From Joins to Aggregates

and Optimization Problems



Dan Olteanu (Oxford & Turing)

https://fdbresearch.github.io

Data Science Class Alan Turing Institute Jan 29, 2018

Acknowledgements

Some work reported in this tutorial has been done in the context of the FDB project, LogicBlox, and RelationalAI by

- Zavodný, Schleich, Kara, Ciucanu, and myself (Oxford)
- Abo Khamis and Ngo (RelationalAI), Nguyen (U. Michigan)

Some of the following slides are derived from presentations by

- Abo Khamis (in-db analytics diagrams)
- Kara (covers and various graphics)
- Ngo (FAQ)
- Schleich (performance and quizzes)

Lastly, Kara and Schleich proofread the slides.

I would like to thank them for their support and for sharing their work!

Goal of This Tutorial

Introduction to a principled, relatively new approach to in-database computation

It starts where mainstream introductory/advanced courses on databases finish.

- Joins
 - Worst-case optimal join algorithms
 - Listing vs. factorized representations of join results
- Aggregates
 - Generalization of join algorithms to aggregates over joins
 - Functional aggregate queries with applications in, e.g., DB, logic, probabilistic graphical models, matrix chain computation
 - New algorithms with low computational complexity
- Optimizations
 - In-database learning of regression and classification models

Quizzes: Test your understanding after class

Outline



Part 1. Joins

Part 2. Aggregates

Part 3. Optimization

Join Queries

Basic building blocks in query languages. Studied extensively.

However, worst-case optimal join algorithms were only proposed recently.
[NPRR12,NRR13,V14,OZ15,ANS17]

Likewise for systematic investigation of redundancy in the computation and representation of join results. [OZ12,OZ15,KO17]

This tutorial highlights recent work on worst-case optimal join algorithms under listing and factorized data representations.

Plan for Part 1 on Joins

- Introduction to join queries via examples
- Size bounds for results of join queries
 - Standard (exhaustive) listing representation
 - Factorized (succinct) representations
- Worst-case optimal join algorithms
 - ► LFTJ (LeapFrog TrieJoin) used by LogicBlox for listing representation
 - ▶ FDB (Factorized Databases) for factorized representations

Introduction to Join Queries

Join Example: Itemized Customer Orders

Orders (O for short)			Dish (D	Dish (D for short)		Items (I for short)		
customer	day	dish	dish	item	item	price		
Elise	Monday	burger	burger	patty	patty	6		
Elise	Friday	burger	burger	onion	onion	2		
Steve	Friday	hotdog	burger	bun	bun	2		
Joe	Friday	hotdog	hotdog	bun	sausage	4		
			hotdog	onion				
			hotdog	sausage				

Consider the natural join of the above relations:

O(customer, day, dish), D(dish, item), I(item, price)

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2

Join Example: Listing the Triangles in the Database

R_1		R_2		F	R_3		$R_1(\boldsymbol{A},\boldsymbol{B}), R_2(\boldsymbol{A},\boldsymbol{C}), R_3(\boldsymbol{B},\boldsymbol{C})$			
Α	В	A	С	В	С	-	Α	В	С	-
a 0	b_0	a 0	<i>C</i> 0	b 0	C 0	-	a 0	b 0	<i>C</i> 0	•
a 0		a 0		b 0			a_0	b 0		
a 0	bm	a 0	Cm	b 0	Cm		a 0	b 0	Cm	
a 1	b 0	a 1	<i>c</i> ₀	b_1	<i>c</i> ₀	-	a 0	b_1	C 0	-
	b_0		<i>C</i> 0		<i>C</i> 0		a 0		C 0	
a _m	b_0	am	<i>c</i> ₀	b_m	<i>C</i> ₀		a_0	b_m	<i>c</i> ₀	
						-	a_1	b 0	C 0	-

b0

 a_1

 b_0

*c*₀

C0

Join Hypergraphs

We associate a hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ with every join query Q

- Each variable in Q corresponds to a node in $\mathcal V$
- Each relation symbol in Q corresponds to a (hyper)edge in $\mathcal E$

Example: Triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$



•
$$\mathcal{V} = \{A, B, C\}$$

• $\mathcal{E} = \{\{A, B\}, \{A, C\}, \{B, C\}\}$

Join Hypergraphs

We associate a hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ with every join query Q

- Each variable in Q corresponds to a node in $\mathcal V$
- Each relation symbol in Q corresponds to a (hyper)edge in \mathcal{E}

Example: Order query O(cust, day, dish), D(dish, item), I(item, price)



- $\mathcal{V} = \{ \texttt{cust}, \texttt{day}, \texttt{dish}, \texttt{item}, \texttt{price} \}$
- $\blacksquare \ \mathcal{E} = \{ \{ \texttt{cust}, \texttt{day}, \texttt{dish} \}, \{ \texttt{dish}, \texttt{item} \}, \{ \texttt{item}, \texttt{price} \} \}$

Hypertree Decompositions

Definition[GLS99]: A (hypertree) decomposition \mathcal{T} of the hypergraph (\mathcal{V}, \mathcal{E}) of a query Q is a pair (\mathcal{T}, χ), where

- T is a tree
- χ is a function mapping each node in T to a subset of \mathcal{V} called *bag*.

Properties of a decomposition \mathcal{T} :

- Coverage: $\forall e \in \mathcal{E}$, there must be a node $t \in \mathcal{T}$ such that $e \subseteq \chi(t)$.
- Connectivity: $\forall v \in \mathcal{V}$, $\{t \mid t \in T, v \in \chi(t)\}$ forms a connected subtree.

The hypergraph of the query $R_1(A, B), R_2(B, C), R_3(C, D)$

A hypertree decomposition



Hypertree Decompositions

Definition[GLS99]: A (hypertree) decomposition \mathcal{T} of the hypergraph (\mathcal{V}, \mathcal{E}) of a query Q is a pair (\mathcal{T}, χ), where

T is a tree

• χ is a function mapping each node in T to a subset of \mathcal{V} called *bag*.

Properties of the decomposition \mathcal{T} :

• Coverage: $\forall e \in \mathcal{E}$, there must be a node $t \in \mathcal{T}$ such that $e \subseteq \chi(t)$.

■ Connectivity: $\forall v \in V$, $\{t \mid t \in T, v \in \chi(t)\}$ forms a connected subtree.

The hypergraph of the triangle query A hypertree decomposition $R_1(A, B), R_2(A, C), R_3(B, C)$





Variable Orders

Definition[OZ15]: A variable order Δ for a query Q is a pair (F, key), where

- F is a rooted forest with one node per variable in Q
- key is a function mapping each variable A to a subset of its ancestor variables in F.

Properties of a variable order Δ for Q:

For each relation symbol, its variables lie along the same root-to-leaf path in *F*. For any such variables *A* and *B*, $A \in key(B)$ if *A* is an ancestor of *B*.

• For every child B of A, $key(B) \subseteq key(A) \cup \{A\}$.

Possible variable orders for the path query $R_1(A, B), R_2(B, C), R_3(C, D)$:

$$A \quad key(A) = \emptyset$$

$$B \quad key(B) = \{A\}$$

$$key(A) = \{B\} A$$

$$C \quad key(C) = \{B\}$$

$$D \quad key(D) = \{C\}$$

$$B \quad key(A) = \{B\} A$$

$$C \quad key(C) = \{B\}$$

$$D \quad key(D) = \{C\}$$

Variable Orders

Definition[OZ15]: A variable order Δ for a query Q is a pair (F, key), where

- F is a rooted forest with one node per variable in Q
- key is a function mapping each variable A to a subset of its ancestor variables in F.

Properties of a variable order Δ for Q:

- For each relation symbol, its variables lie along the same root-to-leaf path in *F*. For any such variables *A* and *B*, *A* ∈ key(*B*) if *A* is an ancestor of *B*.
- For every child B of A, $key(B) \subseteq key(A) \cup \{A\}$.

Possible variable orders for the triangle query $R_1(A, B)$, $R_2(A, C)$, $R_3(B, C)$:

$$\begin{array}{ccccccc} A & key(A) = \emptyset & B & key(B) = \emptyset & C & key(C) = \emptyset \\ & & & & \\ B & key(B) = \{A\} & A & key(A) = \{B\} & B & key(B) = \{C\} \\ & & & & \\ C & key(C) = \{A, B\} & C & key(C) = \{A, B\} & A & key(A) = \{B, C\} \end{array}$$

Hypertree Decompositions ⇔ Variable Orders

From variable order Δ to hypertree decomposition \mathcal{T} : [OZ15]

- For each node A in Δ , create a bag $key(A) \cup \{A\}$.
- The bag for A is connected to the bags for its children and parent.
- Optionally, remove redundant bags

Example: Triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$

$$A \quad key(A) = \emptyset \qquad A$$

$$B \quad key(B) = \{A\} \Rightarrow A, B \Rightarrow A, B, C$$

$$C \quad key(C) = \{A, B\} \qquad A, B, C$$

From variable order Δ to hypertree decomposition \mathcal{T} :

- For each node A in Δ , create a bag $key(A) \cup \{A\}$.
- The bag for A is connected to the bags for its children and parent.
- Optionally, remove redundant bags

Example: Path query $R_1(A, B), R_2(B, C), R_3(C, D)$

$$A \quad key(A) = \emptyset$$

$$B \quad key(B) = \{A\} \Rightarrow A, B \Rightarrow A, B$$

$$C \quad key(C) = \{B\}$$

$$D \quad key(D) = \{C\}$$

$$C, D$$

$$A \quad A, B \Rightarrow A, B$$

$$B, C \quad B, C$$

$$C, D$$

From hypertree decomposition $\mathcal T$ to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in $\mathcal T$
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$

From hypertree decomposition $\mathcal T$ to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in \mathcal{T}
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$

Step 1: A is removed from Tand inserted into Δ



A $kev(A) = \emptyset$

From hypertree decomposition $\mathcal T$ to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in \mathcal{T}
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$



From hypertree decomposition \mathcal{T} to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in \mathcal{T}
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$



From hypertree decomposition $\mathcal T$ to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in $\mathcal T$
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Path query $R_1(A, B), R_2(B, C), R_3(C, D)$



From hypertree decomposition \mathcal{T} to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in $\mathcal T$
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Path query $R_1(A, B), R_2(B, C), R_3(C, D)$



From hypertree decomposition \mathcal{T} to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in \mathcal{T}
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Path query $R_1(A, B), R_2(B, C), R_3(C, D)$



From hypertree decomposition \mathcal{T} to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in $\mathcal T$
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Path query $R_1(A, B), R_2(B, C), R_3(C, D)$



From hypertree decomposition \mathcal{T} to variable order Δ :

- Create a node A in Δ for a variable A in the top bag in \mathcal{T}
- Recurse with \mathcal{T} where A is removed from all bags in \mathcal{T} .
- If top bag empty, then recurse independently on each of its child bags and create children of A in Δ
- Update *key* for each variable at each step.

Example: Path query $R_1(A, B), R_2(B, C), R_3(C, D)$



Size Bounds for Listing Representation of Join Results

Example: the path query $R_1(A, B), R_2(B, C), R_3(C, D)$

Assumption: All relations have size N.

- The result is included in the result of $R_1(A, B), R_3(C, D)$
 - ▶ Its size is upper bounded by $N^2 = |R_1| \times |R_3|$
 - ▶ All variables are "covered" by the relations R₁ and R₃

• There are databases for which the result size is at least N^2

• Let $R_1 = [N] \times \{1\}, R_2 = \{1\} \times [N], R_3 = [N] \times \{1\}.$

Example: the path query $R_1(A, B), R_2(B, C), R_3(C, D)$

Assumption: All relations have size N.

- The result is included in the result of $R_1(A, B), R_3(C, D)$
 - ▶ Its size is upper bounded by $N^2 = |R_1| \times |R_3|$
 - ▶ All variables are "covered" by the relations R₁ and R₃

• There are databases for which the result size is at least N^2

• Let
$$R_1 = [N] \times \{1\}, R_2 = \{1\} \times [N], R_3 = [N] \times \{1\}.$$

• Conclusion: Size of the query result is $\Theta(N^2)$ for some inputs

Example: the triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$

Assumption: All relations have size N.

• The result is included in the result of $R_1(A, B), R_3(B, C)$

- Its size is upper bounded by $N^2 = |R_1| \times |R_3|$
- ▶ All variables are "covered" by the relations R₁ and R₃

• There are databases for which the result size is at least N

▶ Let $R_1 = [N] \times \{1\}, R_2 = [N] \times \{1\}, R_3 \supset \{(1,1)\}$

Example: the triangle query $R_1(A, B), R_2(A, C), R_3(B, C)$

Assumption: All relations have size N.

• The result is included in the result of $R_1(A, B), R_3(B, C)$

- Its size is upper bounded by $N^2 = |R_1| \times |R_3|$
- All variables are "covered" by the relations R₁ and R₃

There are databases for which the result size is at least N

• Let $R_1 = [N] \times \{1\}, R_2 = [N] \times \{1\}, R_3 \supset \{(1,1)\}$

• Conclusion: Size gap between the N^2 upper bound and the N lower bound

Question: Can we close this gap and give tight size bounds?

Edge Covers and Independent Sets

We can generalize the previous examples as follows:

For the size upper bound:

- Cover all nodes (variables) by k edges (relations) \Rightarrow size $\leq N^k$.
- This is an edge cover of the query hypergraph!

For the size lower bound:

- *m* independent nodes \Rightarrow construct database such that size $\ge N^m$.
- This is an independent set of the query hypergraph!

 $\max_m = |\text{IndependentSet}(Q)| \le |\text{EdgeCover}(Q)| = \min_k$

 \max_m and \min_k do not necessarily meet!

Can we further refine this analysis?

The Fractional Edge Cover Number $\rho^*(Q)$

The two bounds meet if we take their fractional versions

- Fractional edge cover of Q with weight $k \Rightarrow \text{size} \le N^k$.
- Fractional independent set with weight $m \Rightarrow \exists$ database with size $\geq N^m$.

