

Do the middle letters of “OLAP” stand for Linear Algebra (“LA”)?

J.N. Oliveira
(joint work with H. Macedo)

HASLab/Universidade do Minho
Braga, Portugal

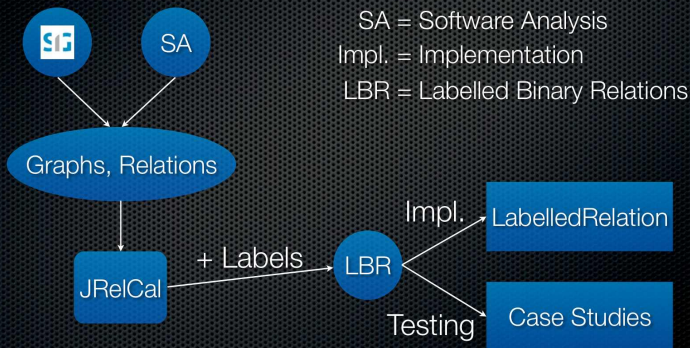
SIG
Amsterdam, NL
26th May 2011

Context

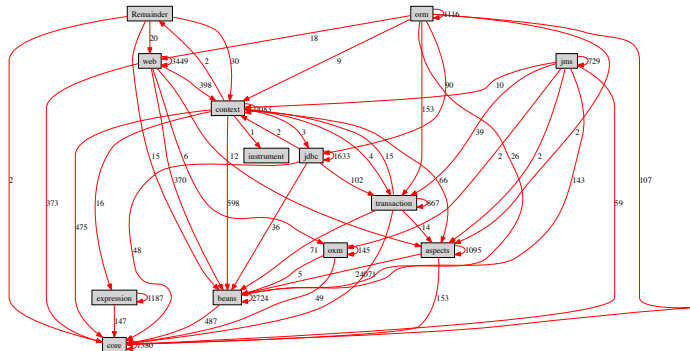
- **HASLab** is a research group at Minho University in Braga, Portugal
- The group has been concerned with developing techniques for *high assurance software*
- **HASLab** - **SIG** collaboration on a regular basis since 2007
- **SIG** contributes with knowledge transfer in the area of **software quality**
- **HASLab** does so in the area of formal analysis, modelling and verification.

A previous collaborative project

Introduction



On-going collaborative project



Mining call-graphs for software **architecture** quality profiling.

Motivation

Clearly:

- Need to **quantify** over relationships
- Raw information too fine-grained
- Too much information involved in data mining
- Need to make sense of **huge** data banks

Need for data summarizing techniques.

Motivation

In fact:

- Data summaries unveil **trends** hidden in raw data
- One sees the “**big picture**”

State-of-the-art:

- “OLAP” stands for *On Line Analytical Processing*
- Proprietary solutions (IBM, Oracle, MS)
- Calls for parallelism
- Expensive.

Can **parallel OLAP** be made more widely accessible?

OLAP's "Hello World"

As generated in MS Excel (choose Data > PivotTableReport):

Raw data

Model	Year	Color	Sales
Chevy	1990	Red	5
Chevy	1990	Blue	87
Ford	1990	Green	64
Ford	1990	Blue	99
Ford	1991	Red	8
Ford	1991	Blue	7

"How many vehicles were sold per color and model?"

Sum of Sales	Model		
Color	Chevy	Ford	Grand Total
Blue	87	106	193
Green		64	64
Red	5	8	13
Grand Total	92	178	270

CTAB

Pivot table

Three dimensions — *Model*, *Year*, *Color* — and one measure — *Sales*. Summarizing over *Year*.

OLAP Cubes

For huge raw data sets

- such **cross tabulation** summaries (vulg. “**pivot tables**”) take too long to generate
- The same tabulation likely to be required by different people in the organization
- Solution: generate all possible summaries overnight so that businessmen can have all pivot tables afresh the day after.
- Build a “**cube**” with all such multi-dimension projections.

```
Chevy 1990 Blue 87
Chevy 1990 Red 5
Ford 1990 Blue 99
Ford 1990 Green 64
Ford 1991 Blue 7
Ford 1991 Red 8
```

```
-----
Chevy 1990 ALL 92
Ford 1990 ALL 163
Ford 1991 ALL 15
Chevy ALL Blue 87
Chevy ALL Red 5
Ford ALL Blue 106
Ford ALL Green 64
Ford ALL Red 8
ALL 1990 Blue 186
ALL 1990 Green 64
ALL 1990 Red 5
ALL 1991 Blue 7
ALL 1991 Red 8
```

```
-----
Chevy ALL ALL 92
Ford ALL ALL 178
ALL 1990 ALL 255
ALL 1991 ALL 15
ALL ALL Blue 193
ALL ALL Green 64
ALL ALL Red 13
```

```
-----
ALL ALL ALL 270
```


OLAP — which theory behind?

