have tried to do so for 50 years...

Still Oettinger (already in 1967):

"It is a matter of **complexity**. Once you start putting thousands of these instructions together you create a **monster** which is **unintelligible** to anyone save its creator and, most of the time, unfortunately even to the creator."



Famous textbook by Niklaus Wirth (1934-) that inspired **structured programming** and made history in its time.

Published in 1976.

PASCAL — the language "par exellence" (mostly in Europe).

No **Java** yet. Object orientation stimulated by MODULA, a PASCAL derivative.



・ロト ・ 雪 ト ・ ヨ ト ・ ヨ ト



The "textbook of all textbooks".

Emphasis on elegance and abstraction.

Also published in 1976.

Programming as a scientific activity.

E.W. Dijkstra (1930-2002), Turing Award, one of the most eminent computer scientists ever.



L'enfant terrible

Golden age

ata mining

Fyped LA

ata cube

Summary

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

References



## Shortly afterwards



Famous "dragon book" by A.V. Aho (1941-) e J.D. Ullman (1942-).

Published in 1977.

Inspiration for generations of compiler writers.

Compiler writing made systematic.



Scientific basis: *The Theory of Parsing, Translating & Compiling* (1972) by the same authors.

L'enfant terrible

Golden age

ita mining

yped LA

ata cube

Reference



## Still people complain (1976)





E.W. Dijkstra (1920-2002)

(...) the term "Database Technology", although sometimes used, is immature, for there is hardly any underlying "science" that could justify the use of the term "technology".

(...) They seem to form an inbred crowd with very little knowledge of computing science in general (...) Often they seemed to be mentally trapped by the intricacies of early, rather **ad hoc** solutions (...)

(E.W. Dijkstra, EWD577, 1976)



Relational database science — books by Ullman (1982) and Maier (1983).

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ



Have scientific standards kept up with so higher demands in technology?

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 … のへで

L'enfant terrible

Golden age

Data mining

Typed LA

ata cube

Summary

References



## A much older concern



"Only by taking infinitesimally small units for observation (the **differential** of history, that is, the individual tendencies of men) and attaining to the art of **integrating** them (that is, finding the sum of these infinitesimals) can we hope to arrive at the laws of history."

> Leo Tolstoy, "War and Peace" - Book XI, Chap.II (1869)



L. Tolstoy (1828-1910)

The oldest definition I know of data-mining.

L'enfant terrible

MAP i DOCTORAL PROSPAM

Golden age

Data mining

Typed LA

Data cube

Summary

References



## The age of data

Note the **maths** flavour of Tolstoy's text

How maths-strong is data-mining?

Have we attained the "art of integration" foreseen by Tolstoy?

Again the **glorious 1970s** — since the early days of psychometrics in the **social sciences** — **linear algebra** (LA) central to data analysis (e.g. Tucker tensor decompositions etc).

All done? Anything else to do?



L.R. Tucker (1910-2004)



Data mining seems to be ignoring the role of **types** and **type parametricity** in software — one of the most significant advances in CS.

Nice theory called **parametric polymorphism** (John Reynolds, CMU).

So nice that you can derive **properties** of your (typed) operations **before** writing them (!)

What can **parametricity** possibly mean in the **data-analysis** setting?



J.C. Reynolds (1935-2013)



Back to square one — some raw data:

t

	#	Model	Year	Color	Sale
	1	Chevy	1990	Red	5
	2	Chevy	1990	Blue	87
=	3	Ford	1990	Green	64
	4	Ford	1990	Blue	99
	5	Ford	1991	Red	8
	6	Ford	1991	Blue	7

**Columns** — attributes — the **observables Rows** — records (*n*-many) — the **infinitesimals** 

**Column-orientation** — each column (attribute) A represented by a function  $t_A : n \to A$  such that  $a = t_A(i)$  means "a is the value of attribute A in record nr i".



Can records be rebuilt from such attribute projection functions?

Yes — by **tupling** them.