By duality of linear programming:

 $\max_{m} = |\operatorname{FractionalIndependentSet}(Q)| = |\operatorname{FractionalEdgeCover}(Q)| = \min_{k}$

■ This is the fractional edge cover number $\rho^*(Q)!$

For query Q and database of size N, the query result has size $O(N^{\rho^*(Q)})$.

[AGM08]

The Fractional Edge Cover Number $\rho^*(Q)$

For a join query $Q(\mathbf{A}_1 \cup \cdots \cup \mathbf{A}_n) = R_1(\mathbf{A}_1), \ldots, R_n(\mathbf{A}_n),$ $\rho^*(Q)$ is the cost of an optimal solution to the linear program:

$$\begin{array}{ll} \text{minimize} & \sum_{i \in [n]} x_{R_i} \\ \text{subject to} & \sum_{i: \text{edge } R_i \text{ covers node } A} x_{R_i} \geq 1 \ \forall A \in \bigcup_{j \in [n]} A_j, \\ & x_{R_i} \geq 0 \qquad \qquad \forall i \in [n]. \end{array}$$

- x_{R_i} is the weight of edge (relation) R_i in the hypergraph of Q
- Each node (variable) has to be covered by edges with sum of weights ≥ 1
- In the integer program variant for the edge cover, $x_{R_i} \in \{0,1\}$

Example of Fractional Edge Cover Computation (1)

Consider the join query Q: R(A, B, C), S(A, B, D), T(A, E), U(E, F).



- The three edges *R*, *S*, *U* to cover all nodes. FractionalEdgeCover(*Q*) ≤ 3
- Each node C, D, and F must be covered by a distinct edge. FractionalIndependentSet $(Q) \ge 3$

 $\Rightarrow \rho^*(Q) = 3$

 \Rightarrow Size $\leq N^3$ and for some inputs is $\Theta(N^3)$.

Example of Fractional Edge Cover Computation (2)

Consider the triangle query Q: $R_1(A, B), R_2(A, C), R_3(B, C)$.



Our previous size upper bound was N^2 :

• This is obtained by setting any two of $x_{R_1}, x_{R_2}, x_{R_3}$ to 1.

What is the fractional edge cover number for the triangle query?
Example of Fractional Edge Cover Computation (2)

Consider the triangle query Q: $R_1(A, B), R_2(A, C), R_3(B, C)$.



Our previous size upper bound was N^2 :

• This is obtained by setting any two of $x_{R_1}, x_{R_2}, x_{R_3}$ to 1.

What is the fractional edge cover number for the triangle query?

We can do better: $x_{R_1} = x_{R_2} = x_{R_3} = 1/2$. Then, $\rho^* = 3/2$.

Lower bound reaches $N^{3/2}$ for $R_1 = R_2 = R_3 = [\sqrt{N}] \times [\sqrt{N}]$.

Example of Fractional Edge Cover Computation (3)

Consider the (4-cycle) join: $R(A_1, A_2), S(A_2, A_3), T(A_3, A_4), W(A_4, A_1).$

The linear program for its fractional edge cover number:



Possible solution: $x_R = x_T = 1$. Another solution: $x_S = x_W = 1$. Then, $\rho^* = 2$.

Lower bound reaches N^2 for $R = T = [N] \times \{1\}$ and $S = W = \{1\} \times [N]$.

Common case in practice:

- Relations have different sizes
- Small-size projections of relations may be added to the join query

Recall the linear program for computing the fractional edge cover number $\rho^*(Q)$ of a join query $Q(A_1 \cup \cdots \cup A_n) = R_1(A_1), \ldots, R_n(A_n)$:

$$\begin{array}{ll} \text{minimize} & \sum_{i \in [n]} x_{R_i} \\ \text{subject to} & \sum_{i: \text{edge } R_i \text{ covers node } A} x_{R_i} \geq 1 \ \forall A \in \bigcup_{j \in [n]} A_j, \\ & x_{R_i} \geq 0 \qquad \qquad \forall i \in [n]. \end{array}$$

Common case in practice:

- Relations have different sizes
- Small-size projections of relations may be added to the join query

Add relation sizes into the linear program that computes the result size of a join query $Q(\mathbf{A}_1 \cup \cdots \cup \mathbf{A}_n) = R_1(\mathbf{A}_1), \ldots, R_n(\mathbf{A}_n)$:

$$\begin{array}{ll} \text{minimize} & N^{\sum_{i \in [n]} x_{R_{i}}} \\ \text{subject to} & \sum_{i: \text{edge } R_{i} \text{ covers node } A} x_{R_{i}} \geq 1 \ \forall A \in \bigcup_{j \in [n]} A_{j}, \\ & x_{R_{i}} \geq 0 \qquad \quad \forall i \in [n]. \end{array}$$

Assumption: All relations have the same size N.

Common case in practice:

- Relations have different sizes
- Small-size projections of relations may be added to the join query

Add relation sizes into the linear program that computes the result size of a join query $Q(\mathbf{A}_1 \cup \cdots \cup \mathbf{A}_n) = R_1(\mathbf{A}_1), \ldots, R_n(\mathbf{A}_n)$:

 $\begin{array}{ll} \text{minimize} & \prod_{i \in [n]} N^{x_i} \\ \text{subject to} & \sum_{i: \text{edge } R_i \text{ covers node } A} x_{R_i} \geq 1 \quad \forall A \in \bigcup_{j \in [n]} A_j, \\ & x_{R_i} \geq 0 \qquad \qquad \forall i \in [n]. \end{array}$

Assumption: All relations have the same size N.

Common case in practice:

- Relations have different sizes
- Small-size projections of relations may be added to the join query

Add relation sizes into the linear program that computes the result size of a join query $Q(\mathbf{A}_1 \cup \cdots \cup \mathbf{A}_n) = R_1(\mathbf{A}_1), \ldots, R_n(\mathbf{A}_n)$:

 $\begin{array}{ll} \text{minimize} & \prod_{i \in [n]} N_i^{x_i} \\ \text{subject to} & \sum_{i: \text{edge } R_i \text{ covers node } A} x_{R_i} \geq 1 \quad \forall A \in \bigcup_{j \in [n]} A_j, \\ & x_{R_i} \geq 0 \qquad \quad \forall i \in [n]. \end{array}$

Assumption: Relation R_i has size N_i , $\forall i \in [n]$.

Size Bounds for Factorized Representations of Join Results

Recall the Itemized Customer Orders Example

Orders (O for short)			Dish (D	for short)	Items (I for short)	
customer	day	dish	dish	item	item	price
Elise	Monday	burger	burger	patty	patty	6
Elise	Friday	burger	er burger onion		onion	2
Steve	Friday	hotdog	burger	bun	bun	2
Joe	Friday	hotdog	hotdog	bun	sausage	4
			hotdog	onion		
		hotdog	sausage			

Consider the natural join of the above relations:

O(customer, day, dish), D(dish, item), I(item, price)

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2

Factor Out Common Data Blocks

O(custom	er, day, dis	n), D(aisn	, item), i(item, price)
customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2

O(customer, day, dish), D(dish, item), I(item, price)

The listing representation of the above query result is:

$\langle \textit{Elise} \rangle$	×	$\langle Monday \rangle$	×	$\langle burger \rangle$	×	$\langle patty \rangle$	×	$\langle 6 \rangle$	U
$\langle \textit{Elise} \rangle$	×	$\langle \mathit{Monday} \rangle$	×	$\langle burger \rangle$	×	$\langle onion \rangle$	×	$\langle 2 \rangle$	U
$\langle \textit{Elise} \rangle$	×	$\langle Monday \rangle$	×	$\langle burger \rangle$	×	$\langle bun \rangle$	×	$\langle 2 \rangle$	U
$\langle \textit{Elise} \rangle$	×	$\langle Friday angle$	×	$\langle burger \rangle$	×	$\langle patty \rangle$	×	$\langle 6 \rangle$	U
$\langle Elise \rangle$	×	(<i>Friday</i>)	×	$\langle burger \rangle$	×	$\langle onion \rangle$	×	$\langle 2 \rangle$	U
$\langle Elise \rangle$	×	$\langle Friday angle$	×	$\langle burger \rangle$	×	(bun)	×	$\langle 2 \rangle$	υ

It uses relational product (\times), union (\cup), and data (singleton relations).

The attribute names are not shown to avoid clutter.

This is How A Factorized Join Looks Like!



Factorized representation of the join result

There are several algebraically equivalent factorized representations defined:

- by distributivity of product over union and their commutativity;
- as groundings of variable orders.

.. Now with Further Compression using Caching



Observation:

- price is under item, which is under dish, but only depends on item,
- .. so the same price appears under an item *regardless* of the dish.
- Idea: Cache price for a specific item and avoid repetition!

Same Data, Different Factorization



.. and Further Compressed using Caching



Which factorization should we choose?

The size of a factorization is the number of its values.

Example:

$$F_{1} = (\langle 1 \rangle \cup \cdots \cup \langle n \rangle) \times (\langle 1 \rangle \cup \cdots \cup \langle m \rangle)$$

$$F_{2} = \langle 1 \rangle \times \langle 1 \rangle \cup \cdots \cup \langle 1 \rangle \times \langle m \rangle$$

$$\cup \cdots \cup$$

$$\langle n \rangle \times \langle 1 \rangle \cup \cdots \cup \langle n \rangle \times \langle m \rangle.$$

*F*₁ is factorized, *F*₂ is a listing representation
 *F*₁ ≡ *F*₂
 BUT |*F*₁| = *m* + *n* ≪ |*F*₂| = *m* * *n*.

How much space does factorization save over the listing representation?

Given a join query Q, for any database of size N, the join result admits

• a listing representation of size $O(N^{\rho^*(Q)})$. [AGM08]

Given a join query Q, for any database of size N, the join result admits

• a listing representation of size $O(N^{\rho^*(Q)})$. [AGM08]

• a factorization without caching of size $O(N^{s(Q)})$. [OZ12]

Given a join query Q, for any database of size N, the join result admits

- a listing representation of size $O(N^{\rho^*(Q)})$. [AGM08]
- a factorization without caching of size $O(N^{s(Q)})$. [OZ12]
- a factorization with caching of size $O(N^{fhtw(Q)})$. [OZ15]

Given a join query Q, for any database of size N, the join result admits

- a listing representation of size $O(N^{\rho^*(Q)})$. [AGM08]
- a factorization without caching of size $O(N^{s(Q)})$. [OZ12]

• a factorization with caching of size $O(N^{fhtw(Q)})$. [OZ15]

$$1 \leq \mathit{fhtw}(Q) \underbrace{\leq}_{ ext{up to } \log |Q|} \mathit{s}(Q) \underbrace{\leq}_{ ext{up to } |Q|}
ho^*(Q) \leq |Q|$$

- |Q| is the number of relations in Q
- $\rho^*(Q)$ is the fractional edge cover number of Q
- s(Q) is the factorization width of Q
- fhtw(Q) is the fractional hypertree width of Q

[M10]

Given a join query Q, for any database of size N, the join result admits

- a listing representation of size $O(N^{\rho^*(Q)})$. [AGM08]
- a factorization without caching of size $O(N^{s(Q)})$. [OZ12]
- a factorization with caching of size $O(N^{fhtw(Q)})$. [OZ15]

These size bounds are asymptotically tight!

Best possible size bounds for factorized representations over variable orders of Q and for listing representation, but not <u>database</u> optimal!

There exists arbitrarily large databases for which

- the listing representation has size $\Omega(N^{\rho^*(Q)})$
- ▶ the factorization with/without caching over any variable order of Q has size $\Omega(N^{s(Q)})$ and $\Omega(N^{fhtw(Q)})$ respectively.

Example: The Factorization Width s



The structure of the factorization over the above variable order Δ :

$$\bigcup_{\mathbf{a}\in\mathbf{A}} \left(\langle \mathbf{a} \rangle \times \bigcup_{\mathbf{b}\in\mathbf{B}} \left(\langle \mathbf{b} \rangle \times \left(\bigcup_{c\in C} \langle \mathbf{c} \rangle \right) \times \left(\bigcup_{d\in D} \langle \mathbf{d} \rangle \right) \right) \times \bigcup_{e\in\mathbf{E}} \left(\langle \mathbf{e} \rangle \times \left(\bigcup_{f\in F} \langle f \rangle \right) \right) \right)$$

The number of values for a variable is dictated by the number of valid tuples of values for its ancestors in Δ :

• One value $\langle f \rangle$ for each tuple (a, e, f) in the join result.

Size of factorization = sum of sizes of results of subqueries along paths.

Example: The Factorization Width s



 \blacksquare The factorization width for Δ is the largest ρ^* over subqueries defined by root-to-leaf paths in Δ

• s(Q) is the minimum factorization width over all variable orders of QIn our example:

• Path A-E-F has fractional edge cover number 2.

 \Rightarrow The number of *F*-values is $\leq N^2$, but can be $\sim N^2$.

- All other root-to-leaf paths have fractional edge cover number 1.
 - \Rightarrow The number of other values is $\leq N$.

$$s(Q) = 2$$
 \Rightarrow Factorization size is $O(N^2)$

Recall that $\rho^*(Q) = 3 \implies$ Listing representation size is $O(N^3)$

Example: The Fractional Hypertree Width *fhtw*

Idea: Avoid repeating identical expressions, store them once and use pointers.



$$\bigcup_{\mathbf{a}\in\mathbf{A}} \left[\langle \mathbf{a} \rangle \times \cdots \times \bigcup_{\mathbf{e}\in\mathbf{E}} \left(\langle \mathbf{e} \rangle \times \left(\bigcup_{f\in F} \langle f \rangle \right) \right) \right]$$

Observation:

- Variable *F* only depends on **E** and not on **A**: key(*F*) = {**E**}
- A value (e) maps to the same union U_{(e,f)∈U}(f) regardless of its pairings with A-values.