OLAP:

- Cross tabulations are matrices (2-dim)
- What about OLAP cubes?
- Why SQL, GROUPBY, and so on?

Parallel solutions:

- MS Excel spreadsheet users may legitimately ask:

Is the generation of pivot tables in Excel actually taking advantage of the underlying multi-core hardware?

How parallel is such a construction?

Inspiration: relational algebra

Projecting over some attributes is a standard operation in relation algebra:

Table

A	B	C	D
a1	b1	c1	d1
a2	b1	c2	d2
a3	b1	c1	d3

"Only interested in B and C"

B	C
b1	c1
b1	c2

$\pi_{B,C}$

Projected table

Relational projection

Given T , a set of tuples:

$$\pi_{B,C}T = \{(t[B], t[C]) \mid t \in T\}$$

Pointwise relational:

$$b(\pi_{B,C}T)c \Leftrightarrow \langle \exists t : t \in T : b = t[B] \wedge c = t[C] \rangle$$

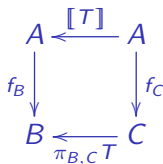
Pointfree relational:

$$\pi_{B,C}T = f_B \cdot [T] \cdot f_C^\circ$$

where $f_X t = t[X]$ and $[T] = \{(t, t) \mid t \in T\}$ — a coreflexive binary relation.

Types (“Relations as arrows”)

For $T \in \text{set } A$, a set of tuples:



In general: R **any** binary relation and f, g **arbitrary** functions in

$$\pi_{g,f} R = g \cdot R \cdot f^\circ \quad (1)$$

A commutative diagram with four nodes: A (top-left), B (top-right), C (bottom-left), and D (bottom-right). Arrows are: $A \xleftarrow{R} B$ (top), $A \xrightarrow{g} C$ (left), $B \xrightarrow{f} D$ (right), and $D \xleftarrow{\pi_{g,f} R} C$ (bottom).

Question

How to project data without losing **quantitative** information in measure columns such as eg. *Sales* in

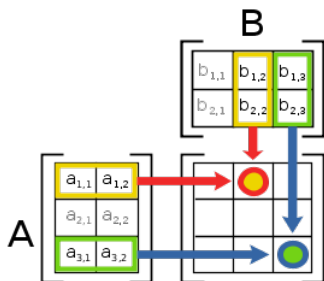
Model	Year	Color	Sales
Chevy	1990	Red	5
Chevy	1990	Blue	87
Ford	1990	Green	64
Ford	1990	Blue	99
Ford	1991	Red	8
Ford	1991	Blue	7

Clearly:

- Relational projection needs to take quantities into account
- Weighted graphs?
- Call them a proper name: we need **matrices**!

MMM — matrix matrix multiplication

From the Wikipedia:

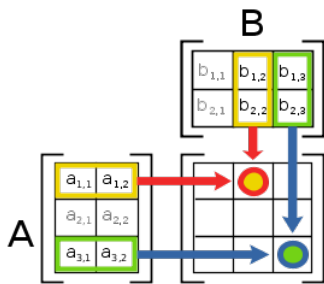


Index-wise definition

$$C_{ij} = \sum_{k=1}^{2,3} A_{ik} \times B_{kj}$$

MMM — matrix matrix multiplication

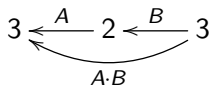
From the Wikipedia:



Index-wise definition

$$C_{ij} = \sum_{k=1}^{2,3} A_{ik} \times B_{kj}$$

Hiding indices i, j, k :



Index-free

$$C = A \cdot B$$

“Matrices as Arrows”

Given

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}_{m \times n}$$

$$m \xleftarrow{A} n$$

$$B = \begin{bmatrix} b_{11} & \dots & b_{1k} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{nk} \end{bmatrix}_{n \times k}$$

$$n \xleftarrow{B} k$$

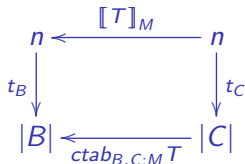
define

$$m \xleftarrow{A} n \xleftarrow{B} k$$

$$\xleftarrow{A \cdot B}$$

Projection (relations) \rightarrow tabulation (matrices)

Types:



where

- n — number of rows in raw data collection T
- $|B|$ — set of values in column B of T
- $|C|$ — set of values in column C of T
- $[T]_M$ — diagonal matrix storing all measures in column M of T
- t_X — “Membership matrix” of (non-metric) column X .