**Tupling**: Given functions  $f : A \rightarrow B$  and  $g : A \rightarrow C$ , their tupling is the function  $f \lor g$  such that  $(f \lor g) a = (f a, g a)$ 

For instance,

 $(t_{Color} \circ t_{Model}) = (Blue, Chevy),$  $(t_{Year} \circ (t_{Color} \circ t_{Model})) = (1990, (Green, Ford))$ 

and so on.



Represent functions by Boolean matrices.

Given (finite) types A and B, any function  $f: A \rightarrow B$ 

can be represented by a matrix  $\llbracket f \rrbracket$  with *A*-many columns and *B*-many rows such that, for any  $b \in B$  and  $a \in A$ , matrix cell

 $b \llbracket f \rrbracket a = \begin{cases} 1 \Leftarrow b = f \\ 0 \text{ otherwise} \end{cases}$ 

**NB**: Following the **infix** notation usually adopted for relations (which are Boolean matrices) — for instance  $y \leq x$  — we write  $y \ M x$  to denote the contents of the cell in matrix M addressed by row y and column x.



#### One projection function (matrix) per dimension attribute:

t <sub>Model</sub>	1	2	3	4	5	6					
Chevy	1	1	0	0	0	0	-				
Ford	0	0	1	1	1	1	#	Model	Year	Color	Sale
tu I	1	2	2	Λ	Б	6	1	Chevy	1990	Red	5
L Year	1	2	3	4	5	0	2	Chevy	1990	Blue	87
1990	1	1	1	1	0	0	3	Ford	1990	Green	64
1991	0	0	0	0	1	1	4	Ford	1990	Blue	99
tcolor	1	2	3	4	5	6	5	Ford	1991	Red	8
Blue	0	1	0	1	0	1	6	Ford	1991	Blue	7
Green	0	0	1	0	0	0					
Red	1	0	0	0	1	0					

**NB**: we tend to abbreviate [f] by f where the context is clear.



Note how the inverse of a function is also represented by a Boolean matrix, e.g.

$t^{\circ}_{Model}$	Chevy	Ford								
1	1	0								
2	1	0		t <sub>Model</sub>	1	2	3	4	5	6
3	0	1	versus	Chevy	1	1	0	0	0	0
4	0	1		Ford	0	0	1	1	1	1
5	0	1								
6	0	1								

- no longer a function, but still a relation, a Boolean matrix :-)

Clearly,

 $j t_{Model}^{\circ} a = a t_{Model} j$ 

Given a matrix M,  $M^{\circ}$  is known as the **transposition** of M.



We **type** matrices in the same way as functions:  $M : A \rightarrow B$  means a matrix M with A-many columns and B-many rows.

Matrices are **arrows**:  $A \xrightarrow{M} B$  denotes a matrix from A (source) to B (target), where A, B are (finite) types.

Writing  $B \stackrel{M}{\longleftarrow} A$  means the same as  $A \stackrel{M}{\longrightarrow} B$ .

**Composition** — *aka* matrix multiplication:



 $b(M \cdot N)c = \langle \sum a :: (b M a) \times (a N c) \rangle$ 

L'enfant terrible

Data cube

References



## The typed LA approach



(1)

**Function composition** implemented by matrix multiplication,  $\llbracket f \cdot g \rrbracket = \llbracket f \rrbracket \cdot \llbracket g \rrbracket$ 

 $\ensuremath{\textbf{ldentity}}$  — the identity matrix  $\ensuremath{\textit{id}}$  corresponds to the identity function and is such that

 $M \cdot id = M = id \cdot M$ 

Function tupling corresponds to the so-called Khatri-Rao product  $M \lor N$  defined index-wise by

 $(b,c) (M \lor N) a = (b M a) \times (c N a)$  (2)

Khatri-Rao is a "column-wise" version of the well-known **Kronecker product**  $M \otimes N$ :

$$(y,x) (M \otimes N) (b,a) = (y M b) \times (x N a)$$
(3)



The raw data given above is represented in typed LA by the expression

$$v = (t_{Year} \circ (t_{Color} \circ t_{Model})) \cdot (t^{Sale})^{\circ}$$

of type

$$v: 1 \rightarrow (Year \times (Color \times Model))$$

depicted aside.

v is a **multi-dimensional** column vector — a **tensor**. Datatype  $1 = \{ALL\}$  is the so-called **singleton** type.

