

Logic & Algorithms \in DB & AI – Simons 2023

Output Cardinality Bounds and Information Theory

Hung Q. Ngo



Outline

The Query Optimization (and Evaluation) Problem

Cardinality Bounds and Worst-Case Optimal Joins

The Bound Hierarchy Under Degree Constraints

Research Questions

References

Listing Triangles

In English

Given a (directed/undirected) graph G , find all triangles in G .

In Logic

$$Q(a, b, c) \leftarrow E(a, b) \wedge E(b, c) \wedge E(c, a)$$

$$Q(a, b, c) \leftarrow E(a, b) \wedge E(b, c) \wedge E(c, a) \wedge a < b \wedge b < c$$

Or, more generally:

$$Q(a, b, c) \leftarrow R(a, b) \wedge S(b, c) \wedge T(c, a)$$

4-Cycle Detection

In English

Given a graph G , does it have a 4-cycle?

In Logic

$$Q() \leftarrow \exists a, b, c, d \ E(a, b) \wedge E(b, c) \wedge E(c, d) \wedge E(d, a)$$

In English

Given a graph G , how many 3-walks are there in G ?

In Sum-Product Form

$$Q() = \sum_{a,b,c,d} E(a,b) \cdot E(b,c) \cdot E(c,d)$$

In English

Given a graph G , how many 3-paths are there in G ?

In Sum-Product Form

$$Q() = \sum_{a,b,c,d} E(a,b) \cdot E(b,c) \cdot E(c,d) \cdot \mathbf{1}_{a \neq b} \cdot \mathbf{1}_{a \neq c} \cdot \mathbf{1}_{a \neq d} \cdot \mathbf{1}_{b \neq c} \cdot \mathbf{1}_{b \neq d} \cdot \mathbf{1}_{c \neq d}$$

All-Pairs Shortest Paths (APSP)

In English

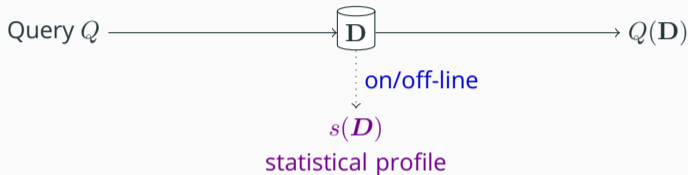
Given a graph G , compute the shortest path lengths between every pair of vertices.

In Datalogo (recursive query!)

$$Q[x, y] = \min \left(E[x, y], \min_z \{ Q[x, z] + E[z, y] \} \right) \quad \text{linear-form}$$

$$Q[x, y] = \min \left(E[x, y], \min_z \left\{ Q[x, z] + Q[z, y] \right\} \right) \quad \text{binary-form}$$

The Main Query Optimization / Evaluation Problem



Given Q and D , compute $Q(D)$ in the most efficient (optimal!?) way possible.

Incremental View Maintenance (IVM) (a.k.a. Dynamic Algorithms)

Given an update to D , how do we update $Q(D)$ efficiently?

Precise Problem Formulation

- What do you mean by “query”?
 - Full conjunctive queries
 - Sum-product queries
 - ... First-Order, Second-Order, Rank-Enumeration
- What is in the statistical profile $s(D)$?
 - Degree constraints
 - ... Frequency moment constraints, histograms, samples, ML models
- What do you mean by “optimality”?
 - Worst-Case Optimality
 - Instance Optimality
 - ... Fine-grained complexity.
- Optimizer designed to work before seeing Q .
 - Optimizer = *Meta-Algorithm* (input: problem, output: algorithm)
 - Tutorial on 3 meta-algorithms: join, variable elimination, tensor decomposition

Outline

The Query Optimization (and Evaluation) Problem

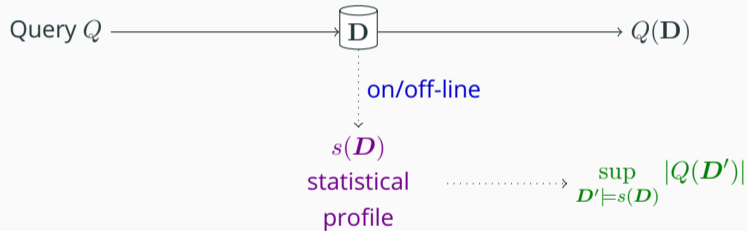
Cardinality Bounds and Worst-Case Optimal Joins