 $\Rightarrow \text{ Define } U_e = \bigcup_{(e,f) \in U} \langle f \rangle \text{ for each value } \langle e \rangle \text{ and use } U_e \text{ instead of the union } \bigcup_{(e,f) \in U} \langle f \rangle.$

Example: The Fractional Hypertree Width *fhtw*

Idea: Avoid repeating identical expressions, store them once and use pointers.



A factorization with caching would be:

$$\bigcup_{a \in \mathbf{A}} \left[\langle a \rangle \times \cdots \times \bigcup_{e \in \mathbf{E}} \left(\langle e \rangle \times U_e \right) \right]; \qquad \left\{ U_e = \bigcup_{(e,f) \in U} \langle f \rangle \right\}$$

- *fhtw* for Δ is the largest ρ^{*}(Q_{key(X)∪{X}}) over subqueries Q_{key(X)∪{X}} defined by the variables key(X)∪{X} for each variable X in Δ
- fhtw(Q) is the minimum fhtw over all variable orders of Q

In our example: $fhtw(Q) = 1 < s(Q) = 2 < \rho^*(Q) = 3$.

Alternative Characterizations of *fhtw*

The fractional hypertree width *fhtw* has been originally defined for hypertree decompositions. [M10]

- Given a join query Q.
- Let \mathbf{T} be the set of hypertree decompositions of the hypergraph of Q.

 $\mathit{fhtw}(Q) = \min_{(\mathcal{T},\chi) \in \mathbf{T}} \max_{n \in \mathcal{T}} \rho^*(Q_{\chi(n)})$

Alternative Characterizations of *fhtw*

The fractional hypertree width *fhtw* has been originally defined for hypertree decompositions. [M10]

- Given a join query Q.
- Let \mathbf{T} be the set of hypertree decompositions of the hypergraph of Q.

$$\mathit{fhtw}(\mathit{Q}) = \min_{(\mathit{T},\chi)\in \mathbf{T}} \max_{\mathit{n}\in \mathit{T}} \rho^*(\mathit{Q}_{\chi(\mathit{n})})$$

Alternative characterization of the fractional hypertree width *fhtw* using the mapping between hypertree decompositions and variable orders [OZ15]

- Given a join query Q.
- Let **VO** be the set of variable orders of *Q*.

 $\mathit{fhtw}(Q) = \min_{(F,\mathit{key}) \in \mathbf{VO}} \max_{v \in F} \rho^*(Q_{\mathit{key}(v) \cup \{v\}})$

Relational Counterpart of Factorized Representation

Covers: Relational Counterparts of Factorizations

Factorized representations are not relational :(

This makes it difficult to integrate them into relational data systems

- Covers of Query Results
 - Relations that are lossless representations of query results, yet are as succinct as factorized representations
 - For a join query Q and any database of size N, a cover has size $O(N^{fhtw(Q)})$ and can be computed in time $\tilde{O}(N^{fhtw(Q)})$

How to get a cover?

- Construct a hypertree decomposition of the query
- Project query result onto the bags of the hypertree decomposition
- Construct on these projections the hypergraph of the query result
- Take a minimal edge cover of this hypergraph

[KO17]

Recall the Itemized Customer Orders Example

Orders (O for short)			Dish (D	Dish (D for short)		Items (I for short)	
customer	day	dish	dish	item	item	price	
Elise	Monday	burger	burger	patty	patty	6	
Elise	Friday	burger	burger	onion	onion	2	
Steve	Friday	hotdog	burger	bun	bun	2	
Joe	Friday	hotdog	hotdog	bun	sausage	4	
			hotdog	onion			
			hotdog	sausage			



O(customer, day, disit), D(disit, item), i(item, price	O(customer,	day,	dish)	, D(dish,	item).	l(item,	price
--	-------------	------	-------	------	-------	--------	---------	-------

			. , (
customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2

Elise Monday burger

Elise Friday burger

<	customer,day,dish
	dish,item
	item,price

O(customer, day, dish), D(dish, item), I(item, price)

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2

burger patty

Elise Monday burger

burger onion

Elise Friday burger

burger bun

<	customer,day,dish	,
	dish,item	
	item,price	

O(customer, day, dish), D(dish, item), I(item, price)

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2

					burger	patty			patty	6
	Elise	Monday	burger							
					burger	onion			onion	2
	Elise	Friday	burger							
					burger	bun			bun	2
\subset	custom	er,day,dish	$\overline{}$				5/11		1/1-	
		- <u>,</u>		O(custo	mer, day,	dish),	D(dish	, item),	I(item,	price)
	dis	hitem		customer	d	ay	dish	item		price
				Elise	Mond	ay b	ourger	patty	(6
	iter	n,price		Elise	Mond	ay b	ourger	onion		2
		·		Elise	Mond	ay b	ourger	bun		2
				Elise	Frid	ay b	ourger	patty	,	6
				Elise	Frid	ay b	ourger	onion		2
				Elise	Frid	ay b	ourger	bun		2









O(customer, day,	dish), D(dish	, item), I	(item, price)
------------------	---------------	------------	---------------

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2





O(Customer, uay, uisir), D(uisir, item), i(item, pi	price	- 1
---	-------	-----

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2





O(customer, day, dish), D(dish, item), I(item,	I(item, price)
--	----------------

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2
The Hypergraph of the Query Result





O(customer, da	/, dish),	D(dish,	item), I	(item,	price)
----------------	-----------	---------	----------	--------	--------

customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Monday	burger	onion	2
Elise	Monday	burger	bun	2
Elise	Friday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2

A Minimal Edge Cover of the Hypergraph



A Cover of (a part of) the Query Result

O(custon	ner, day, dis	h), D(dish	, item),	I(item, price)
customer	day	dish	item	price
Elise	Monday	burger	patty	6
Elise	Friday	burger	onion	2
Elise	Friday	burger	bun	2



O(customer, day, dish), D(dish, item), I(item, price)					
customer	day	dish	item	price	
Elise	Monday	burger	patty	6	
Elise	Monday	burger	onion	2	
Elise	Monday	burger	bun	2	
Elise	Friday	burger	patty	6	
Elise	Friday	burger	onion	2	
Elise	Friday	burger	bun	2	

Compression by Factorization in Practice

Compression Contest: Factorized vs. Zipped Relations



Result of query $\operatorname{Orders} \bowtie \operatorname{Dish} \bowtie \operatorname{Items}$



- Tabular = listing representation in CSV text format
- Gzip (compression level 6) outputs binary format
- Factorized representation in text format (each digit takes one character)

Observations:

- Gzip does not exploit distant repetitions!
- **Factorizations** can be arbitrarily more succinct than gzipped relations.
- **Gzipping factorizations improves the compression by 3x.**

Factorization Gains in Practice (1/3)

Retailer dataset used for LogicBlox analytics

- Relations: Inventory (84M), Sales (1.5M), Clearance (368K), Promotions (183K), Census (1K), Location (1K).
- Compression factors (caching not used):
 - > 26.61x for natural join of Inventory, Census, Location.
 - ▶ 159.59x for natural join of Inventory, Sales, Clearance, Promotions

Factorization Gains in Practice (2/3)

LastFM public dataset

- Relations: UserArtists (93K), UserFriends (25K), TaggedArtists (186K).
- Compression factors:
 - ▶ 143.54x for joining two copies of Userartists and Userfriends

With caching: 982.86x

- 253.34x when also joining on TaggedArtists
- > 2.53x/ 3.04x/ 924.46x for triangle/4-clique/bowtie query on UserFriends
- ▶ 9213.51x/ 552Kx/ ≥86Mx for versions of triangle/4-clique/bowtie queries with copies for UserArtists for each UserFriend copy

Factorization Gains in Practice (3/3)

Twitter public dataset

- Relation: Follower-Followee (1M)
- Compression factors:
 - 2.69x for triangle query
 - ▶ 3.48x for 4-clique query
 - 4918.73x for bowtie query

Worst-Case Optimal Join Algorithms

How Fast Can We Compute Join Results?

Given a join query Q, for any database of size N, the join result can be computed in time

- $\tilde{O}(N^{\rho^*(Q)})$ as listing representation [NPRR12,V14] • $\tilde{O}(N^{s(Q)})$ as factorization without caching [OZ15]
- $\widetilde{O}(N^{fhtw(Q)})$ as factorization with caching [OZ15]

These upper bounds essentially follow the succinctness gap. They are:

- worst-case optimal (modulo log *N*) within the given representation model
- with respect to data complexity
 - \blacktriangleright additional quadratic factor in the number of variables and linear factor in the number of relations in Q

Example: Computing the Factorized Join Result with FDB

Our join: O(customer, day, dish), D(dish, item), I(item, price) can be grounded to a factorized representation as follows:



This computation follows the variable order given below:



Example: Computing the Factorized Join Result with FDB



Relations are sorted following any topological order of the variable order

- The intersection of relations O and D on *dish* takes time $O(N_{\min} \log(N_{\max}/N_{\min})) = \widetilde{O}(N_m in)$, where $N_m = m(|\pi_{dish}O|, |\pi_{dish}D|)$.
- The remaining operations are lookups in the relations, where we first fix the *dish* value and then the *day* and *item* values.

LeapFrog TrieJoin Algorithm

- Much acclaimed worst-case optimal join algorithm used by LogicBlox [V14]
- Computes a listing representation of the join result
 - \Rightarrow It does not exploit factorization
- Glorified multi-way sort-merge join with an efficient list intersection
- Several generalizations, e.g., PANDA

[NRR13,ANS17]

LeapFrog TrieJoin is a special case of FDB, where

- \blacksquare the input variable order Δ is a path, and
- for each variable A, key(A) consists of all ancestors of A in Δ .

Experiment: Factorized vs. Listing Computation

		Retailer (3B)	LastFM (5.8M)
Join	Factorization	169M	316K
Size	Listing	3.6B	591M
(values)	Compression	21.4×	1870.7×
Join	FDB	30	10
Time	PostgreSQL	217	61
(sec)	Speedup	7×	6.1×



Both FDB and PostgreSQL list the records in the results of the join queries.

Outline



Part 1. Joins

Part 2. Aggregates

Part 3. Optimization

Aggregates

Important operators in database query languages and essential for applications.

Natural generalization of aggregates over joins can express a host of problems across Computer Science. [ANR16]

We highlight recent work on aggregate computation with lowest known computational complexity. This extends the work from Part 1.

[BKOZ13,ANR16]

Part 3 later discusses an extension of this work to state-of-the-art machine learning inside the database. [SOC16,ANNOS17]

Plan for Part 2 on Aggregates

 Computation of aggregates over factorized joins using the FDB algorithm [BKOZ13]

Factorized computation of aggregates using optimized relational queries.
[SOC16,OS16]

Functional Aggregate Queries (FAQs) [ANR16]

- Generalize aggregate-join queries to many semirings, e.g., sum-product, max-product, Boolean
- FAQ computation is factorized and has the computational complexity of aggregates over factorized joins
- FAQ computation using the InsideOut algorithm [ANR16]

Examples: Aggregates over Factorized Joins

Example 1: COUNT Aggregate over Factorized Join



- COUNT(*):
 - ▶ values \mapsto 1,
 - $\blacktriangleright \cup \mapsto +,$
 - $\blacktriangleright \times \mapsto *.$

Example 1: COUNT Aggregate over Factorized Join



- COUNT(*):
 - ▶ values \mapsto 1,
 - $\blacktriangleright \cup \mapsto +,$
 - $\blacktriangleright \times \mapsto *.$

Example 2: SumProd Aggregate over Factorized Join



- SUM(dish * price):
 - ▶ Assume there is a function *f* that turns dish into reals or indicator vectors.
 - ▶ All values except for dish & price \mapsto 1,
 - $\blacktriangleright \cup \mapsto +,$
 - $\blacktriangleright \times \mapsto *.$

Example 2: SumProd Aggregate over Factorized Join



- SUM(dish * price):
 - Assume there is a function f that turns dish into reals.
 - ▶ All values except for dish & price \mapsto 1,
 - \blacktriangleright \cup \mapsto +,
 - $\blacktriangleright \times \mapsto *.$

Computing Aggregates over Factorized Joins using FDB

Given an aggregate-join query ${\it Q}$

[BKOZ13]

- Construct a variable order Δ where the group-by (free) variables are above the other (bound) variables of Q
 - A new width w measure that is at least *fhtw*
- Compute the factorized join over Δ
 - The complexity now depends on the width w
- Finally compute the aggregates in one pass over the factorized join.

Computing Aggregates over Factorized Joins using FDB

Given an aggregate-join query Q

[BKOZ13]

- Construct a variable order Δ where the group-by (free) variables are above the other (bound) variables of Q
 - A new width w measure that is at least *fhtw*
- Compute the factorized join over Δ
 - The complexity now depends on the width w
- Finally compute the aggregates in one pass over the factorized join.

Is it necessary to first compute the factorized join?

Computing Aggregates over Factorized Joins using FDB

Given an aggregate-join query ${\it Q}$

[BKOZ13]

- Construct a variable order Δ where the group-by (free) variables are above the other (bound) variables of Q
 - A new width w measure that is at least *fhtw*
- Compute the factorized join over Δ
 - The complexity now depends on the width w
- Finally compute the aggregates in one pass over the factorized join.

Is it necessary to first compute the factorized join?

Aggregates can be computed without materializing the factorized join

- The factorized join becomes the *trace* of the aggregate computation
- This is the *factorized computation* of the query *Q*.

Example: Factorized Aggregate Computation

The 4-path query Q_4 on a graph with the edge relation $E(E_i)$'s are copies of E:

 $V_1(A), E_1(A, B), E_2(B, C), E_3(C, D), E_4(D, E), V_2(E)$



Example: Factorized Aggregate Computation

The 4-path query Q_4 on a graph with the edge relation $E(E_i)$'s are copies of E:

 $V_1(A), E_1(A, B), E_2(B, C), E_3(C, D), E_4(D, E), V_2(E)$



Recall sizes for factorized results of path queries

- $\rho^*(Q_4) = 3 \Rightarrow$ listing representation has size $O(|E|^3)$.
- $fhtw(Q_4) = 1 \Rightarrow$ factorization with caching has size O(|E|).