Year x (	Color x	Model)	ALL
	Blue	Chevy	87
		Ford	99
1990	Casar	Chevy	0
	Green	Ford	64
	Pod	Chevy	5
	Keu	Ford	0
	Plue	Chevy	0
	Diue	Ford	7
1001	Croon	Chevy	0
1991	Green	Ford	0
	Pod	Chevy	0
	Rea	Ford	8

HASLab



## Dimensions and measures

**Sale** is a special kind of data — a **measure**. Measures are encoded as **row** vectors, e.g.

recall

#	Model	Year	Color	Sale
1	Chevy	1990	Red	5
2	Chevy	1990	Blue	87
3	Ford	<i>1990</i>	Green	64
4	Ford	1990	Blue	<i>99</i>
5	Ford	1991	Red	8
6	Ford	1991	Blue	7



Measures provide for integration in Tolstoy's sense — aka consolidation

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙



There is a unique function in type  $A \rightarrow 1$ , usually named  $A \xrightarrow{!} 1$ . This corresponds to a row vector wholly filled with 1s. Example:  $2 \xrightarrow{!} 1 = \begin{bmatrix} 1 & 1 \end{bmatrix}$ 

Given  $M: B \to A$ , the expression  $! \cdot M$  (where  $A \xrightarrow{!} 1$ ) is the row vector (of type  $B \to 1$ ) that contains all column **totals** of M,

 $\begin{bmatrix} 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 50 & 40 & 85 & 115 \\ 50 & 10 & 85 & 75 \end{bmatrix} = \begin{bmatrix} 100 & 50 & 170 & 190 \end{bmatrix}$ 

Given type A, define its **totalizer** matrix  $A \xrightarrow{\tau_A} A + 1$  by

$$\tau_{A} : A \to A + 1$$
  
$$\tau_{A} = \left[\frac{id}{!}\right]$$
(5)

Thus  $\tau_A \cdot M$  yields a copy of M on top of the corresponding totals.



Data cubes can be obtained from products of totalizers.

Recall the Kronecker (tensor) product  $M \otimes N$  of two matrices  $A \xrightarrow{M} B$  and  $C \xrightarrow{N} D$ , which is of type  $A \times C \xrightarrow{M \otimes N} B \times D$ .

The matrix

$$A \times B \xrightarrow{\tau_A \otimes \tau_B} (A+1) \times (B+1)$$

provides for totalization on the two dimensions A and B.

Indeed, type  $(A + 1) \times (B + 1)$  is isomorphic to  $A \times B + A + B + 1$ , whose four parcels represent the four elements of the "dimension powerset of  $\{A, B\}$ ".



Recalling

$$v = (t_{Year} \circ (t_{Color} \circ t_{Model})) \cdot (t^{Sale})^{\circ}$$

build

 $c = ( au_{Year} \otimes ( au_{Color} \otimes au_{Model})) \cdot v$ 

This is the multidimensional vector (tensor) representing the **data cube** for

- dimensions Year, Color, Model
- measure Sale

depicted aside.

(Year+1	L) x ((Co	olor+1) x (Model+1))	ALL
		Chevy	87
	Blue	Ford	99
		ALL	186
		Chevy	0
	Green	Ford	64
		ALL	64
1990		Chevy	5
	Red	Ford	0
		ALL	5
		Chevy	92
	ALL	Ford	163
		ALL	255
		Chevy	0
	Blue	Ford	7
		ALL	7
		Chevy	0
	Green	Ford	0
		ALL	0
1991		Chevy	0
	Red	Ford	8
		ALL	8
		Chevy	0
	ALL	Ford	15
		ALL	15
		Chevy	87
	Blue	Ford	106
		ALL	193
		Chevy	0
	Green	Ford	64
ALL		ALL	64
rach		Chevy	5
	Red	Ford	8
		ALL	13
		Chevy	92
	ALL	Ford	178
		ALL	270

・ロト・西ト・山田・山田・山市・山口・



In our approach a **cube** is not necessarily one such column vector.