The Bound Hierarchy Under Degree Constraints

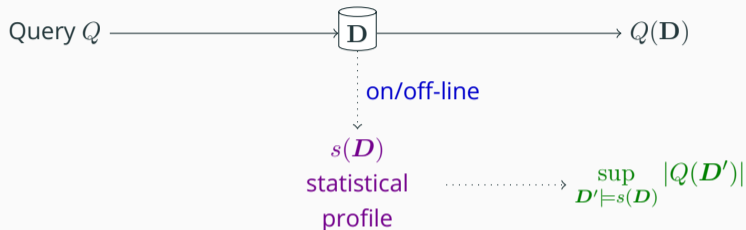
Research Questions

References

Worst-Case Cardinality Bound



Worst-Case Optimal Join (WCOJ) Algorithm



Definition

A “worst-case optimal” join algorithm is an algorithm computing $Q(D)$ in time

$$\tilde{O} \left(|D| + \sup_{D' \models s(D)} |Q(D')| \right)$$

\tilde{O} hides log and query dependent factors

(For Now) Q is a Full Conjunctive Queries

In a movie database

$$\begin{aligned} Q(\text{director}, \text{actor}, \text{movie}, \text{actor_age}, \text{name}) \leftarrow & \\ & \text{parent}(\text{director}, \text{actor}) \\ & \wedge \text{acted_in}(\text{actor}, \text{movie}) \\ & \wedge \text{director_of}(\text{director}, \text{movie}) \\ & \wedge \text{age}(\text{actor}, \text{actor_age}) \wedge (20 < \text{actor_age} \vee \text{actor_age} \neq 10) \\ & \wedge \text{person_name}(\text{director}, \text{name}) \wedge \text{regex_match}(".*\text{spiel}.*", \text{name}) \end{aligned}$$

In a graph database with edge relation E ,

$$Q(a, b, c) \leftarrow E(a, b) \wedge E(a, c) \wedge E(b, c)$$

(For Now) Q is a Full Conjunctive Queries

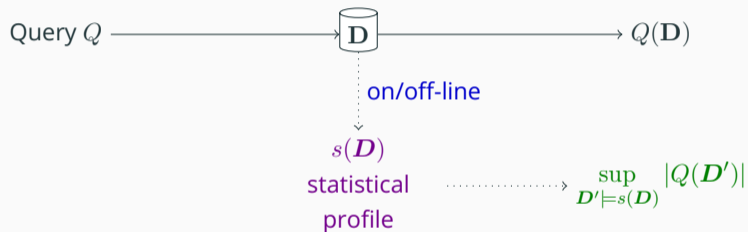
More generally, $\mathcal{H} = (V, \mathcal{E})$ is the hypergraph of a query:

$$Q(\mathbf{X}_V) \leftarrow \bigwedge_{S \in \mathcal{E}} R_S(\mathbf{X}_S)$$

For example $Q(a, b, c) \leftarrow E(a, b) \wedge E(a, c) \wedge E(b, c)$

- $V = \{a, b, c\}$
- $\mathcal{H} = (V, \mathcal{E}) = (V, \{ab, ac, bc\})$
- $R_F = E$ for all $F \in \mathcal{E}$.

What is in the Statistical Profile $s(D)$?