Example: Factorized Aggregate Computation

We would like to compute $COUNT(Q_4)$:

- in O(|E|) time (no free variables, so use best variable order for Q_4)
- using optimized queries that are derived from the variable order of Q_4
- without materializing the factorized result of the path query

Convention:

- View the relations as functions mapping tuples to numbers.
- The functions for input relations map their tuples to 1.















This computation strategy corresponds to the following query rewriting:

$$\begin{split} &\sum_{a \in \text{Dom}(A)} \sum_{b \in \text{Dom}(B)} \sum_{c \in \text{Dom}(C)} \sum_{d \in \text{Dom}(D)} \sum_{e \in \text{Dom}(E)} V_1(a) \cdot E_1(b, a) \cdot E_2(c, b) \cdot E_3(c, d) \cdot E_4(d, e) \cdot V_2(e) \\ = \\ &\sum_{c \in \text{Dom}(C)} \left(\sum_{b \in \text{Dom}(B)} E_2(c, b) \cdot \left(\sum_{a \in \text{Dom}(A)} V_1(a) \cdot E_1(b, a) \right) \right) \cdot \\ &\left(\sum_{d \in \text{Dom}(D)} E_3(c, d) \cdot \left(\sum_{e \in \text{Dom}(E)} E_4(d, e) \cdot V_2(e) \right) \right) \end{split}$$

Is Factorized Aggregate Computation Practical?

Experiments published in several papers, here a quick glimpse from [ANNOS17]

Retailer dataset (records)	excerpt (17M)	full (86M)
PostgreSQL computing the join	50.63 sec	216.56 sec
FDB computing both the join and the aggregates Number of aggregates (scalar+group-by)	25.51 sec 595+2,418	<mark>380.31 sec</mark> 595+145k
FDB computing both the join and the aggregates	132.43 sec	1,819.80 sec
Number of aggregates (scalar+group-by)	158k+742k	158k+37M

In this experiment:

- FDB only used one core of a commodity machine
- For both PostgreSQL and FDB, the dataset was entirely in memory
- The aggregates represent gradients (or parts thereof) used for learning degree 1 and 2 polynomial regression models
Functional Aggregate Queries

Functional Aggregate Query

FAQ generalizes factorized aggregate computation to a host of problems.

We use the following notation $(i \in [n] = \{1, \dots, n\})$:

X_i are variables,

• x_i are values in discrete domain $Dom(X_i)$

- $\mathbf{x} = (x_1, \ldots, x_n) \in \mathsf{Dom}(X_1) \times \cdots \times \mathsf{Dom}(X_n)$
- For any $S \subseteq [n]$,

$$\mathbf{x}_{5} = (x_{i})_{i \in S} \in \prod_{i \in S} \text{Dom}(X_{i})$$

e.g. $\mathbf{x}_{\{2,5,8\}} = (x_{2}, x_{5}, x_{8}) \in \text{Dom}(X_{2}) \times \text{Dom}(X_{5}) \times \text{Dom}(X_{8})$

Functional Aggregate Query: The Problem





All functions have the same range D

- n variables X₁,..., X_n
- a multi-hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$
 - Each vertex is a variable (notation overload: $\mathcal{V} = [n]$)
 - ▶ To each hyperedge $S \in \mathcal{E}$ there corresponds a *factor* ψ_S

$$\psi_{\mathcal{S}}:\prod_{i\in\mathcal{S}}\mathsf{Dom}(X_i)\to \mathbf{D}$$



- n variables X₁,...,X_n
- a multi-hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$
 - Each vertex is a variable (notation overload: $\mathcal{V} = [n]$)
 - ▶ To each hyperedge $S \in \mathcal{E}$ there corresponds a *factor* ψ_S

$$\psi_{\mathcal{S}}:\prod_{i\in\mathcal{S}}\mathsf{Dom}(X_i)\to \mathbf{D}$$



- n variables X₁,...,X_n
- a multi-hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$
 - Each vertex is a variable (notation overload: $\mathcal{V} = [n]$)
 - ▶ To each hyperedge $S \in \mathcal{E}$ there corresponds a *factor* ψ_S

$$\psi_{\mathcal{S}}:\prod_{i\in\mathcal{S}}\mathsf{Dom}(X_i)\to \mathbf{D}$$



- *n* variables X_1, \ldots, X_n
- a multi-hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$
 - Each vertex is a variable (notation overload: $\mathcal{V} = [n]$)
 - ▶ To each hyperedge $S \in \mathcal{E}$ there corresponds a *factor* ψ_S

$$\psi_{\mathcal{S}}:\prod_{i\in\mathcal{S}}\mathsf{Dom}(X_i)\to \mathbf{D}$$



- n variables X₁,...,X_n
- a multi-hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$
 - Each vertex is a variable (notation overload: $\mathcal{V} = [n]$)
 - ▶ To each hyperedge $S \in \mathcal{E}$ there corresponds a *factor* ψ_S

$$\psi_{S}: \prod_{i \in S} \mathsf{Dom}(X_{i}) o \underbrace{\mathsf{D}}_{\uparrow}$$
 $\mathbf{R}_{+}, \{\mathsf{true}, \mathsf{false}\}, \{0, 1\}, 2^{\mathcal{U}}, \mathsf{etc.}$



All functions have the same range ${\bf D}$

- n variables X₁,...,X_n
- a multi-hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$
 - Each vertex is a variable (notation overload: $\mathcal{V} = [n]$)
 - ▶ To each hyperedge $S \in \mathcal{E}$ there corresponds a *factor* ψ_S

$$\psi_{S} : \prod_{i \in S} \mathsf{Dom}(X_{i}) \to \bigwedge_{i \in S} \mathsf{R}_{+}, \{\mathsf{true}, \mathsf{false}\}, \{0, 1\}, 2^{\mathcal{U}}, \mathsf{etc}$$

• a set $F \subseteq \mathcal{V}$ of free variables (wlog, $F = [f] = \{1, \dots, f\}$)



All functions have the same range D

• Compute the function $\varphi : \prod_{i \in F} \text{Dom}(X_i) \to \mathbf{D}$.



- Compute the function $\varphi : \prod_{i \in F} \text{Dom}(X_i) \to \mathbf{D}$.
- φ defined by the FAQ-expression

$$\varphi(\mathbf{x}_{[f]}) = \bigoplus_{x_{f+1} \in \mathsf{Dom}(X_{f+1})}^{(f+1)} \cdots \bigoplus_{x_{n-1} \in \mathsf{Dom}(X_{n-1})}^{(n-1)} \bigoplus_{x_n \in \mathsf{Dom}(X_n)}^{(n)} \bigotimes_{S \in \mathcal{E}} \psi_S(\mathbf{x}_S)$$



- Compute the function $\varphi : \prod_{i \in F} \text{Dom}(X_i) \to \mathbf{D}$.
- φ defined by the FAQ-expression

$$\varphi(\mathbf{x}_{[f]}) = \bigoplus_{x_{f+1} \in \text{Dom}(X_{f+1})}^{(f+1)} \cdots \bigoplus_{x_{n-1} \in \text{Dom}(X_{n-1})}^{(n-1)} \bigoplus_{x_n \in \text{Dom}(X_n)}^{(n)} \bigotimes_{S \in \mathcal{E}} \psi_S(\mathbf{x}_S)$$

• For each $\bigoplus^{(i)}$



- Compute the function $\varphi : \prod_{i \in F} \text{Dom}(X_i) \to \mathbf{D}$.
- φ defined by the FAQ-expression

$$\varphi(\mathbf{x}_{[f]}) = \bigoplus_{x_{f+1} \in \text{Dom}(X_{f+1})}^{(f+1)} \cdots \bigoplus_{x_{n-1} \in \text{Dom}(X_{n-1})}^{(n-1)} \bigoplus_{x_n \in \text{Dom}(X_n)}^{(n)} \bigotimes_{S \in \mathcal{E}} \psi_S(\mathbf{x}_S)$$

■ For each ⊕⁽ⁱ⁾
 ▶ Either (D, ⊕⁽ⁱ⁾, ⊗) is a commutative semiring



- Compute the function $\varphi : \prod_{i \in F} \text{Dom}(X_i) \to \mathbf{D}$.
- φ defined by the FAQ-expression

$$\varphi(\mathbf{x}_{[f]}) = \bigoplus_{x_{f+1} \in \text{Dom}(X_{f+1})}^{(f+1)} \cdots \bigoplus_{x_{n-1} \in \text{Dom}(X_{n-1})}^{(n-1)} \bigoplus_{x_n \in \text{Dom}(X_n)}^{(n)} \bigotimes_{S \in \mathcal{E}} \psi_S(\mathbf{x}_S)$$

Semirings

($\mathbf{D}, \oplus, \otimes$ **) is a** commutative semiring when

Additive identity $0 \in D : 0 \oplus e = e \oplus 0 = e$ Multiplicative identity $1 \in D : 1 \otimes e = e \otimes 1 = e$

Annihilation by **0**

$$\otimes e = e \otimes \mathbf{0} = \mathbf{0}$$

Distributive law $a \otimes b \oplus a \otimes c = a \otimes (b \oplus c)$

Semirings

($\mathbf{D}, \oplus, \otimes$ **)** is a commutative semiring when

Additive identity $0 \in D : 0 \oplus e = e \oplus 0 = e$ Multiplicative identity $1 \in D : 1 \otimes e = e \otimes 1 = e$ Annihilation by 0 $0 \otimes e = e \otimes 0 = 0$ Distributive law $a \otimes b \oplus a \otimes c = a \otimes (b \oplus c)$

Common examples (there are many more!)

Boolean ({true, false}, \lor , \land) sum-product (\mathbb{R} , +, \times) max-product (\mathbb{R} _+, max, \times) set (2^U, \cup , \cap)

$\mathsf{SumProduct} \subset \mathsf{FAQ}$

Problem (SumProduct)

Given a commutative semiring $(\mathbf{D},\oplus,\otimes)$, compute the function

$$\varphi(x_1,\ldots,x_f) = \bigoplus_{x_{f+1}} \bigoplus_{x_{f+2}} \cdots \bigoplus_{x_n} \bigotimes_{S \in \mathcal{E}} \psi_S(\mathbf{x}_S)$$

$\mathsf{SumProduct} \subset \mathsf{FAQ}$

Problem (SumProduct)

Given a commutative semiring $(\mathbf{D},\oplus,\otimes)$, compute the function

$$\varphi(\mathbf{x}_1,\ldots,\mathbf{x}_f) = \bigoplus_{\mathbf{x}_{f+1}} \bigoplus_{\mathbf{x}_{f+2}} \cdots \bigoplus_{\mathbf{x}_n} \bigotimes_{S \in \mathcal{E}} \psi_S(\mathbf{x}_S)$$

SumProduct

- Rina Dechter (Artificial Intelligence 1999 and earlier)
- $\blacksquare \equiv$ Marginalize a Product Function
 - ▶ Aji and McEliece (IEEE Trans. Inform. Theory 2000)

Many examples for SumProduct

- ({true, false}, \lor , \land)
 - Constraint satisfaction problems
 - Boolean conjunctive query evaluation
 - SAT
 - k-colorability
 - etc.
- (*U*, ∪, ∩)
 - Conjunctive query evaluation
- (ℝ, +, ×)
 - Permanent
 - DFT
 - Inference in probabilistic graphical models
 - ► #CSP
 - Matrix chain multiplication
 - Aggregates in DB
- **(R** $_+, \max, \times$)
 - MAP queries in probabilistic graphical models

SumProduct Example 1: Boolean Query Evaluation

Boolean Conjunctive Queries:

- Boolean query Φ with set $rels(\Phi)$ of relation symbols
- Each relation symbol $R \in rels(\Phi)$ has variables vars(R)

$$\Phi = \exists X_1 \dots \exists X_n : \bigwedge_{R \in rels(\Phi)} R(vars(R))$$

FAQ encoding:

$$\phi = \bigvee_{\mathbf{x}} \bigwedge_{S \in \mathcal{E}} \psi_{S}(\mathbf{x}_{S}), \text{ where }$$

- Φ has the hypergraph $(\mathcal{V}, \mathcal{E})$ with
- $\mathcal{V} = \bigcup_{R \in rels(\Phi)} vars(R)$ and $\mathcal{E} = \{vars(R) \mid R \in rels(\Phi)\}$
- For each $S \in \mathcal{E}$, there is a factor ψ_S such that $\psi_S(\mathbf{x}_S) = (\mathbf{x}_S \in R)$

SumProduct Example 2: Matrix Chain Multiplication

Compute the product $\mathbf{A} = \mathbf{A}_1 \cdots \mathbf{A}_n$ of *n* matrices

Each matrix \mathbf{A}_i is over field \mathbb{F} and has dimensions $p_i \times p_{i+1}$

FAQ encoding:

- We use n + 1 variables X_1, \ldots, X_{n+1} with domains $Dom(X_i) = [p_i]$
- Each matrix **A**_i can be viewed as a function of two variables:

$$\psi_{i,i+1}$$
: Dom (X_i) × Dom (X_{i+1}) \rightarrow \mathbb{F} , where $\psi_{i,i+1}(x, y) = (\mathbf{A}_i)_{xy}$

The problem is now to compute the FAQ expression

$$\phi(\mathbf{x}_1, \mathbf{x}_{n+1}) = \sum_{\mathbf{x}_2 \in \mathsf{Dom}(\mathbf{X}_2)} \cdots \sum_{\mathbf{x}_n \in \mathsf{Dom}(\mathbf{X}_n)} \prod_{i \in [n]} \psi_{i,i+1}(\mathbf{x}_i, \mathbf{x}_{i+1}).$$

SumProduct Example 3: Queries in Graphical Models

- Discrete undirected graphical model represented by a hypergraph $(\mathcal{V}, \mathcal{E})$
- $\mathcal{V} = \{X_1, \ldots, X_n\}$ consists of *n* discrete random variables
- There is a factor $\psi_S : \prod_{i \in S} \mathsf{Dom}(X_i) \to \mathbb{R}_+$ for each edge $S \in \mathcal{E}$

FAQ expression to compute the marginal Maximum A Posteriori estimates:

$$\phi(x_1,\ldots,x_f) = \max_{x_{f+1} \in \text{Dom}(X_{f+1})} \cdots \max_{x_n \in \text{Dom}(X_n)} \prod_{S \in \mathcal{S}} \psi_S(\mathbf{x}_S)$$

FAQ expression to compute the marginal distribution of variables X_1, \ldots, X_f :

$$\phi(x_1,\ldots,x_f) = \sum_{x_{f+1} \in \text{Dom}(X_{f+1})} \cdots \sum_{x_n \in \text{Dom}(X_n)} \prod_{S \in \mathcal{E}} \psi_S(\mathbf{x}_S)$$

For conditional distributions $p(\mathbf{x}_A \mid \mathbf{x}_B)$, the variables \mathbf{X}_B are set to values \mathbf{x}_B .