The key to generic data cubes is (generic) **vectorization**, a kind of "**matrix currying**": given  $A \times B \xrightarrow{M} C$  with  $A \times B$ -many columns and *C*-many rows, reshape *M* into its **vectorized** version  $B \xrightarrow{\text{vec}_A M} A \times C$  with *B*-many columns and  $A \times C$ -many rows.

Matrices M and  $\operatorname{vec}_A M$  are **isomorphic** — they contain the same information in different formats, cf.

 $c M (a, b) = (a, c) (\operatorname{vec}_A M) b$ (6)

for every *a*, *b*, *c*.



Vectorization thus has an inverse operation — unvectorization:



That is, M can be retrieved back from  $\operatorname{vec}_A M$  by unvectorizing it:

$$N = \operatorname{vec}_A M \quad \Leftrightarrow \quad \operatorname{unvec}_A N = M \tag{7}$$

Vectorization has a rich algebra, e.g. a fusion-law

 $(\mathbf{vec}\,M)\cdot N = \mathbf{vec}\,(M\cdot(id\otimes N)) \tag{8}$ 

and an **absorption**-law:

$$\operatorname{vec}(M \cdot N) = (id \otimes M) \cdot \operatorname{vec} N \tag{9}$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・



There is room for further unvectorizing the outcome, this time across *Color* — next slide:

・ロト ・ 雪 ト ・ ヨ ト

э



#### Further unvectorization:



	BI	Blue		een	Red		
	1990	1990 1991		1991	1990	1991	
Chevy	87	0	0	0	5	0	
Ford	99	7	64	0	0	8	

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

and so on.



It turns out **that** cubes can be calculated for any such two-dimensional versions of our original data tensor, for instance,

> cube N:  $Model + 1 \leftarrow (Color + 1) \times (Year + 1)$ cube  $N = \tau_{Model} \cdot N \cdot (\tau_{Color} \otimes \tau_{Year})^{\circ}$

where N stands for the second matrix of the previous slide, yielding

	Blue			Green			Red			ALL		
	1990	1991	ALL	1990	1991	ALL	1990	1991	ALL	1990	1991	ALL
Chevy	87	0	87	0	0	0	5	0	5	92	0	92
Ford	99	7	106	64	0	64	0	8	8	163	15	178
ALL	186	7	193	64	0	64	5	8	13	255	15	270

See how the 36 entries of the original cube have been rearranged in a 3\*12 rectangular layout, as dictated by the **dimension** cardinalities.



Definition (Cube) Let *M* be a matrix of type

$$\prod_{j=1}^{n} B_j \stackrel{M}{\longleftarrow} \prod_{i=1}^{m} A_i \tag{10}$$

We define matrix **cube** M, the cube of M, as follows

cube 
$$M = (\bigotimes_{j=1}^{n} \tau_{B_j}) \cdot M \cdot (\bigotimes_{i=1}^{m} \tau_{A_i})^{\circ}$$
 (11)

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

where  $\bigotimes$  is finite Kronecker product.

So **cube** *M* has type  $\prod_{j=1}^{n} (B_j + 1) \leftarrow \prod_{i=1}^{m} (A_i + 1)$ .



Linearity:

$$cube (M + N) = cube M + cube N$$
(12)

**Proof:** Immediate by bilinearity of matrix composition:

$$M \cdot (N+P) = M \cdot N + M \cdot P$$
(13)  
(N+P) \cdot M = N \cdot M + P \cdot M (14)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

This can be taken advantage of not only in **incremental** data cube construction but also in **parallelizing** data cube generation.