Degree constraints

Relation $R(\text{actor}, \text{movie}, \text{role})$, imagine a **frequency vector** d_{actor} (billions of entries)

| actor | movie | role |
|-------|-------|------|
| alice | | |
| bob | | |
| bob | | |
| bob | | |
| bob | | |
| carol | | |
| carol | | |

$$d_{\text{actor}}(\text{alice}) = 1 \quad d_{\text{actor}}(\text{bob}) = 4$$

$$d_{\text{actor}}(\text{carol}) = 2$$

$$d_{\text{actor}}(v) = 0 \quad v \notin \{\text{alice}, \text{bob}, \text{carol}\}$$

The profile $s(D)$ contains **degree constraints**:

- $\|d_{\text{actor}}\|_{\infty} = 4$ (degree constraint!)
- $\|d_{\emptyset}\|_{\infty} = 7 = |R|$ (cardinality constraint!)
- $\|d_{\text{actor}, \text{movie}}\|_{\infty} = 1$ (functional dependency)

General DC : (X, Y, N) in relation R means $|\pi_Y \sigma_{X=x} R| \leq N, \forall x$

Outline

The Query Optimization (and Evaluation) Problem

Cardinality Bounds and Worst-Case Optimal Joins

The Bound Hierarchy Under Degree Constraints

Research Questions

References

Hierarchy of Set Functions

$h : 2^{[n]} \rightarrow \mathbb{R}_+$, non-negative, monotone, $h(\emptyset) = 0$, $h(X) \leq h(Y)$ if $X \subseteq Y$

$SA_n := \{h \mid h \text{ is sub-additive}\} \quad h(X \cup Y) \leq h(X) + h(Y)$

$\Gamma_n := \{h \mid h \text{ is submodular}\} = \text{polymatroids}$

$h(X \cup Y) + h(X \cap Y) \leq h(X) + h(Y)$

$\bar{\Gamma}_n^*$: topological closure of Γ_n^* , *almost entropic*

$\Gamma_n^* = \{h : h \text{ is entropic}\}$

N_n : Normal convex-hull of step functions
(weighted coverage functions)
(non-negative multivariate mutual information)

M_n : Modular $h(X) = \sum_{x \in X} h(x)$

$$\begin{aligned} \log \sup_{\mathbf{D}' \models s(\mathbf{D})} |Q(\mathbf{D}')| &= \text{combinatorial-bound}(Q, s) && \text{computable but impractical} \\ &\leq \text{entropic-bound}(Q, s) \\ &\leq \text{polymatroid-bound}(Q, s) \\ &\leq \text{flow-bound}(Q, s, \sigma) \\ &\leq \text{chain-bound}(Q, s, \sigma) \\ &\leq \text{agm-bound}(Q, s) \\ &\leq \text{integral-edge-cover}(Q, s) \end{aligned}$$

- $Q(a, b, c) = R(a, b) \wedge S(b, c) \wedge T(a, c)$
- $D = \{R, S, T\}$
- $s(D) = \{|R|, |S|, |T|\}$

$|R|, |S|, |T|$ are integers, in the order of 10^9 or more

$$Q(a, b, c) \leftarrow R(a, b) \wedge S(b, c) \wedge T(a, c)$$

$$s(D) = \{|R|, |S|, |T|\}$$

$$\max \sum_{a,b,c} \mathbf{1}_{R(a,b)} \mathbf{1}_{S(a,c)} \mathbf{1}_{T(b,c)}$$

$$\mathbf{1}_X \in \{0, 1\}$$

$$\text{s.t.} \sum_{a,b} \mathbf{1}_{R(a,b)} \leq |R|$$

$$\sum_{b,c} \mathbf{1}_{S(b,c)} \leq |S|$$

$$\sum_{a,c} \mathbf{1}_{T(a,c)} \leq |T|.$$

Can turn this into a *linear integer program*, but does not help much.

$$Q(a, b, c) \leftarrow R(a, b) \wedge S(b, c) \wedge T(a, c)$$

$$s(\mathbf{D}) = \{|R|, |S|, |T|\}$$

- Fix a *worst-case* input R, S, T . Select tuples $(a, b, c) \in Q$ uniformly at random
- Let h be the entropy function of this 3D-distribution, then

$$\log \sup_{\mathbf{D}' \models s(\mathbf{D})} |Q(\mathbf{D}')| = h(a, b, c)$$

$$h(a, b) \leq \log |R|$$

$$h(b, c) \leq \log |S|$$

$$h(a, c) \leq \log |T|$$

$$h \in \Gamma_3^*$$

h is entropic

$$Q(a, b, c) \leftarrow R(a, b) \wedge S(b, c) \wedge T(a, c)$$

$$s(\mathbf{D}) = \{|R|, |S|, |T|\}$$

$$\log \sup_{\mathbf{D}' \models s(\mathbf{D})} |Q(\mathbf{D}')| \leq \max h(a, b, c)$$