 $\varphi(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4) = \sum \psi_1(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \cdot \psi_2(\mathbf{x}_2, \mathbf{x}_4, \mathbf{x}_5) \cdot \psi_3(\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6) \cdot \psi_4(\mathbf{x}_6, \mathbf{x}_8) \cdot \psi_5(\mathbf{x}_5, \mathbf{x}_7)$ x_3, x_5, x_6, x_7, x_8 X_2 $key(X_2) = \emptyset$ $key(X_1) = \{X_2\}$ ψ_2 ′X₄ X_4 ψ_1 X_1 X_1 $key(X_3) = \{X_1, X_2\}$ $key(X_4) = \{X_2\}$ X_5 X_5 ψ_3 X_3 $key(X_5) = \{X_2, X_4\}$ $key(X_6) = \{X_4, X_5\}$ X_6 X_7 X_6 $key(X_8) = \{X_6\}$ ψ_5 $key(X_7) = \{X_5\}$ ψ_4 X_8 X_8



• $\rho^*(\varphi) = 4$, $s(\varphi) = 2$, $fhtw(\varphi) = 1$. The above variable order Δ has the free variables x_1, x_2, x_4 on top of the others and $fhtw(\Delta) = 1$.

• The query result has size: O(N) when factorized; $O(N^2)$ when listed

$$\begin{array}{c} X_2 \\ X_1 \\ X_3 \\ X_7 \\ X_7 \\ X_6 \\ X_8 \end{array}$$

$$\varphi(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4) = \sum_{\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_6, \mathbf{x}_7, \mathbf{x}_8} \psi_1(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \cdot \psi_2(\mathbf{x}_2, \mathbf{x}_4, \mathbf{x}_5) \cdot \psi_3(\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6) \cdot \psi_4(\mathbf{x}_6, \mathbf{x}_8) \cdot \psi_5(\mathbf{x}_5, \mathbf{x}_7)$$

$$\begin{split} \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{3}, x_{5}, x_{6}, x_{7}, x_{8}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}, x_{8}} \left(\underbrace{\sum_{\mathbf{x}_{3}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3})}_{\psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2})} \right) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \end{split}$$

$$\begin{split} \varphi(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4) &= \sum_{\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_6, \mathbf{x}_7, \mathbf{x}_8} \psi_1(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \cdot \psi_2(\mathbf{x}_2, \mathbf{x}_4, \mathbf{x}_5) \cdot \psi_3(\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6) \cdot \psi_4(\mathbf{x}_6, \mathbf{x}_8) \cdot \psi_5(\mathbf{x}_5, \mathbf{x}_7) \\ \varphi(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4) &= \sum_{\mathbf{x}_5, \mathbf{x}_6, \mathbf{x}_7, \mathbf{x}_8} \psi_6(\mathbf{x}_1, \mathbf{x}_2) \cdot \psi_2(\mathbf{x}_2, \mathbf{x}_4, \mathbf{x}_5) \cdot \psi_3(\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6) \cdot \psi_4(\mathbf{x}_6, \mathbf{x}_8) \cdot \psi_5(\mathbf{x}_5, \mathbf{x}_7) \quad \widetilde{O}(N) \end{split}$$

$$\begin{split} \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{3}, x_{5}, x_{6}, x_{7}, x_{8}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}, x_{8}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \left(\sum_{x_{8}} \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \right) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \psi_{7}(\mathbf{x}_{6}) \end{split}$$



$$\begin{split} \varphi(x_1, x_2, x_4) &= \sum_{x_3, x_5, x_6, x_7, x_8} \psi_1(x_1, x_2, x_3) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_3(x_4, x_5, x_6) \cdot \psi_4(x_6, x_8) \cdot \psi_5(x_5, x_7) \\ \varphi(x_1, x_2, x_4) &= \sum_{x_5, x_6, x_7, x_8} \psi_6(x_1, x_2) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_3(x_4, x_5, x_6) \cdot \psi_4(x_6, x_8) \cdot \psi_5(x_5, x_7) \quad \widetilde{O}(N) \\ \varphi(x_1, x_2, x_4) &= \sum_{x_5, x_6, x_7} \psi_6(x_1, x_2) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_3(x_4, x_5, x_6) \cdot \psi_7(x_6) \cdot \psi_5(x_5, x_7) \quad \widetilde{O}(N) \end{split}$$



$$\begin{split} \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{3}, x_{5}, x_{6}, x_{7}, x_{8}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}, x_{8}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \left(\sum_{x_{7}} \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7})\right) \\ \underbrace{\psi_{8}(\mathbf{x}_{5})} \psi_{8}(\mathbf{x}_{5}) \cdot \psi_{8}(\mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \left(\sum_{x_{7}} \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}\right)\right) \\ \psi_{8}(\mathbf{x}_{5}) \cdot \psi_{8}(\mathbf{x}_{5}) \cdot \psi_{8}(\mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{8}(\mathbf{x}_{6}, \mathbf{x}_{6}) \cdot \psi_{8}(\mathbf{x}_{6}, \mathbf{x}_{6}, \mathbf{x}_{6}, \mathbf{x}_{7}, \mathbf{x}_{8}) \cdot \psi_{8}(\mathbf{x}_{8}, \mathbf{x}_{8}) \cdot \psi_{8}(\mathbf{x}_{$$

$$\begin{split} \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{3}, x_{5}, x_{6}, x_{7}, x_{8}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}, x_{8}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{\mathbf{x}_{5}, x_{6}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{8}(\mathbf{x}_{5}) \quad \widetilde{O}(N) \end{split}$$

$$\begin{split} \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{\substack{x_{3}, x_{5}, x_{6}, x_{7}, x_{8}}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{\substack{x_{5}, x_{6}, x_{7}, x_{8}}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{\substack{x_{5}, x_{6}, x_{7}}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{\substack{x_{5}, x_{6}}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{8}(\mathbf{x}_{5}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{\substack{x_{5}, x_{6}}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot (\sum_{\substack{x_{6}, x_{7}, x_{9}, x_$$

$$\begin{split} \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{3}, x_{5}, x_{6}, x_{7}, x_{8}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}, x_{8}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \tilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \tilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{8}(\mathbf{x}_{5}) \quad \tilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{9}(\mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{8}(\mathbf{x}_{5}) \quad \tilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{9}(\mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{8}(\mathbf{x}_{5}) \quad \tilde{O}(N) \end{split}$$



$$\begin{split} \varphi(x_1, x_2, x_4) &= \sum_{\substack{x_3, x_5, x_6, x_7, x_8}} \psi_1(x_1, x_2, x_3) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_3(x_4, x_5, x_6) \cdot \psi_4(x_6, x_8) \cdot \psi_5(x_5, x_7) \\ \varphi(x_1, x_2, x_4) &= \sum_{\substack{x_5, x_6, x_7, x_8}} \psi_6(x_1, x_2) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_3(x_4, x_5, x_6) \cdot \psi_4(x_6, x_8) \cdot \psi_5(x_5, x_7) \quad \widetilde{O}(N) \\ \varphi(x_1, x_2, x_4) &= \sum_{\substack{x_5, x_6, x_7}} \psi_6(x_1, x_2) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_3(x_4, x_5, x_6) \cdot \psi_7(x_6) \cdot \psi_5(x_5, x_7) \quad \widetilde{O}(N) \\ \varphi(x_1, x_2, x_4) &= \sum_{\substack{x_5, x_6}} \psi_6(x_1, x_2) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_3(x_4, x_5, x_6) \cdot \psi_7(x_6) \cdot \psi_8(x_5) \quad \widetilde{O}(N) \\ \varphi(x_1, x_2, x_4) &= \sum_{\substack{x_5, x_6}} \psi_6(x_1, x_2) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_9(x_4, x_5) \cdot \psi_8(x_5) \quad \widetilde{O}(N) \\ \varphi(x_1, x_2, x_4) &= \sum_{\substack{x_5, x_6}} \psi_6(x_1, x_2) \cdot \psi_2(x_2, x_4, x_5) \cdot \psi_9(x_4, x_5) \cdot \psi_8(x_5) \quad \widetilde{O}(N) \\ \varphi(x_1, x_2, x_4) &= \psi_6(x_1, x_2) \cdot \left(\sum_{\substack{x_5}} \psi_2(x_2, x_4, x_5) \cdot \psi_9(x_4, x_5) \cdot \psi_8(x_5) \right) \\ \psi_{10}(x_2, x_4) &= \psi_{10}(x_2, x_4) \\ \varphi(x_1, x_2, x_4) &= \psi_{10}(x_1, x_2) \cdot \left(\sum_{\substack{x_5}} \psi_2(x_2, x_4, x_5) \cdot \psi_9(x_4, x_5) \cdot \psi_8(x_5) \right) \\ \psi_{10}(x_2, x_4) &= \psi_{10}(x_1, x_2) \cdot \left(\sum_{\substack{x_5}} \psi_2(x_2, x_4, x_5) \cdot \psi_9(x_4, x_5) \cdot \psi_8(x_5) \right) \\ \psi_{10}(x_2, x_4) &= \psi_{10}(x_2, x_4) \\ \varphi(x_1, x_2, x_4) &= \psi_{10}(x_1, x_2) \cdot \left(\sum_{\substack{x_5}} \psi_2(x_2, x_4, x_5) \cdot \psi_9(x_4, x_5) \cdot \psi_8(x_5) \right) \\ \psi_{10}(x_2, x_4) &= \psi_{10}(x_2, x_4) \\ \varphi(x_1, x_2, x_4) &= \psi_{10}(x_1, x_2) \cdot \left(\sum_{\substack{x_5}} \psi_2(x_2, x_4, x_5) \cdot \psi_9(x_4, x_5) \cdot \psi_8(x_5) \right) \\ \psi_{10}(x_2, x_4) &= \psi_{10}(x_2, x_4) \\ \psi_{10}(x_2, x_4) \\ \varphi(x_1, x_2, x_4) &= \psi_{10}(x_2, x_4) \\ \psi_{10}(x_2$$

115 / 171

$$\begin{split} \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{3}, x_{5}, x_{6}, x_{7}, x_{8}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}, x_{8}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{4}(\mathbf{x}_{6}, \mathbf{x}_{8}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}, x_{7}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{5}(\mathbf{x}_{5}, \mathbf{x}_{7}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}, x_{6}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{3}(\mathbf{x}_{4}, \mathbf{x}_{5}, \mathbf{x}_{6}) \cdot \psi_{7}(\mathbf{x}_{6}) \cdot \psi_{8}(\mathbf{x}_{5}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \sum_{x_{5}} \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{9}(\mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{8}(\mathbf{x}_{5}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{4}) &= \psi_{6}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{10}(\mathbf{x}_{2}, \mathbf{x}_{4}) \quad \widetilde{O}(N) \end{split}$$

Example 2: FAQ Computation with Indicator Projections

$$\varphi(x_1) = \sum_{x_2, x_3, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_2(x_2, x_3) \cdot \psi_3(x_3, x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1)$$



$$\begin{array}{cccc} X_1 & key(X_1) = \emptyset \\ & key(X_2) = \{X_1\} \\ X_2 & X_4 & key(X_3) = \{X_1, X_2\} \\ & & key(X_4) = \{X_1\} \\ & & key(X_4) = \{X_1\} \\ & key(X_5) = \{X_1, X_4\} \end{array}$$
$$\varphi(x_1) = \sum_{x_2, x_3, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_2(x_2, x_3) \cdot \psi_3(x_3, x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1)$$



• $\rho^*(\varphi) = 2.5$, $s(\varphi) = 1.5$, $fhtw(\varphi) = 1.5$. The above variable order Δ has the free variable x_1 on top of the others and $fhtw(\Delta) = 1.5$.