Details

Each column A in T “is” a function which tells, for each row, which value of $|A|$ can be found in such column, which *matricizes* into:

$$t_A : |A| \leftarrow n$$

$$a \ t_A \ r = \begin{cases} 1 & \text{if } T(r, A) = a \\ 0 & \text{otherwise} \end{cases}$$

Diagonal construction for measure column M :

$$[T]_M : n \leftarrow n$$

$$j[T]_M i = \begin{cases} T(j, M) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

“Hello World” illustration

Recall raw data example:

Model	Year	Color	Sales
Chevy	1990	Red	5
Chevy	1990	Blue	87
Ford	1990	Green	64
Ford	1990	Blue	99
Ford	1991	Red	8
Ford	1991	Blue	7

“Hello World” illustration

$$|Model| \xleftarrow{t_{Model}} 6$$

$$t_{Model} = \begin{array}{c} \\ \end{array} \begin{array}{c} 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \\ \hline Chevy \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \\ Ford \quad 0 \quad 0 \quad 1 \quad 1 \quad 1 \quad 1 \end{array}$$

$$|Color| \xleftarrow{t_{Color}} 6$$

$$t_{Color} = \begin{array}{c} \\ \\ \end{array} \begin{array}{c} 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \\ \hline Blue \quad 0 \quad 1 \quad 0 \quad 1 \quad 0 \quad 1 \\ Green \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0 \\ Red \quad 1 \quad 0 \quad 0 \quad 0 \quad 1 \quad 0 \end{array}$$

Counting

Typewise, the composition of matrices $|Color| \xleftarrow{t_{Color}} 6$ and $6 \xleftarrow{t_{Model}^\circ} |Model|$ makes sense and yields

$$t_{Color} \cdot t_{Model}^\circ = \begin{array}{c} \phantom{t_{Color} \cdot t_{Model}^\circ} \\ \phantom{t_{Color} \cdot t_{Model}^\circ} \\ \phantom{t_{Color} \cdot t_{Model}^\circ} \end{array} \begin{array}{cc} & \begin{array}{cc} Chevy & Ford \end{array} \\ \hline \begin{array}{c} Blue \\ Green \\ Red \end{array} & \begin{array}{cc} 1 & 2 \\ 0 & 1 \\ 1 & 1 \end{array} \end{array} \quad (2)$$

Matrix $t_{Model} \cdot t_{Color}^\circ$ (counting)

counting sale records — corresponds to formula $t_A \cdot [T] \cdot t_B^\circ$ where the middle matrix is the identity.

Pivot Table Calculation

The outcome of cross-tabulation

$$ctab_{Color,Model;Sales} = t_{Color} \cdot [T]_{Sales} \cdot t_{Model}^{\circ} \quad (3)$$

solely using matrix operations is the desired pivot table:

$$t_{Color} \cdot [T]_{Sales} \cdot t_{Model}^{\circ} = \begin{array}{c} \begin{array}{cc} & \begin{array}{cc} Chevy & Ford \end{array} \\ \hline \begin{array}{c} Blue \\ Green \\ Red \end{array} & \begin{array}{cc} 87 & 106 \\ 0 & 64 \\ 5 & 8 \end{array} \end{array} \end{array} \quad (4)$$

Grand Totals (ALL) still missing

- Easily obtained via “bang” matrices ($!_A$)
- Matrix counterpart of the “bang” function (unique function to singleton type), that is, matrix $1 \xleftarrow{!_A} |A|$ wholly filled up with 1s.

Using matrix **block notation** cross tabulation (with totals) becomes:

$$\begin{aligned}
 ctab_{A,B;M} &: |A| + 1 \leftarrow |B| + 1 \\
 ctab_{A,B;M} &= \begin{bmatrix} t_A \\ \hline 1 \end{bmatrix} \cdot [T]_M \cdot \begin{bmatrix} t_B \\ \hline 1 \end{bmatrix}^\circ
 \end{aligned} \tag{5}$$

“Hello World” illustration

$$\begin{bmatrix} t_{Color} \\ \vdots \end{bmatrix} \cdot \llbracket T \rrbracket_{Sales} \cdot \begin{bmatrix} t_{Model} \\ \vdots \end{bmatrix}^{\circ} = \begin{array}{rcccl} & \textit{Chevy} & \textit{Ford} & \textit{ALL} & \\ \textit{Blue} & 87 & 106 & 193 & \\ \textit{Green} & 0 & 64 & 64 & \\ \textit{Red} & 5 & 8 & 13 & \\ \textit{ALL} & 92 & 178 & 270 & \end{array}$$