$$\begin{aligned} \text{s.t.} \quad & h(a, b) \leq \log |R|, \\ & h(b, c) \leq \log |S|, \\ & h(a, c) \leq \log |T|, \\ & h \in \Gamma_3^* \end{aligned}$$

h is entropic

$$Q(a, b, c) \leftarrow R(a, b) \wedge S(b, c) \wedge T(a, c)$$

$$s(\mathbf{D}) = \{|R|, |S|, |T|\}$$

$$\log \sup_{\mathbf{D}' \models s(\mathbf{D})} |Q(\mathbf{D}')| \leq \max_{a, b, c} h(a, b, c)$$

$$\begin{aligned} \text{s.t.} \quad & h(a, b) \leq \log |R|, \\ & h(b, c) \leq \log |S|, \\ & h(a, c) \leq \log |T|, \\ & h \in \Gamma_3 \end{aligned}$$

h is a polymatroid

$$\max\{h(a, b, c) \mid h(a, b) \leq \log |R|, h(b, c) \leq \log |S|, h(a, c) \leq \log |T|, h \in \Gamma_3\}$$

Define a **modular** $g \in M_3$ as follows, then g satisfies all constraints:

$$g(a) = h(a) \quad g(ab) = g(a) + g(b) \quad = h(ab)$$

$$g(b) = h(b|a) \quad g(bc) = g(b) + g(c) = h(abc) - h(a) \quad \leq h(bc)$$

$$g(c) = h(c|ab) \quad g(ac) = g(a) + g(c) = h(abc) + h(a) - h(ab) \quad \leq h(ac)$$

$$g(abc) = h(abc)$$

Problem can be reformulated, optimized over *modular* functions g :

$$\max \quad g(a) + g(b) + g(c)$$

$$g(a) + g(b) \leq \log |R|, \quad g(b) + g(c) \leq \log |S|, \quad g(a) + g(c) \leq \log |T|,$$

$$g(a), g(b), g(c) \geq 0$$

Vertex packing LP

$$\max \quad g(a) + g(b) + g(c)$$

$$g(a) + g(b) \leq \log |R|, \quad g(b) + g(c) \leq \log |S|, \quad g(a) + g(c) \leq \log |T|,$$

$$g(a), g(b), g(c) \geq 0$$

Fractional edge cover LP

[AGM 2008] took the dual of the above LP:

$$\min \quad \lambda_{ab} \log |R| + \lambda_{bc} \log |S| + \lambda_{ac} \log |T|$$

$$\lambda_{ab} + \lambda_{ac} \geq 1$$

$$\lambda_{ab} + \lambda_{bc} \geq 1$$

$$\lambda_{bc} + \lambda_{ac} \geq 1$$

$$\boldsymbol{\lambda} \geq \mathbf{0}.$$

$$Q(a, b, c) \leftarrow R(a, b) \wedge S(b, c) \wedge T(a, c)$$

$$s(\mathbf{D}) = \{|R|, |S|, |T|\}$$

$$\log \sup_{\mathbf{D}' \models s(\mathbf{D})} |Q(\mathbf{D}')| \leq \log \min\{|R| \cdot |S|, |R| \cdot |T|, |S| \cdot |T|\}$$

We used the example to illustrate the following bounds / concepts:

$$\begin{aligned} \log \sup_{D' \models s(D)} |Q(D')| &= \text{combinatorial-bound}(Q, s) \quad (\text{computable but impractical}) \\ &\leq \text{entropic-bound}(Q, s) \\ &= \text{polymatroid-bound}(Q, s) \\ &= \text{modular-bound}(Q, s) \\ &= \text{agm-bound}(Q, s) \\ &\leq \text{integral-edge-cover}(Q, s). \end{aligned}$$

