

Aggregation and Ordering in Factorized Databases



Bakibayev, **K**očiský, **O**lteanu, and **Z**ávodný
University of Oxford

VLDB Sept 2, 2014

<http://www.cs.ox.ac.uk/projects/FDB/>

Outline



What are Factorized Databases?

Applications

A Glimpse at Aggregating Factorized Data

Factorized Databases by Example

Orders			Pizzas		Items	
customer	day	pizza	pizza	item	item	price
Mario	Monday	Capricciosa	Capricciosa	base	base	6
Mario	Friday	Capricciosa	Capricciosa	ham	ham	1
Pietro	Friday	Hawaii	Capricciosa	mushrooms	mushrooms	1
Lucia	Friday	Hawaii	Hawaii	base	pineapple	2
			Hawaii	ham		
			Hawaii	pineapple		

Consider the natural join of the three relations above:

Orders \bowtie Pizzas \bowtie Items					
customer	day	pizza	item	price	
Mario	Monday	Capricciosa	base	6	
Mario	Monday	Capricciosa	ham	1	
Mario	Monday	Capricciosa	mushrooms	1	
Mario	Friday	Capricciosa	base	6	
Mario	Friday	Capricciosa	ham	1	
Mario	Friday	Capricciosa	mushrooms	1	
...	

Factorized Databases by Example

Orders \bowtie Pizzas \bowtie Items				
customer	day	pizza	item	price
Mario	Monday	Capricciosa	base	6
Mario	Monday	Capricciosa	ham	1
Mario	Monday	Capricciosa	mushrooms	1
Mario	Friday	Capricciosa	base	6
Mario	Friday	Capricciosa	ham	1
Mario	Friday	Capricciosa	mushrooms	1
...

A *flat* relational algebra expression encoding the above query result is:

$$\begin{aligned} &\langle \text{Mario} \rangle \times \langle \text{Monday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{base} \rangle \times \langle 6 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Monday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{ham} \rangle \times \langle 1 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Monday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{mushrooms} \rangle \times \langle 1 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Friday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{base} \rangle \times \langle 6 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Friday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{ham} \rangle \times \langle 1 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Friday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{mushrooms} \rangle \times \langle 1 \rangle \cup \dots \end{aligned}$$

It uses relational product (\times), union (\cup), and singleton relations (e.g., $\langle 1 \rangle$).

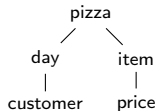
- The attribute names are not shown to avoid clutter.

Factorized Databases by Example

The previous relational expression entails lots of redundancy due to the joins:

$$\begin{aligned} &\langle \text{Mario} \rangle \times \langle \text{Monday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{base} \rangle \times \langle 6 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Monday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{ham} \rangle \times \langle 1 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Monday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{mushrooms} \rangle \times \langle 1 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Friday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{base} \rangle \times \langle 6 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Friday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{ham} \rangle \times \langle 1 \rangle \cup \\ &\langle \text{Mario} \rangle \times \langle \text{Friday} \rangle \times \langle \text{Capricciosa} \rangle \times \langle \text{mushrooms} \rangle \times \langle 1 \rangle \cup \dots \end{aligned}$$

We can *factorize* the expression following the join structure, e.g.:

$$\begin{aligned} &\langle \text{Capricciosa} \rangle \times ((\langle \text{Monday} \rangle \times \langle \text{Mario} \rangle \cup \langle \text{Friday} \rangle \times \langle \text{Mario} \rangle) \\ &\quad \times (\langle \text{base} \rangle \times \langle 6 \rangle \cup \langle \text{ham} \rangle \times \langle 1 \rangle \cup \langle \text{mushrooms} \rangle \times \langle 1 \rangle)) \\ &\cup \langle \text{Hawaii} \rangle \times \langle \text{Friday} \rangle \times ((\langle \text{Lucia} \rangle \cup \langle \text{Pietro} \rangle) \\ &\quad \times (\langle \text{base} \rangle \times \langle 6 \rangle \cup \langle \text{ham} \rangle \times \langle 1 \rangle \cup \langle \text{pineapple} \rangle \times \langle 2 \rangle)) \end{aligned}$$


There are several *algebraically equivalent* factorized representations defined by distributivity of product over union and commutativity of product and union.