• The (unary) query result has size O(N) when factorized or listed.







$$\begin{split} \varphi(\mathbf{x}_{1}) &= \sum_{x_{2}, x_{3}, x_{4}, x_{5}} \psi_{1}(\mathbf{x}_{1}, x_{2}) \cdot \psi_{2}(x_{2}, x_{3}) \cdot \psi_{3}(x_{3}, x_{1}) \cdot \psi_{4}(\mathbf{x}_{1}, x_{4}) \cdot \psi_{5}(\mathbf{x}_{4}, x_{5}) \cdot \psi_{6}(\mathbf{x}_{5}, x_{1}) \\ \varphi(\mathbf{x}_{1}) &= \sum_{x_{2}, x_{3}, x_{4}, x_{5}} \psi_{1}(\mathbf{x}_{1}, x_{2}) \cdot \left(\underbrace{\sum_{x_{3}} \psi_{1}'(\mathbf{x}_{1}, x_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, x_{3}) \cdot \psi_{3}(\mathbf{x}_{3}, x_{1}) \cdot \psi_{4}'(\mathbf{x}_{1}) \cdot \psi_{6}'(\mathbf{x}_{1})}_{\psi_{7}(\mathbf{x}_{1}, \mathbf{x}_{2})} \right) \cdot \underbrace{\psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, x_{3}) \cdot \psi_{3}(\mathbf{x}_{3}, x_{1}) \cdot \psi_{4}'(\mathbf{x}_{1}) \cdot \psi_{6}'(\mathbf{x}_{1})}_{\psi_{7}(\mathbf{x}_{1}, \mathbf{x}_{2})} \end{split}$$



$$\begin{split} \varphi(\mathbf{x}_1) &= \sum_{x_2, x_3, x_4, x_5} \psi_1(\mathbf{x}_1, \mathbf{x}_2) \cdot \psi_2(\mathbf{x}_2, \mathbf{x}_3) \cdot \psi_3(\mathbf{x}_3, \mathbf{x}_1) \cdot \psi_4(\mathbf{x}_1, \mathbf{x}_4) \cdot \psi_5(\mathbf{x}_4, \mathbf{x}_5) \cdot \psi_6(\mathbf{x}_5, \mathbf{x}_1) \\ \varphi(\mathbf{x}_1) &= \sum_{x_2, x_4, x_5} \psi_1(\mathbf{x}_1, \mathbf{x}_2) \cdot \psi_7(\mathbf{x}_1, \mathbf{x}_2) \cdot \psi_4(\mathbf{x}_1, \mathbf{x}_4) \cdot \psi_5(\mathbf{x}_4, \mathbf{x}_5) \cdot \psi_6(\mathbf{x}_5, \mathbf{x}_1) \quad \widetilde{O}(N^{1.5}) \end{split}$$



$$\begin{split} \varphi(\mathbf{x}_{1}) &= \sum_{x_{2}, x_{3}, x_{4}, x_{5}} \psi_{1}(\mathbf{x}_{1}, x_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, x_{3}) \cdot \psi_{3}(\mathbf{x}_{3}, \mathbf{x}_{1}) \cdot \psi_{4}(\mathbf{x}_{1}, x_{4}) \cdot \psi_{5}(\mathbf{x}_{4}, x_{5}) \cdot \psi_{6}(\mathbf{x}_{5}, \mathbf{x}_{1}) \\ \varphi(\mathbf{x}_{1}) &= \sum_{x_{2}, x_{4}, x_{5}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{7}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{4}(\mathbf{x}_{1}, \mathbf{x}_{4}) \cdot \psi_{5}(\mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{6}(\mathbf{x}_{5}, \mathbf{x}_{1}) \quad \widetilde{O}(N^{1.5}) \\ \varphi(\mathbf{x}_{1}) &= \sum_{x_{4}, x_{5}} \left(\underbrace{\sum_{x_{2}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{7}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{4}'(\mathbf{x}_{1}) \cdot \psi_{6}'(\mathbf{x}_{1})}_{\psi_{8}(\mathbf{x}_{1})} \right) \cdot \underbrace{\psi_{8}(\mathbf{x}_{1})}_{\psi_{8}(\mathbf{x}_{1})} \end{split}$$



$$\begin{split} \varphi(x_1) &= \sum_{x_2, x_3, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_2(x_2, x_3) \cdot \psi_3(x_3, x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \\ \varphi(x_1) &= \sum_{x_2, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_7(x_1, x_2) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \quad \widetilde{O}(N^{1.5}) \\ \varphi(x_1) &= \sum_{x_4, x_5} \psi_8(x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \quad \widetilde{O}(N) \end{split}$$





$$\begin{split} \varphi(x_1) &= \sum_{x_2, x_3, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_2(x_2, x_3) \cdot \psi_3(x_3, x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \\ \varphi(x_1) &= \sum_{x_2, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_7(x_1, x_2) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \quad \widetilde{O}(N^{1.5}) \\ \varphi(x_1) &= \sum_{x_4, x_5} \psi_8(x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \quad \widetilde{O}(N) \\ \varphi(x_1) &= \sum_{x_4} \psi_8(x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_9(x_1, x_4) \quad \widetilde{O}(N^{1.5}) \end{split}$$



$$\begin{split} \varphi(\mathbf{x}_{1}) &= \sum_{x_{2}, x_{3}, x_{4}, x_{5}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{2}(\mathbf{x}_{2}, \mathbf{x}_{3}) \cdot \psi_{3}(\mathbf{x}_{3}, \mathbf{x}_{1}) \cdot \psi_{4}(\mathbf{x}_{1}, \mathbf{x}_{4}) \cdot \psi_{5}(\mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{6}(\mathbf{x}_{5}, \mathbf{x}_{1}) \\ \varphi(\mathbf{x}_{1}) &= \sum_{x_{2}, x_{4}, x_{5}} \psi_{1}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{7}(\mathbf{x}_{1}, \mathbf{x}_{2}) \cdot \psi_{4}(\mathbf{x}_{1}, \mathbf{x}_{4}) \cdot \psi_{5}(\mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{6}(\mathbf{x}_{5}, \mathbf{x}_{1}) \quad \widetilde{O}(N^{1.5}) \\ \varphi(\mathbf{x}_{1}) &= \sum_{x_{4}, x_{5}} \psi_{8}(\mathbf{x}_{1}) \cdot \psi_{4}(\mathbf{x}_{1}, \mathbf{x}_{4}) \cdot \psi_{5}(\mathbf{x}_{4}, \mathbf{x}_{5}) \cdot \psi_{6}(\mathbf{x}_{5}, \mathbf{x}_{1}) \quad \widetilde{O}(N) \\ \varphi(\mathbf{x}_{1}) &= \sum_{x_{4}} \psi_{8}(\mathbf{x}_{1}) \cdot \psi_{4}(\mathbf{x}_{1}, \mathbf{x}_{4}) \cdot \psi_{9}(\mathbf{x}_{1}, \mathbf{x}_{4}) \quad \widetilde{O}(N^{1.5}) \\ \varphi(\mathbf{x}_{1}) &= \psi_{8}(\mathbf{x}_{1}) \cdot \left(\sum_{x_{4}} \frac{\psi_{8}'(\mathbf{x}_{1}) \cdot \psi_{4}(\mathbf{x}_{1}, \mathbf{x}_{4}) \cdot \psi_{9}(\mathbf{x}_{1}, \mathbf{x}_{4})}{\psi_{10}(\mathbf{x}_{1})}\right) \\ \end{split}$$



$$\begin{split} \varphi(x_1) &= \sum_{x_2, x_3, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_2(x_2, x_3) \cdot \psi_3(x_3, x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \\ \varphi(x_1) &= \sum_{x_2, x_4, x_5} \psi_1(x_1, x_2) \cdot \psi_7(x_1, x_2) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \quad \widetilde{O}(N^{1.5}) \\ \varphi(x_1) &= \sum_{x_4, x_5} \psi_8(x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_5(x_4, x_5) \cdot \psi_6(x_5, x_1) \quad \widetilde{O}(N) \\ \varphi(x_1) &= \sum_{x_4} \psi_8(x_1) \cdot \psi_4(x_1, x_4) \cdot \psi_9(x_1, x_4) \quad \widetilde{O}(N^{1.5}) \\ \varphi(x_1) &= \psi_8(x_1) \cdot \psi_{10}(x_1) \quad \widetilde{O}(N) \end{split}$$

Outline



Part 1. Joins

Part 2. Aggregates

Part 3. Optimization

Optimization Inside the Database

Why solving optimization problems aka analytics inside the database?

- 1. Bring analytics close to data
 - \Rightarrow Save non-trivial export/import time
- 2. Large chunks of analytics code can be rewritten into SumProduct FAQs
 - $\Rightarrow \mathsf{Use \ scalable}/\mathsf{factorized \ query \ processing}$

Hot topic in the current DB research & industry landscape:

- Very recent tutorials and research agenda [A17,KBY17,PRWZ17]
- This tutorial highlights our recent work

[SOC16,ANNOS17]

In-database vs. Out-of-database Analytics



h and g are functions over features and respectively model parameters

• $heta^*$ are the parameters of the learned model

Plan for Part 3 on Optimization

• We will first introduce the main technical ideas via an example

- Train a linear regression model using batch gradient descent
- Express gradient computation as database queries
- Re-parameterize the model under functional dependencies
- We will then discuss a generalization
 - Polynomial regression, factorization machines, classification
- We will conclude with complexity & experimental analysis
 - Model training faster than computing the input to external ML library!

In-Database Analytics Approach in This Tutorial

Unified in-database analytics solution for a host of optimization problems.

Deployed in industrial retail-planning and forecasting applications

- Typical databases have weekly sales, promotions, and products
- Training dataset = Result of a feature extraction query over the database
- Task = Train parameterized model to predict, e.g., additional demand generated for a product due to promotion
- Training algorithm = First-order optimization algorithm, e.g., batch or stochastic gradient descent

Retail Example

Simplified Retail Example

- Database I = (R₁, R₂, R₃, R₄, R₅)
- Feature selection query *Q*:

Q(sku, store, color, city, country, *unitsSold*) =

 $R_1(sku, store, day, unitsSold), R_2(sku, color),$

 $R_3(\text{day}, \text{quarter}), R_4(\text{store}, \text{city}), R_5(\text{city}, \text{country}).$

- Free variables
 - Categorical (qualitative): F = {sku, store, color, city, country}.
 - Continuous (quantitative): unitsSold.
- Bounded variables
 - Categorical (qualitative): B = {day, quarter}
- We learn the ridge linear regression model $\langle \theta, \mathbf{x} \rangle = \sum_{f \in F} \langle \theta_f, \mathbf{x}_f \rangle$ over D = Q(I) with feature vector \mathbf{x} and response $y_{unitsSold}$.
- The parameters θ are obtained by minimizing the square loss function:

$$J(oldsymbol{ heta}) = rac{1}{2|D|} \sum_{(\mathbf{x},y)\in D} (\langle oldsymbol{ heta}, \mathbf{x}
angle - y_{unitsSold})^2 + rac{\lambda}{2} \|oldsymbol{ heta}\|_2^2$$

Recap: One-hot encoding of categorical variables

Continuous variables are mapped to scalars

• $y_{unitsSold} \in \mathbb{R}$.

Categorical variables are mapped to indicator vectors

Say variable country has categories vietnam and england.

• The variable country is then mapped to an indicator vector $\mathbf{x}_{\text{country}} = [x_{\text{vietnam}}, x_{\text{england}}]^{\top} \in (\{0, 1\}^2)^{\top}.$

• $\mathbf{x}_{\text{country}} = [0, 1]^{\top}$ for a tuple with country = ''england''

One-hot encoding leads to very wide training datasets and many 0-values.

Recap: Role of the Least Square Loss Function

Goal: Describe a linear relationship $fun(x) = \theta_1 x + \theta_0$ between variables x and y = fun(x), so we can estimate new y values given new x values.



- We are given n (black) data points $(x_i, y_i)_{i \in [n]}$
- We would like to find a (red) regression line fun(x) such that the (green) error $\sum_{i \in [n]} (fun(x_i) y_i)^2$ is minimized
- The role of the l₂-regularization ||θ||₂² = θ₀² + θ₁² is to avoid over/under-fitting. It gives preference to functions *fun* with smaller norms.

From Optimization to SumProduct FAQ Queries

We can solve $\theta^* := \arg \min_{\theta} J(\theta)$ by repeatedly updating θ in the direction of the gradient until convergence:

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \boldsymbol{\alpha} \cdot \boldsymbol{\nabla} J(\boldsymbol{\theta}).$$

From Optimization to SumProduct FAQ Queries

We can solve $\theta^* := \arg \min_{\theta} J(\theta)$ by repeatedly updating θ in the direction of the gradient until convergence:

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \boldsymbol{\alpha} \cdot \boldsymbol{\nabla} J(\boldsymbol{\theta}).$$

Define the matrix $\mathbf{\Sigma} = (\boldsymbol{\sigma}_{ij})_{i,j \in [|F|]}$, the vector $\mathbf{c} = (c_i)_{i \in [|F|]}$, and the scalar s_Y :

$$\boldsymbol{\sigma}_{ij} = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \mathbf{x}_i \mathbf{x}_j^\top \qquad \mathbf{c}_i = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} y \cdot \mathbf{x}_i \qquad \mathbf{s}_Y = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} y^2.$$

From Optimization to SumProduct FAQ Queries

We can solve $\theta^* := \arg \min_{\theta} J(\theta)$ by repeatedly updating θ in the direction of the gradient until convergence:

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \cdot \boldsymbol{\nabla} J(\boldsymbol{\theta}).$$

Define the matrix $\mathbf{\Sigma} = (\sigma_{ij})_{i,j \in [|F|]}$, the vector $\mathbf{c} = (c_i)_{i \in [|F|]}$, and the scalar s_Y :

$$\boldsymbol{\sigma}_{ij} = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \mathbf{x}_i \mathbf{x}_j^\top \qquad \quad \mathbf{c}_i = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} y \cdot \mathbf{x}_i \qquad \quad \mathbf{s}_Y = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} y^2.$$

Then,

$$J(\boldsymbol{\theta}) = \frac{1}{2|D|} \sum_{(\mathbf{x}, y) \in D} \left(\langle \boldsymbol{\theta}, \mathbf{x} \rangle - y \right)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

$$= \frac{1}{2} \boldsymbol{\theta}^{\top} \boldsymbol{\Sigma} \boldsymbol{\theta} - \langle \boldsymbol{\theta}, \mathbf{c} \rangle + \frac{s_{Y}}{2} + \frac{\lambda}{2} \| \boldsymbol{\theta} \|_{2}^{2}$$

$$\nabla J(\theta) = \Sigma \theta - \mathbf{c} + \lambda \theta$$

Expressing Σ , **c**, s_Y as SumProduct FAQ Queries FAQ queries for $\sigma_{ij} = \frac{1}{|D|} \sum_{(\mathbf{x}, \mathbf{y}) \in D} \mathbf{x}_i \mathbf{x}_j^{\top}$ (w/o factor $\frac{1}{|D|}$):

• x_i , x_j continuous \Rightarrow FAQ query with no free variable

$$\psi_{ij} = \sum_{f \in F: a_f \in \mathsf{Dom}(x_f)} \sum_{b \in B: a_b \in \mathsf{Dom}(x_b)} a_i \cdot a_j \cdot \prod_{k \in [5]} \mathbf{1}_{R_k(\mathbf{a}_{\mathcal{S}(R_k)})}$$

• x_i categorical, x_j continuous \Rightarrow FAQ query with one free variable

$$\psi_{ij}[a_i] = \sum_{f \in F - \{i\}: a_f \in \mathsf{Dom}(\mathsf{x}_f)} \sum_{b \in B: a_b \in \mathsf{Dom}(\mathsf{x}_b)} a_j \cdot \prod_{k \in [5]} \mathbf{1}_{R_k(\mathsf{a}_{\mathcal{S}(R_k)})}$$

• x_i , x_j categorical \Rightarrow FAQ query with two free variables

$$\psi_{ij}[a_i, a_j] = \sum_{f \in F - \{i, j\}: a_f \in \mathsf{Dom}(x_f)} \sum_{b \in B: a_b \in \mathsf{Dom}(x_b)} \prod_{k \in [5]} \mathbf{1}_{R_k(\mathsf{a}_{\mathcal{S}(R_k)})}$$

 $S(R_k)$ is the set of variables of R_k ; $\mathbf{a}_{S(R_k)}$ is a tuple in relation R_k ; $\mathbf{1}_E$ is the Kronecker delta that is 1 (0) whenever the event E holds (does not hold).