Updatability: by Khatri-Rao product linearity,

 $(M+N) \circ P = M \circ P + N \circ P$  $P \circ (M+N) = P \circ M + P \circ N$ 

the **cube** operator commutes with the usual CRUDE operations, namely record **updating**. For instance, suppose record

	#	Model	Year	Color	Sale		t <sub>Model</sub>	1	2	3	4	5	6
		E.u.l	1001	D	0	cf	Chevy	1	1	0	0	0	0
	5	Fora	1991	Rea	ð		Ford	0	0	1	1	1	1
is up	date	ed to											
	#	Model	Voar	Color	Sala		t' <sub>Model</sub>	1	2	3	4	5	6
	#	widder	1001		Oulc	cf	Chevy	1	1	0	0	1	0
	5	Chevy	1991	Ked	8		Ford	0	0	1	1	0	1





**Commutativity**: cube commutes with **vectorization**:

Let  $X < \frac{M}{M} Y \times C$  and  $Y \times X < \frac{\text{vec } M}{C}$  be its Y-vectorization. Then

 $\operatorname{vec}\left(\operatorname{cube}\,M\right) = \operatorname{cube}\,\left(\operatorname{vec}\,M\right) \tag{15}$ 

holds. 🗆

Type diagrams:



・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Proof in (Oliveira and Macedo, 2017).



The following theorem shows that changing the dimensions of a data cube does not change its totals.

### Theorem (Free theorem)

Let  $B \stackrel{M}{\leftarrow} A$  be cubed into  $B + 1 \stackrel{\text{cube } M}{\leftarrow} A + 1$ , and  $r : C \rightarrow A$ and  $s : D \rightarrow B$  be arbitrary functions. Then

**cube**  $(s^{\circ} \cdot M \cdot r) = (s^{\circ} \oplus id) \cdot ($ **cube**  $M) \cdot (r \oplus id)$  (16) holds, where  $M \oplus N = \begin{bmatrix} M & 0 \\ 0 & N \end{bmatrix}$  is matrix **direct sum**.

The proof given in the same paper resorts to the **free theorem** of **polymorphic** operators due to J. Reynolds and popularized by Wadler (1989) under the heading **Theorems for free!**.



**Slicing** is a specialized filter for a particular value in a dimension.

Suppose that from our starting cube

 $c: 1 \rightarrow (Year + 1) \times ((Color + 1) \times (Model + 1))$ 

one is only interested in the data concerning year 1991.

It suffices to regard data values as (categorial) **points**: given  $p \in A$ , constant function  $\underline{p}: 1 \to A$  is said to be a *point* of A, for instance

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

$$\underline{1991}: 1 \rightarrow Year + 1$$
$$\underline{1991} = \begin{bmatrix} 0\\1\\0 \end{bmatrix}$$





Gray et al. (1997) say that going up the levels [of aggregated data] is called rolling-up.

In this sense, a **roll-up** operation over dimensions A, B and C could be the following form of (increasing) summarization:

```
A \times (B \times C)A \times BA1
```

How does this work over a data cube? We take the simpler case of two dimensions A, B as example.



The dimension powerset for *A*, *B* is captured by the corresponding matrix **injections** onto the cube target type  $(A + 1) \times (B + 1)$ :



where

 $\begin{aligned} \theta &= i_1 \otimes i_1 \\ \alpha &= i_1 \lor i_2 \cdot ! \\ \beta &= i_1 \cdot ! \lor i_2 \\ \omega &= i_2 \lor i_2 \end{aligned}$ 

**NB**: the injections  $i_1$  and  $i_2$  are such that  $[i_1|i_2] = id$ , where [M|N] denotes the horizonal gluing of two matrices.



One can build compound injections, for instance

 $\rho: (A+1) \times (B+1) \leftarrow A \times B + (A+1)$  $\rho = [\theta | [\alpha | \omega]]$ 

Then, for  $M : C \rightarrow A \times B$ :

$$\rho^{\circ} \cdot (\mathbf{cube} \ M) = \left[\frac{M}{\left[\frac{fst \cdot M}{1 \cdot M}\right]}\right] \cdot \tau_{C}^{\circ}$$

extracts from **cube** *M* the corresponding **roll-up**.