Properties of Factorized Representations

Factorized representations of results of queries with select, project, join, aggregate, groupby, and orderby operators:

- **Very high compression rate**

- ▶ Can be exponentially more succinct than the relations they encode.
- ▶ Arbitrarily better than generic compression schemes, e.g., bzip2
- ▶ Factorized representations of asymptotically-tight size bounds computable directly from input database and query

- **Querying in the compressed domain**

- ▶ Factorizations are relational expressions
- ▶ We developed the FDB in-memory query engine for this purpose

- **Constant-delay enumeration of represented tuples**

- ▶ Tuple iteration as fast as listing them from equivalent flat relations

Outline



What are Factorized Databases?

Applications

A Glimpse at Aggregating Factorized Data

Spot the Factorized Database!

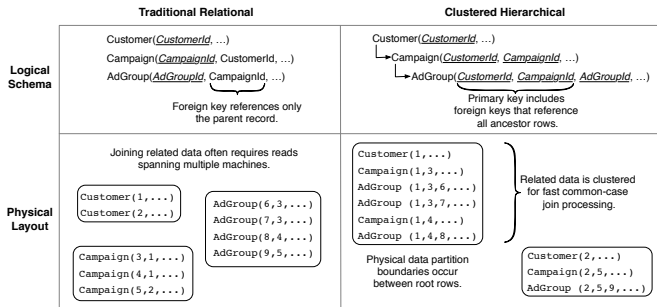


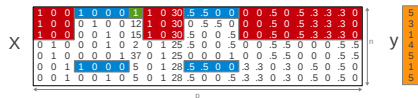
Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

Excerpt from *F1: A Distributed SQL Database That Scales*. PVLDB'13.

- Google's DB supporting their lucrative AdWords business
- Database factorization increases data locality for common access patterns
 - ▶ Tables pre-joined using a nesting structure defined by key-fkey constraints
- Data partitioned across servers into factorization fragments

Spot the Factorized Database!

(a) Training Data in Numeric Format (Design Matrix)



(b) Block Structure Representation of Design Matrix

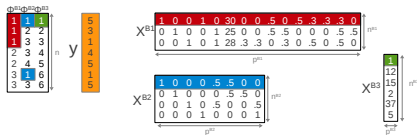


Figure 3: (a) In relational domains, design matrices X have large blocks of repeating patterns (example from Figure 2). (b) Repeating patterns in X can be formalized by a block notation (see section 2.3) which stems directly from the relational structure of the original data. Machine learning methods have to make use of repeating patterns in X to scale to large relational datasets.

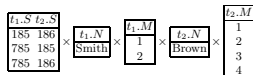
Excerpt from *Scaling Factorization Machines to Relational Data*. PVLDB'13.

- Feature vectors for predictive modelling represented as very large design matrices (= relations with high cardinality)
- Standard learning algorithms cannot scale on design matrix representation
- Use repeating patterns in the design matrix as key to scalability

Spot the Factorized Database!

Social Security Number: <u>785</u> Name: <u>Smith</u> Marital Status: (1) single <input checked="" type="checkbox"/> (2) married <input checked="" type="checkbox"/> (3) divorced <input type="checkbox"/> (4) widowed <input type="checkbox"/>	$t_{1,S}$ $t_{1,N}$ $t_{1,M}$ $t_{2,S}$ $t_{2,N}$ $t_{2,M}$ 185 Smith 1 186 Brown 1 185 Smith 1 186 Brown 2 185 Smith 1 186 Brown 3 185 Smith 1 186 Brown 4 185 Smith 2 186 Brown 1 185 Smith 2 186 Brown 2 185 Smith 2 186 Brown 3 185 Smith 2 186 Brown 4 ⋮ 785 Smith 2 186 Brown 4
Social Security Number: <u>185</u> Name: <u>Brown</u> Marital Status: (1) single <input type="checkbox"/> (2) married <input type="checkbox"/> (3) divorced <input type="checkbox"/> (4) widowed <input type="checkbox"/>	

Fig. 1. Two completed survey forms and a world-set relation representing the possible worlds with unique social security numbers.