Expressing Σ , c, s_Y as SQL Queries

SQL queries for
$$\sigma_{ij} = \frac{1}{|D|} \sum_{(\mathbf{x},y) \in D} \mathbf{x}_i \mathbf{x}_j^{\top}$$
 (w/o factor $\frac{1}{|D|}$):

• x_i , x_j continuous \Rightarrow SQL query with no group-by attribute

SELECT SUM $(x_i * x_j)$ **FROM** *D*;

• x_i categorical, x_j continuous \Rightarrow SQL query with one group-by attribute

SELECT x_i , **SUM** (x_i) **FROM** *D* **GROUP BY** x_i ;

• x_i , x_j categorical \Rightarrow SQL query with two free variables

SELECT x_i, x_j , **SUM**(1) **FROM** *D* **GROUP BY** x_i, x_j ;

- **\Sigma**, **c**, s_Y are all aggregates that can be computed inside the database!
- We avoid one-hot/sparse encoding of the input data.

Consider the functional dependency city $\,\rightarrow\,$ country

There is one country for each city.

Assume we have:

- vietnam, england as categories for country
- saigon, hanoi, oxford, leeds,bristol as categories for city

The one-hot encoding enforces the following identities:

$$X_{ ext{vietnam}} = X_{ ext{saigon}} + X_{ ext{hanoi}}$$

That is: If country is vietnam, then city is either saigon or hanoi if $x_{\text{vietnam}} = 1$ then either $x_{\text{saigon}} = 1$ or $x_{\text{hanoi}} = 1$

That is: If country is england, then city is either oxford, leeds, or bristol if $x_{\text{england}} = 1$ then either $x_{\text{oxford}} = 1$ or $x_{\text{leeds}} = 1$ or $x_{\text{bristol}} = 1$

Identities due to one-hot encoding

 $x_{\texttt{vietnam}} = x_{\texttt{saigon}} + x_{\texttt{hanoi}}$

 $x_{\text{england}} = x_{\text{oxford}} + x_{\text{leeds}} + x_{\text{bristol}}$

 \blacksquare Encode $x_{\texttt{country}}$ as $x_{\texttt{country}} = Rx_{\texttt{city}},$ where

	saigon	hanoi	oxford	leeds	bristol	
$\mathbf{R} =$	1	1	0	0	0	vietnam
	0	0	1	1	1	england

For instance, if city is saigon, i.e., $\mathbf{x}_{\text{city}} = [1, 0, 0, 0, 0]^{\top}$, then country is vietnam, i.e., $\mathbf{x}_{\text{country}} = \mathbf{R}\mathbf{x}_{\text{city}} = [1, 0]^{\top}$.

$$\begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

 \blacksquare Functional dependency: city \rightarrow country

Replace all occurrences of $\mathbf{x}_{\text{country}}$ by $\mathbf{R}\mathbf{x}_{\text{city}}$:

$$\sum_{f \in F - \{\text{city, country}\}} \langle \boldsymbol{\theta}_{f}, \mathbf{x}_{f} \rangle + \langle \boldsymbol{\theta}_{\text{country}}, \mathbf{x}_{\text{country}} \rangle + \langle \boldsymbol{\theta}_{\text{city}}, \mathbf{x}_{\text{city}} \rangle$$
$$= \sum_{f \in F - \{\text{city, country}\}} \langle \boldsymbol{\theta}_{f}, \mathbf{x}_{f} \rangle + \langle \boldsymbol{\theta}_{\text{country}}, \mathbf{R} \mathbf{x}_{\text{city}} \rangle + \langle \boldsymbol{\theta}_{\text{city}}, \mathbf{x}_{\text{city}} \rangle$$
$$= \sum_{f \in F - \{\text{city, country}\}} \langle \boldsymbol{\theta}_{f}, \mathbf{x}_{f} \rangle + \left\langle \underbrace{\mathbf{R}^{\top} \boldsymbol{\theta}_{\text{country}} + \boldsymbol{\theta}_{\text{city}}}_{\gamma_{\text{city}}}, \mathbf{x}_{\text{city}} \right\rangle$$

 \blacksquare Functional dependency: city \rightarrow country

Replace all occurrences of $\mathbf{x}_{\text{country}}$ by $\mathbf{R}\mathbf{x}_{\text{city}}$:

$$\sum_{f \in F - \{\text{city, country}\}} \langle \boldsymbol{\theta}_{f}, \mathbf{x}_{f} \rangle + \langle \boldsymbol{\theta}_{\text{country}}, \mathbf{x}_{\text{country}} \rangle + \langle \boldsymbol{\theta}_{\text{city}}, \mathbf{x}_{\text{city}} \rangle$$
$$= \sum_{f \in F - \{\text{city, country}\}} \langle \boldsymbol{\theta}_{f}, \mathbf{x}_{f} \rangle + \langle \boldsymbol{\theta}_{\text{country}}, \mathbf{R} \mathbf{x}_{\text{city}} \rangle + \langle \boldsymbol{\theta}_{\text{city}}, \mathbf{x}_{\text{city}} \rangle$$
$$= \sum_{f \in F - \{\text{city, country}\}} \langle \boldsymbol{\theta}_{f}, \mathbf{x}_{f} \rangle + \left\langle \underbrace{\mathbf{R}^{\top} \boldsymbol{\theta}_{\text{country}} + \boldsymbol{\theta}_{\text{city}}}_{\gamma_{\text{city}}}, \mathbf{x}_{\text{city}} \right\rangle$$

- We avoid computing aggregates over **x**_{country}.
- We reparameterize the problem and ignore parameters θ_{country} .
- What about the penalty term in the loss function?

• Functional dependency: city \rightarrow country

x_{country} =
$$\mathbf{R}\mathbf{x}_{city}$$

•
$$\boldsymbol{\gamma}_{ ext{city}} = \boldsymbol{\mathsf{R}}^{ op} \boldsymbol{ heta}_{ ext{country}} + \boldsymbol{ heta}_{ ext{city}}$$

Rewrite the penalty term

$$\|\boldsymbol{\theta}\|_{2}^{2} = \sum_{j \neq \texttt{city}} \|\boldsymbol{\theta}_{j}\|_{2}^{2} + \left\|\boldsymbol{\gamma}_{\texttt{city}} - \boldsymbol{\mathsf{R}}^{\top}\boldsymbol{\theta}_{\texttt{country}}\right\|_{2}^{2} + \|\boldsymbol{\theta}_{\texttt{country}}\|_{2}^{2}$$

• "Optimize out" θ_{country} by expressing it in terms of γ_{city} :

$$\boldsymbol{\theta}_{\texttt{country}} = (\boldsymbol{\mathsf{I}}_{\texttt{country}} + \boldsymbol{\mathsf{R}}\boldsymbol{\mathsf{R}}^\top)^{-1}\boldsymbol{\mathsf{R}}\boldsymbol{\gamma}_{\texttt{city}} = \boldsymbol{\mathsf{R}}(\boldsymbol{\mathsf{I}}_{\texttt{city}} + \boldsymbol{\mathsf{R}}^\top\boldsymbol{\mathsf{R}})^{-1}\boldsymbol{\gamma}_{\texttt{city}}$$

I_{country} is the order-N_{country} identity matrix and similarly for I_{city}.
 The penalty term becomes

$$\|\boldsymbol{\theta}\|_2^2 = \sum_{j \neq \texttt{city}} \|\boldsymbol{\theta}_j\|_2^2 + \left\langle (\mathbf{I}_{\texttt{city}} + \mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\gamma}_{\texttt{city}}, \boldsymbol{\gamma}_{\texttt{city}} \right\rangle$$

The General Picture

General Problem Formulation

A typical machine learning task is to solve $\theta^* := \arg \min_{\theta} J(\theta)$, where

$$J(oldsymbol{ heta}) := \sum_{(\mathbf{x},y)\in D} \mathcal{L}\left(\left\langle g(oldsymbol{ heta}), h(\mathbf{x}) \right\rangle, y\right) + \Omega(oldsymbol{ heta}).$$

 $oldsymbol{ heta} oldsymbol{ heta} = oldsymbol{(} heta_1, \ldots, heta_p oldsymbol) \in oldsymbol{\mathsf{R}}^p$ are parameters

- functions $g: \mathbf{R}^{p} \to \mathbf{R}^{m}$ and $h: \mathbf{R}^{n} \to \mathbf{R}^{m}$ for *n* numeric features, m > 0
 - ▶ g = (g_j)_{j∈[m]} is a vector of multivariate polynomials
 - ▶ h = (h_j)_{j∈[m]} is a vector of multivariate monomials
- \blacksquare ${\cal L}$ is a loss function, Ω is the regularizer
- D is the training dataset with features x and response y.

Example problems: ridge linear regression, degree-*d* polynomial regression, degree-*d* factorization machines; logistic regression, SVM; PCA.

Special Case: Ridge Linear Regression

General problem formulation:

$$J(\boldsymbol{\theta}) := \sum_{(\mathbf{x},y)\in D} \mathcal{L}\left(\left\langle g(\boldsymbol{\theta}), h(\mathbf{x}) \right\rangle, y\right) + \Omega(\boldsymbol{\theta}).$$

Under

- square loss \mathcal{L} , ℓ_2 -regularization,
- data points $\mathbf{x} = (x_0, x_1, \dots, x_n, y)$,

$$oldsymbol{p} = n+1$$
 parameters $oldsymbol{ heta} = (heta_0, \dots, heta_n)$,

- $x_0 = 1$ corresponds to the bias parameter θ_0
- **g** and *h* identity functions $\mathbf{g}(\theta) = \theta$ and $h(\mathbf{x}) = \mathbf{x}$,

we obtain the following formulation for ridge linear regression:

$$J(oldsymbol{ heta}) := rac{1}{2|D|} \sum_{(\mathbf{x},y)\in D} \left(\langle oldsymbol{ heta}, \mathbf{x}
angle - y
ight)^2 + rac{\lambda}{2} \|oldsymbol{ heta}\|_2^2 \,.$$

Special Case: Degree-d Polynomial Regression

General problem formulation:

$$J(\boldsymbol{\theta}) := \sum_{(\mathbf{x}, y) \in D} \mathcal{L}\left(\left\langle g(\boldsymbol{\theta}), h(\mathbf{x}) \right\rangle, y \right) + \Omega(\boldsymbol{\theta}).$$

Under

 \blacksquare square loss $\mathcal L$, $\ell_2\text{-regularization},$

data points
$$\mathbf{x} = (x_0, x_1, \dots, x_n, y)$$
,

- $p = m = 1 + n + n^2 + \dots + n^d$ parameters $\theta = (\theta_a)$, where $\mathbf{a} = (a_1, \dots, a_n)$ is a tuple of non-negative integers such that $\|\mathbf{a}\|_1 \le d$. • $g(\theta) = \theta$,
- the components of h are given by $h_{\mathbf{a}}(\mathbf{x}) = \prod_{i=1}^{n} x_i^{a_i}$.

we obtain the following formulation for polynomial regression:

$$J(oldsymbol{ heta}) := rac{1}{2|D|} \sum_{(\mathbf{x},y)\in D} \left(\langle g(oldsymbol{ heta}), h(\mathbf{x})
angle - y
ight)^2 + rac{\lambda}{2} \left\| oldsymbol{ heta}
ight\|_2^2.$$

Special Case: Factorization Machines

Under

- square loss \mathcal{L} , ℓ_2 -regularization,
- data points $\mathbf{x} = (x_0, x_1, \dots, x_n, y)$,
- $p = m = 1 + n + r \cdot n$ parameters and $m = 1 + n + \binom{n}{2}$ features

we obtain the following formulation for degree-2 rank-r factorization machines:

$$J(\boldsymbol{\theta}) := \frac{1}{2|D|} \sum_{(\mathbf{x}, y) \in D} \left(\sum_{i=0}^{n} \theta_{i} x_{i} + \sum_{\substack{\{i, j\} \in \binom{[n]}{2} \\ \ell \in [r]}} \theta_{i}^{(\ell)} \theta_{j}^{(\ell)} x_{i} x_{j} - y \right)^{2} + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_{2}^{2}.$$

where

$$h_{S}(\mathbf{x}) = \prod_{i \in S} x_{i}, \text{ for } S \subseteq [n], |S| \le 2$$

$$g_{S}(\theta) = \begin{cases} \theta_{0} & \text{when } |S| = 0\\ \theta_{i} & \text{when } S = \{i\}\\ \sum_{\ell=1}^{r} \theta_{i}^{(\ell)} \theta_{j}^{(\ell)} & \text{when } S = \{i, j\} \end{cases}$$

Special Case: Classification methods

Examples: support vector machines, logistic regression, Adaboost

- Typically, the regularizer is $\frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$
- The response is now binary: $y \in \{\pm 1\}$
- The loss function $\mathcal{L}(\gamma, y)$, where $\gamma := \langle g(\theta), h(\mathbf{x}) \rangle$, takes the form:
 - $\mathcal{L}(\gamma, y) = \max\{1 y\gamma, 0\}$ for support vector machines (SVM),
 - $\mathcal{L}(\gamma, y) = \log(1 + e^{-y\gamma})$ for logistic regression, and
 - $\mathcal{L}(\gamma, y) = e^{-y\gamma}$ for Adaboost.