Sum of Sales	Model		
Color	Chevy	Ford	Grand Total
Blue	87	106	193
Green		64	64
Red	5	8	13
Grand Total	92	178	270

Incremental OLAP (proving things)

- Let T be yesterday's raw data and T' be the today's data.
- Assume that T has remained the same (no updates, no deletes).
- Let $T'' = T; T'$ denote the two data sources appended. Then the following facts hold:

$$t''_A = [t_A | t'_A] \quad (6)$$

$$t''_B = [t_B | t'_B] \quad (7)$$

$$[[T; T']]_M = [T]_M \oplus [T']_M \quad (8)$$

where \oplus denotes the direct sum of two matrices.

Let us prove that cross tabulation is incremental:

$$ctab_{A,B;M}(T; T') = ctab_{A,B;M}T + ctab_{A,B;M}T' \quad (9)$$

Calculational proof

$$ctab_{A,B;M}(T; T')$$

$$\Leftrightarrow \{ (5); \text{ totals off with no loss of generality} \}$$

$$t''_A \cdot [T; T']_M \cdot (t''_B)^\circ$$

$$\Leftrightarrow \{ (6); (7) \text{ and } (8) \}$$

$$[t_A|t'_A] \cdot ([T]_M \oplus [T']_M) \cdot [t_B|t'_B]^\circ$$

$$\Leftrightarrow \{ \text{absorption} \}$$

$$[t_A \cdot [T]_M | t'_A \cdot [T']_M] \cdot \left[\frac{t_B^\circ}{(t'_B)^\circ} \right]$$

$$\Leftrightarrow \{ \text{divide \& conquer matrix multiplication} \}$$

$$t_A \cdot [T]_M \cdot t_B^\circ + t'_A \cdot [T']_M \cdot (t'_B)^\circ$$

$$\Leftrightarrow \{ (5) \text{ twice} \}$$

$$ctab_{A,B;M} T + ctab_{A,B;M} T'$$

OLAP Cube (parallel) construction

Thus far

Summary generation in “human readable” format:
cross-tabulations are 2D charts.

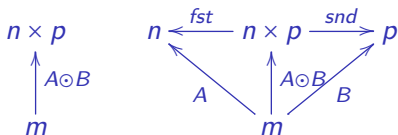
Higher dimensions

Generation of cubic and hypercubic data summaries also
captured by our typed LA approach.

We have to introduce some notion of dimension **product**.

Khatri-Rao matrix product

Given matrices $n \xleftarrow{A} m$ and $p \xleftarrow{B} m$, build the Khatri-Rao product of A and B ,



as follows,

$$\begin{aligned}
 u \odot v &= u \otimes v \\
 [A_1|A_2] \odot [B_1|B_2] &= [A_1 \odot B_1 | A_2 \odot B_2]
 \end{aligned}
 \tag{10}$$

where u, v are column-vectors and A_i, B_i are suitably typed matrices.

Example: measures in Khatri-Rao

As an example of Khatri-Rao product operation, consider row vector

$$s = [5 \quad 87 \quad 64 \quad 99 \quad 8 \quad 7]$$

of type $1 \xleftarrow{s} 6$ capturing the transposition of the *Sales* column. Then Khatri-Rao product $s \odot id$ is the corresponding diagonal matrix:

$$6 \xleftarrow{s \odot id} 6 = \begin{bmatrix} 5 & 0 & 0 & 0 & 0 & 0 \\ 0 & 87 & 0 & 0 & 0 & 0 \\ 0 & 0 & 64 & 0 & 0 & 0 \\ 0 & 0 & 0 & 99 & 0 & 0 \\ 0 & 0 & 0 & 0 & 8 & 0 \\ 0 & 0 & 0 & 0 & 0 & 7 \end{bmatrix} \quad (11)$$

“Hello World” illustration

Back to our running example, recall projections

$$t_{Model} : |Model| \longleftarrow 6$$

$$t_{Model} = \begin{array}{c|cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline Chevy & 1 & 1 & 0 & 0 & 0 & 0 \\ Ford & 0 & 0 & 1 & 1 & 1 & 1 \end{array}$$

and

$$t_{Color} : |Color| \longleftarrow 6$$

$$t_{Color} = \begin{array}{c|cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline Blue & 0 & 1 & 0 & 1 & 0 & 1 \\ Green & 0 & 0 & 1 & 0 & 0 & 0 \\ Red & 1 & 0 & 0 & 0 & 1 & 0 \end{array}$$