Excerpt from 10^{10^6} *Worlds and Beyond: Efficient Representation and Processing of Incomplete Information*. ICDE'07.

Managing a large set of possibilities or choices:

- Configuration problems (space of valid solutions)
- Incomplete information (space of possible worlds)

Spot the Factorized Database!

98 5. INTENSIONAL QUERY EVALUATION

5.1.3 READ-ONCE FORMULAS

An important class of propositional formulas that play a special role in probabilistic databases are read-once formulas. We restrict our discussion to the case when all random variables X are Boolean variables.

Φ is called *read-once* if there is a formula Φ' equivalent to Φ such that every variable occurs at most once in Φ' . For example:

$$\Phi = X_1 Y_1 \vee X_1 Y_2 \vee X_2 Y_3 \vee X_2 Y_4 \vee X_2 Y_5$$

is read-once because it is equivalent to the following formula:

$$\Phi' = X_1 (Y_1 \vee Y_2) \vee X_2 (Y_3 \vee Y_4 \vee Y_5)$$

Excerpt from *Probabilistic Databases*. Morgan & Claypool. 2011.

Provenance and probabilistic data:

- Compact encoding for large provenance
- Factorization of provenance is used for efficient query evaluation in probabilistic databases

Outline



What are Factorized Databases?

Applications

A Glimpse at Aggregating Factorized Data

Aggregating Factorized Data

We only present here COUNT and SUM aggregation functions.

COUNT(F) is the number of tuples in a factorization F :

- $\text{COUNT}(\langle a \rangle) = 1$.
- $\text{COUNT}(F_1 \cup \dots \cup F_k) = \text{COUNT}(F_1) + \dots + \text{COUNT}(F_k)$.
- $\text{COUNT}(F_1 \times \dots \times F_k) = \text{COUNT}(F_1) \cdot \dots \cdot \text{COUNT}(F_k)$.

SUM _{A} (F) is the sum of all values of attribute A in a factorization F :

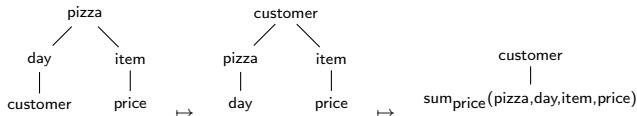
- $\text{SUM}_A(\langle a \rangle) = a$, if the singleton $\langle a \rangle$ has attribute A .
- $\text{SUM}_A(F_1 \cup \dots \cup F_k) = \text{SUM}_A(F_1) + \dots + \text{SUM}_A(F_k)$.
- $\text{SUM}_A(F_1 \times \dots \times F_k) = \text{SUM}_A(F_1) \cdot \text{COUNT}(F_2) \cdot \dots \cdot \text{COUNT}(F_k)$,
where wlog values for attribute A are in expression F_1 .

Aggregation by Example

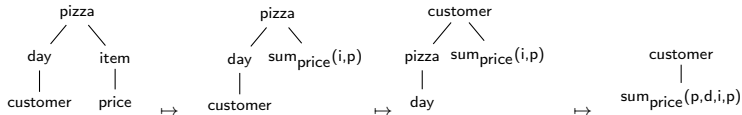
- Recall the natural join of Orders, Pizzas, and Items
- We would like to find the overall sales per customer
- Assume the factorization structure discussed before (leftmost below)

Example of possible evaluation plans:

1. First restructure for GROUP-BY, then aggregate

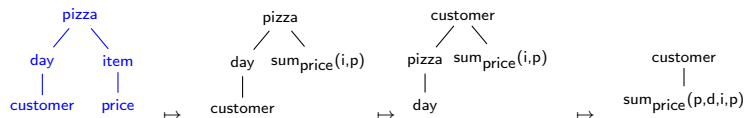


2. Intertwine restructuring for GROUP-BY and partial aggregation



Query Evaluation Step by Step

Let us consider the second evaluation plan:

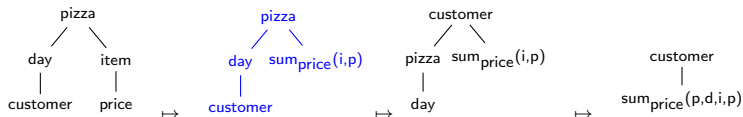


The initial factorization with the structure highlighted above:

$$\begin{aligned} & \langle \text{Capricciosa} \rangle \times (\langle \text{Monday} \rangle \times \langle \text{Mario} \rangle \cup \langle \text{Friday} \rangle \times \langle \text{Mario} \rangle) \\ & \quad \times (\langle \text{base} \rangle \times \langle 6 \rangle \cup \langle \text{ham} \rangle \times \langle 1 \rangle \cup \langle \text{mushrooms} \rangle \times \langle 1 \rangle) \\ \cup & \langle \text{Hawaii} \rangle \times \langle \text{Friday} \rangle \times (\langle \text{Lucia} \rangle \cup \langle \text{Pietro} \rangle) \\ & \quad \times (\langle \text{base} \rangle \times \langle 6 \rangle \cup \langle \text{ham} \rangle \times \langle 1 \rangle \cup \langle \text{pineapple} \rangle \times \langle 2 \rangle) \end{aligned}$$

Query Evaluation Step by Step

Let us consider the second evaluation plan:

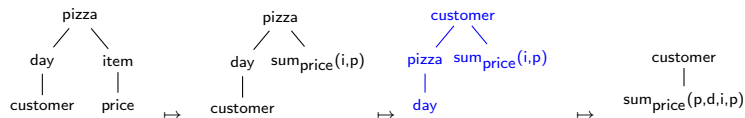


The factorization after partial aggregation with the structure highlighted above:

$$\begin{aligned} & \langle \text{Capricciosa} \rangle \times (\langle \text{Monday} \rangle \times \langle \text{Mario} \rangle \cup \langle \text{Friday} \rangle \times \langle \text{Mario} \rangle) \\ & \quad \times \langle 8 \rangle \\ \cup & \langle \text{Hawaii} \rangle \times \langle \text{Friday} \rangle \times (\langle \text{Lucia} \rangle \cup \langle \text{Pietro} \rangle) \\ & \quad \times \langle 9 \rangle \end{aligned}$$

Query Evaluation Step by Step

Let us consider the second evaluation plan:

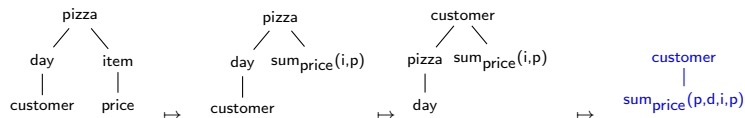


The factorization after restructuring with the structure highlighted above:

$$\begin{aligned} & \langle Lucia \rangle \times \langle Hawaii \rangle \times \langle Friday \rangle \times \langle 9 \rangle \cup \\ & \langle Mario \rangle \times \langle Capricciosa \rangle \times (\langle Monday \rangle \cup \langle Friday \rangle) \times \langle 8 \rangle \cup \\ & \langle Pietro \rangle \times \langle Hawaii \rangle \times \langle Friday \rangle \times \langle 9 \rangle \end{aligned}$$

Query Evaluation Step by Step

Let us consider the second evaluation plan:



The factorization after final aggregation with the structure highlighted above:

$$\begin{aligned} &\langle Lucia \rangle \times \langle 9 \rangle \cup \\ &\langle Mario \rangle \times \langle 16 \rangle \cup \\ &\langle Pietro \rangle \times \langle 9 \rangle \end{aligned}$$

Thank you!