Batch Gradient Descent (BGD)

Repeatedly update heta in the direction of the gradient until convergence

 $\boldsymbol{ heta} \leftarrow \mathsf{a} \mathsf{ random point};$

while not converged yet do

$$\begin{array}{l} \alpha \leftarrow \text{next step size;} \\ \mathbf{d} \leftarrow \boldsymbol{\nabla} J(\boldsymbol{\theta}); \\ \text{while } \left(J(\boldsymbol{\theta} - \alpha \cdot \mathbf{d}) \geq J(\boldsymbol{\theta}) - \frac{\alpha}{2} \cdot \|\mathbf{d}\|_2^2 \right) \ \mathbf{do} \ \alpha \leftarrow \alpha/2; \ // \ \text{line search} \\ \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \cdot \mathbf{d}; \end{array}$$

end

BGD needs:

• Computation of the gradient vector $\nabla J(\theta)$

Its data-dependent component is computed once for all iterations

Point evaluation $J(\theta)$

 \blacktriangleright A few times per iteration to adjust α using line search
Compute Parameters θ using BGD

Immediate extension of the linear regression case discussed before.

Define the matrix $\mathbf{\Sigma} = (\sigma_{ij})_{i,j \in [m]}$, the vector $\mathbf{c} = (c_i)_{i \in [m]}$, and the scalar s_Y by

$$\begin{split} \boldsymbol{\Sigma} &= \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} h(\mathbf{x}) h(\mathbf{x})^{\top} \\ \mathbf{c} &= \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} y \cdot h(\mathbf{x}) \\ s_Y &= \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} y^2. \end{split}$$

Under square loss \mathcal{L} and ℓ_2 -regularization:

$$J(\theta) = \frac{1}{2}g(\theta)^{\top} \mathbf{\Sigma}g(\theta) - \langle g(\theta), \mathbf{c} \rangle + \frac{s_{Y}}{2} + \frac{\lambda}{2} \|\theta\|_{2}^{2}$$
$$\nabla J(\theta) = \frac{\partial g(\theta)^{\top}}{\partial \theta} \mathbf{\Sigma}g(\theta) - \frac{\partial g(\theta)^{\top}}{\partial \theta} \mathbf{c} + \lambda \theta$$

Summing Up

Insight #1:

\Sigma, **c**, s_Y are queries that can be computed **inside the database**!

They can take much less time than computing the feature extraction query

Insight #2:

The training dataset has repeating data blocks as it satisfies the join dependencies given by the feature extraction query.

> A factorized training dataset avoids this redundancy.

Insight #3:

- The training dataset has many functional dependencies in practice.
 - First learn a smaller, reparameterized model whose features functionally determine the left-out features, then map it back to the original model with both functionally determining and determined parameters

Zoom-in: In-database vs. Out-of-database Learning



Complexity & Experimental Analysis

Complexity Analysis: The General Case

Complexity of learning models falls back to factorized computation of aggregates over joins [BKOZ13,OZ15,SOC16,ANR16] Let:

- $(\mathcal{V}, \mathcal{E}) =$ hypergraph of Q
- $N = \max_{R \in I} |R|$
- $|\sigma_{ij}| = size$ of the sparse representation of the σ_{ij} tensor
- faqw(i,j) = FAQ-width of the query that expresses σ_{ij} over Q

The tensors σ_{ij} and c_j can be sparsely represented by queries with group-by variables and can be computed in time

$$\widetilde{O}\left(\left|\mathcal{V}\right|^2 \cdot |\mathcal{E}| \cdot \sum_{i,j \in [m]} (N^{\mathsf{faqw}(i,j)} + |\sigma_{ij}|)\right)$$

Complexity Analysis: Continuous Features Only

Complexity in the general case:

[ANNOS17]

$$\widetilde{O}\left(\left|\mathcal{V}\right|^2 \cdot |\mathcal{E}| \cdot \sum_{i,j \in [m]} (N^{\mathsf{faqw}(i,j)} + |\sigma_{ij}|)\right).$$

Complexity in case all features are continuous:

[SOC16]

 $\widetilde{O}(|\mathcal{V}|^2 \cdot |\mathcal{E}| \cdot m^2 \cdot N^{fhtw}).$

In this case, faqw(i, j) becomes the fractional hypertree width *fhtw* of Q.

Complexity Analysis: Comparison with State of the Art

Let:

- d = degree of polynomial regression model
- $c = \max$ number of variables in any monomial in h; $c \le d$
- $\rho^* =$ fractional edge cover number of query Q

Comparison against state of the art:

[ANNOS17]

- $faqw(i,j) \leq fhtw + c 1$ and $|\sigma_{ij}| \leq \min\{|D|, N^c\}$.
- For any query Q with ρ* > fhtw + c − 1, there are infinitely many database instances of size N for which

$$\lim_{N\to\infty}\frac{|D|}{\sum_{i,j\in[m]}(N^{\mathsf{faqw}(i,j)}+|\boldsymbol{\sigma}_{ij}|)\log N}=\infty.$$

Computing σ_{ij} for degree-d polynomial regression takes

$$\widetilde{O}(|\mathcal{V}|^2 \cdot |\mathcal{E}| \cdot m^2 \cdot N^{\text{fhtw}+2d}).$$

under one-hot encoding of categorical variables.

Factorized Machine Learning in Practice

Experiments published in several papers, here a glimpse from

[ANNOS17]

Retailer dataset (records)	excerpt (17M)	full (86M)
Linear Regression		
Number of features (Cont. + Categ.)	33+55	33+3702
MADIib (ols)	1,898.35 sec	-
PostgreSQL + R (qr decomposition)	798.96 sec	_
FDB	25.53 sec	380.31 sec
Polynomial Regression degree 2		
Number of features (Cont. + Categ.)	562+2363	562+154K
MADlib	> 22h	-
PostgreSQL + R	-	-
FDB	135.7 sec	2039.31 sec

- We measure end-to-end performance: joins + aggregates + convergence
- R: "-" means R's data frame limit is exceeded and cannot run.
- MADlib: "-" means it cannot one-hot encode the data in a relation with more than 1600 columns.

From Joins to Aggregates and Optimization Problems

One idea to rule them all

and at their core FACTORIZE them!

Thank you!

Quizzes

QUIZ 1: Joins (1/3)

For each of the following queries, please show the following:

- 1. Hypertree decomposition and variable order for query.
- 2. The fractional edge cover number and the fractional hypertree width (assume all relations have the same size).

Path Query of length n:

 $P_n(X_1,\ldots,X_{n+1})=R_1(X_1,X_2),R_2(X_2,X_3),R_3(X_3,X_4),\ldots,R_n(X_n,X_{n+1}).$

QUIZ 1: Joins (2/3)

For each of of the following queries, please show the following:

- 1. Hypertree decomposition and variable order for query.
- 2. The fractional edge cover number and the fractional hypertree width (assume all relations have the same size).

Bowtie Query:

 $Q_{\bowtie}(A, B, C, D, E) = R_1(A, C), R_2(A, B), R_3(B, C), R_4(C, E), R_5(E, D), R_6(C, D).$



QUIZ 1: Joins (3/3)

For each of of the following queries, please show the following:

- 1. Hypertree decomposition and variable order for query.
- 2. The fractional edge cover number and the fractional hypertree width (assume all relations have the same size).

Loomis-Whitney Queries of length n: A LW_n query has n variables X_1, \ldots, X_n and n relation symbols such that for every $i \in [n]$ the relation symbol R_i has variables $\{X_1, \ldots, X_n\} - \{X_i\}$:

$$LW_n(X_1,...,X_n) = R_1(X_2,...,X_n),...,R_i(X_1,...,X_{i-1},X_{i+1},...,X_n),...,$$
$$R_n(X_1,...,X_{n-1})$$

 LW_3 is the triangle query.

QUIZ 2: Aggregates (1/2)

For each of of the following functional aggregate queries:

- 1. Give a hypertree decomposition and variable order.
- 2. If you were to compute it as stated below (with all sums done after the products), what would be its time complexity? (Assume all functions have the same size.)
- 3. Is there an equivalent rewriting of φ that would allow for quadratic time complexity? What about linear time?

The *n*-hop query:

$$\varphi(x_1, x_{n+1}) = \sum_{x_2, \dots, x_n} \psi_1(X_1, X_2) \cdot \psi_2(X_2, X_3) \cdot \psi_3(X_3, X_4) \cdot \dots \cdot \psi_n(X_n, X_{n+1}).$$

QUIZ 2: Aggregates (2/2)

For each of of the following functional aggregate queries:

- 1. Give a hypertree decomposition and variable order.
- 2. If you were to compute it as stated below (with all sums done after the products), what would be its time complexity? Assume all functions have the same size.
- 3. Is there an equivalent rewriting of φ that would allow for quadratic time complexity? What about linear time?

Query:

$$\varphi = \sum_{a} \sum_{b} \sum_{c} \sum_{f} \sum_{d} \sum_{e} \psi_1(a, b) \cdot \psi_2(a, c) \cdot \psi_3(c, d) \cdot \psi_4(b, c, e) \cdot \psi_5(e, f).$$

QUIZ 3: Optimization

Assume that the natural join of the following relations provides the features we use to predict revenue:

Sales(store_id, product_id, quantity, revenue),
Product(product_id, color),
Store(store_id, distance_city_center).

Variables revenue, quantity, and distance_city_center stand for continuous features, while product_id and color for categorical features.

- 1. Give the FAQs required to compute the gradient of the squares loss function for learning a ridge linear regression models with the above features.
- 2. We know that product_id functionally determines color. Give a rewriting of the objective function that exploits the functional dependency.
- 3. The FAQs require the computation of a lot of common sub-problems. Can you think of ways to share as much computation as possible?

References

References on Join Computation

- GLS99 Hypertree decompositions and tractable queries. Gottlob, Leone, Scarcello. In PODS 1999. https://arxiv.org/abs/cs/9812022
- AGM08 Size bounds and query plans for relational joins. Atserias, Grohe, Marx. In FOCS 2008 and SIAM J. Comput., 42(4) 2013. http://epubs.siam.org/doi/10.1137/110859440
 - M10 Approximating fractional hypertree width. Marx. In ACM TALG 2010. urlhttp://dl.acm.org/citation.cfm?id=1721845
- NPRR12 Worst-case optimal join algorithms: [extended abstract] Ngo, Porat, Ré, Rudra. In PODS 2012. urhttp://dl.acm.org/citation.cfm?id=2213565
 - OZ12 Factorised representations of query results: size bounds and readability. Olteanu, Zavodny. In ICDT 2012. http://dl.acm.org/citation.cfm?doid=2274576.2274607 Also https://arxiv.org/abs/1104.0867, April 2011.
 - NRR13 Skew Strikes Back: New Developments in the Theory of Join Algorithms. Ngo, Ré, Rudra. In SIGMOD Rec. 2013. https://arxiv.org/abs/1310.3314

References on Join Computation

V14 Triejoin: A Simple, Worst-Case Optimal Join Algorithm. Veldhuizen. In ICDT 2014. http://openproceedings.org/ICDT/2014/paper_13.pdf

OZ15 Size Bounds for Factorised Representations of Query Results. Olteanu, Zavodny. In ACM TODS 2015. http://dl.acm.org/citation.cfm?doid=2656335

CO15 Worst-Case Optimal Join At A Time. Ciucanu, Olteanu. Technical report, Oxford, Nov 2015.

ANS17 What do Shannon-type inequalities, submodular width, and disjunctive Datalog have to do with one another? Abo Khamis, Ngo, Suciu. In PODS 2017. https://arxiv.org/abs/1612.02503

KO17 Covers of Query Results.

Kara, Olteanu. Technical report, Oxford, March 2017. To appear in ICDT 2018. https://arxiv.org/abs/1709.01600

References on Aggregate Computation

BKOZ13 Aggregation and Ordering in Factorised Databases. Bakibayev, Kocisky, Olteanu, Zavodny. In PVLDB 2013. https://arxiv.org/abs/1307.0441

ANR16 FAQ: Questions Asked Frequently. Abo Khamis, Ngo, Rudra. In PODS 2016. https://arxiv.org/abs/1504.04044

References on In-Database Analytics

SOC16 Learning Linear Regression Models over Factorized Joins. Schleich, Olteanu, Ciucanu. In SIGMOD 2016. http://dl.acm.org/citation.cfm?doid=2882903.2882939

A17 Research Directions for Principles of Data Management (Dagstuhl Perspectives Workshop 16151). Abiteboul et al. In SIGMOD Rec. 2017. https://arxiv.org/pdf/1701.09007.pdf

ANNOS17 In-Database Learning with Sparse Tensors. Abo Khamis, Ngo, Nguyen, Olteanu, Schleich. March 2017. To appear in PODS 2018. https://arxiv.org/abs/1703.04780

KBY17 Data Management in Machine Learning: Challenges, Techniques, and Systems.

Kumar, Boehm, Yang. In SIGMOD 2017, Tutorial. https://www.youtube.com/watch?v=U8J0Dd_Z5wo

NO17 Incremental View Maintenance with Triple Lock Factorisation Benefits. Nikolic, Olteanu. March 2017. To appear in SIGMOD 2018. https://arxiv.org/abs/1703.07484

PRWZ17 Data Management Challenges in Production Machine Learning. Polyzotis, Roy, Whang, Zinkevich. In SIGMOD 2017, Tutorial. http://dl.acm.org/citation.cfm?doid=3035918.3054782