Pairing up dimensions

The Khatri-Rao product of t_{Model} and t_{Color} is matrix

	1	2	3	4	5	6
<i>Chevy Blue</i>	0	1	0	0	0	0
<i>Chevy Green</i>	0	0	0	0	0	0
<i>Chevy Red</i>	1	0	0	0	0	0
<i>Ford Blue</i>	0	0	0	1	0	1
<i>Ford Green</i>	0	0	1	0	0	0
<i>Ford Red</i>	0	0	0	0	1	0

of type $|Model| \times |Color| \longleftarrow n$.

It tells in which rows the particular pairs of values turn up.

In other words, this matrix is the projection $t_{Model \times Color}$ of the Cartesian product of the two dimensions. In general:

$$t_{A \times B} = t_A \odot t_B$$

Multi-dimensional summaries

Multidimensional cross-tabulations are obtained via the same formula (5) just by supplying higher-rank projections, for instance

	1990	1991	ALL
<i>Chevy Blue</i>	87	0	87
<i>Chevy Green</i>	0	0	0
<i>Chevy Red</i>	5	0	5
<i>Ford Blue</i>	99	7	106
<i>Ford Green</i>	64	0	64
<i>Ford Red</i>	0	8	8
ALL	255	15	270

corresponding to $A = \text{Model} \times \text{Color}$ and $B = \text{Year}$ in (5).

Computing the whole cube

General formula

$$\bigoplus_{i=0}^{\#D} \left(\bigoplus_{j \in \binom{D}{i}} \left(\bigodot_{d \in j} t_d \cdot [T]_{Unit} \cdot !^\circ \right) \right) \quad (13)$$

for

- *Unit* is the chosen measure (quantitative/numerical attribute),
- $\bigodot_i A_i$ iterates $A_1 \odot A_2$ to more than two arguments — mind that

$$! \odot A = A = A \odot ! \quad (14)$$

holds

- $\bigoplus_i A_i$ is the n -ary extension of the vertical blocking combinator $\begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$

MATLAB script

```
function C = Cube(proj,dnum,ndim,lines)
    C = [];
    for i=1:ndim
        ind = nchoosek(1:ndim,i);
        for j=1:size(ind,1)
            C = [C ; kr(proj{ind(j,:)}) * dnum * bang(lines)'];
        end
    end
    C = lift([C; bang(lines) * dnum * bang(lines)']);
end
```

By running

```
>> Cube({m,y,c},d,3,6)
```

in MATLAB, where variables m , y , c and d respectively hold t_{Model} , t_{Year} , t_{Color} , $[T]_{Sales}$ we will obtain the cube previously shown.

Summary

Summing up:

- A **no-SQL** approach to data mining
- Formal **semantics** implicit in LA encoding
- Other OLAP operations such as **roll-up** easy to implement in typed LA.

Moreover:

- All constructions in the approach **embarrassingly parallel** (Foster, 1995).
- Projection and diagonal matrices are **sparse**, therefore calling for suitably optimization in a parallel environment (Williams et al., 2009).

Promises **inexpensive** parallel implementation of OLAP/data mining in multi-core, **lap-top** machines, eg. on top of MS Excel or OpenOffice.

Putting ideas on paper

Draft paper

Details in http://alfa.di.uminho.pt/~hmacedo/wiki/doku.php?id=blog:2011:0411_do_the_middle_letters

We regard this as a practical application of the typed LA approach we are developing under the “**matrices as arrows**” motto (Macedo and Oliveira, 2010)

Ian Foster. *Designing and Building Parallel Programs: Concepts and Tools for Parallel Software Engineering*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1995. ISBN 0201575949.

H.D. Macedo and J.N. Oliveira. Matrices As Arrows! A Biproduct Approach to Typed Linear Algebra. In *Mathematics of Program Construction*, volume 6120 of *Lecture Notes in Computer Science*, pages 271–287. Springer, 2010.

Samuel Williams, Leonid Oliker, Richard Vuduc, John Shalf, Katherine Yelick, and James Demmel. Optimization of sparse matrix-vector multiplication on emerging multicore platforms. *Parallel Comput.*, 35:178–194, March 2009. ISSN 0167-8191. doi: 10.1016/j.parco.2008.12.006. URL <http://portal.acm.org/citation.cfm?id=1513001.1513318>.