The next slides give a concrete example.



MAP i DOCTORAL PROBRAM

iolden age

ta mining

Typed LA

Data cube

References





#### Let M be the (generalized) data cube

	1990	1991	ALL
Chevy	87	0	87
Blue Ford	99	7	106
ALL	186	7	193
Chevy	0	0	0
Green Ford	64	0	64
ALL	64	0	64
Chevy	5	0	5
Red Ford	0	8	8
ALL	5	8	13
Chevy	92	0	92
ALL <i>Ford</i>	163	15	178
ALL	255	15	270

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ



Building the injection matrix  $\rho = [\theta | [\alpha | \omega]]$  for types  $Color \times Model + Color + 1 \rightarrow (Color + 1) \times (Model + 1)$  we get the following matrix (already transposed):

			Blue		Green			Red			ALL		
		Chevy	Ford	ALL									
Plus	Chevy	1	0	0	0	0	0	0	0	0	0	0	0
Diue	Ford	0	1	0	0	0	0	0	0	0	0	0	0
Green	Chevy	0	0	0	1	0	0	0	0	0	0	0	0
	Ford	0	0	0	0	1	0	0	0	0	0	0	0
Ded	Chevy	0	0	0	0	0	0	1	0	0	0	0	0
Rea	Ford	0	0	0	0	0	0	0	1	0	0	0	0
	Blue	0	0	1	0	0	0	0	0	0	0	0	0
	Green	0	0	0	0	0	1	0	0	0	0	0	0
	Red	0	0	0	0	0	0	0	0	1	0	0	0
	ALL	. 0	0	0	0	0	0	0	0	0	0	0	1

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



Note how a roll-up is a particular "subset" of a cube.

Matrix  $\rho^{\circ}$  performs the (quantitative) selection of such a subset.



There is a **lot** of SCIENCE in CS that CS people simply **ignore**.



Abadir and Magnus (2005) stress on the need for a **standardized** notation for **linear algebra** in the field of econometrics and **statistics**.

This talk suggests such a notation should be **properly typed**.

Since (Macedo and Oliveira, 2013) the author has invested in **typing** linear algebra in a way that makes it closer to modern **typed** languages.

This extends previous efforts on applying LA to **OLAP** (Macedo and Oliveira, 2015)



MAPi is the right **environment** for this kind of cross-breeding:

- 3 universities
- many labs
- different cultures, nice cooperation...



#### (for those who care mostly about efficiency)

Theory pays off!

Plot taken from a recent MSc report on TPC-H benchmarking LA approach:





## References

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

L'enfant terrible

ata mining

Typed LA

ata cube

References

K.M. Abadir and J.R. Magnus. *Matrix algebra. Econometric exercises 1.* C.U.P., 2005.

- Jim Gray, Surajit Chaudhuri, Adam Bosworth, Andrew Layman, Don Reichart, Murali Venkatrao, Frank Pellow, and Hamid Pirahesh. Data cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals. J. Data Mining and Knowledge Discovery, 1(1):29–53, 1997. URL citeseer.nj.nec.com/article/gray95data.html.
- H.D. Macedo and J.N. Oliveira. Typing linear algebra: A biproduct-oriented approach. SCP, 78(11):2160–2191, 2013.
- H.D. Macedo and J.N. Oliveira. A linear algebra approach to OLAP. *FAoC*, 27(2):283–307, 2015.
- J. N. Oliveira and H. D. Macedo. The data cube as a typed linear algebra operator. In *Proc. of the 16th Int. Symposium on Database Programming Languages*, DBPL '17, pages 6:1–6:11, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-5354-0. doi: 10.1145/3122831.3122834.
- P.L. Wadler. Theorems for free! In 4th International Symposium

L'enfant terrible

Fyped LA

ta cube

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

on Functional Programming Languages and Computer Architecture, pages 347–359, London, Sep. 1989. ACM.