Let g be an optimal solution to the modular bound:

$$\begin{aligned} \max \quad & g(a) + g(b) + g(c) \\ & g(a) + g(b) \leq \log |R|, \quad g(b) + g(c) \leq \log |S|, \quad g(a) + g(c) \leq \log |T|, \\ & g(a), g(b), g(c) \geq 0 \end{aligned}$$

Then, $\sup_{D' \models s(D)} |Q(D')| \leq \text{modular-bound} = 2^{g(a)+g(b)+g(c)} = 2^{g(a)} \times 2^{g(b)} \times 2^{g(c)}$

Construct a database instance D'

[AGM 08]

- $R = \llbracket 2^{g(a)} \rrbracket \times \llbracket 2^{g(b)} \rrbracket$, $S = \llbracket 2^{g(b)} \rrbracket \times \llbracket 2^{g(c)} \rrbracket$, $T = \llbracket 2^{g(a)} \rrbracket \times \llbracket 2^{g(c)} \rrbracket$
- Then $D' \models s(D)$ and $|Q| \geq \frac{1}{8} 2^{g(a)+g(b)+g(c)} = \Omega\left(\sup_{D' \models s(D)} |Q(D')|\right)$

- Given a set of *degree constraints* (DCs)
 - Triples (X, Y, N) which says, for each x there are at most N y 's
 - If $X = \emptyset$, then this is a *cardinality constraint* (CC) (e.g. *distinct counts*)
 - If $N = 1$, then this is a *functional dependency* (FD) (very common)
- Entropy argument implies

$$\log \sup_{D' \models s(D)} |Q(D')| \leq \text{entropic-bound}$$

The Entropic and Polymatroid Bounds

Theorem (ANS 17)

If $s(\mathcal{D})$ contains only degree constraints, then

$$\log \sup_{\mathcal{D}' \models s(\mathcal{D})} |Q(\mathcal{D}')| \leq \sup_{h \in \bar{\Gamma}_n^* \cap DC} h(V) \leq \max_{h \in \Gamma_n \cap DC} h(V)$$

where DC is the set of linear constraints of the form

$$h(Y|X) \leq \log N$$

for each degree constraint (X, Y, N) .

Example: the Triangle Query

- $R(a, b) \wedge S(b, c) \wedge T(a, c)$ $D = \{R, S, T\}$
- $s(D) = \{|R|, |S|, |T|\}$
- Constraint set:

$$DC = \{h \mid h(ab) \leq \log |R| \wedge h(bc) \leq \log |S| \wedge h(ac) \leq \log |T|\}$$

- Polymatroid bound:

$$\max\{h(abc) \mid h \in \Gamma_3 \cap DC\} = \log \min\{|R| \cdot |S|, |S| \cdot |T|, |T| \cdot |R|, \sqrt{|R| \cdot |S| \cdot |T|}\}$$

e.g. if $|R|, |S|, |T| = N$, then $|Q| \leq N^{3/2}$ (Loomis-Whitney inequality!)

Example: the Triangle Query with Extra FD Information

- $R(a, b) \wedge S(b, c) \wedge T(a, c)$ $D = \{R, S, T\}$
- $s(D) = \{|R|, |S|, |T|, b \rightarrow c\}$ (b is a key in S)
- Constraint set:

$$DC = \{h \mid h(ab) \leq \log |R| \wedge h(bc) \leq \log |S| \wedge h(ac) \leq \log |T| \wedge h(c|b) = 0\}$$

- Polymatroid bound:

$$\max\{h(abc) \mid h \in \Gamma_3 \cap DC\} = \log \min\{|R|, |S| \cdot |T|\}$$

e.g. $|R|, |S|, |T| = N$, then $|Q| \leq N$

Example: Builtins and FDs

- $R(a) \wedge S(b) \wedge a + b = 5$ $D = \{R, S\}$
- $s(D) = \{|R|, |S|, a \rightarrow b, b \rightarrow a\}$
- Constraint set:

$$DC = \{h(a) \leq \log |R| \wedge h(b) \leq \log |S| \wedge h(a|b) = h(b|a) = 0\}$$

- Polymatroid bound:

$$\max\{h(ab) \mid h \in \Gamma_2 \cap DC\} = \log \min\{|R|, |S|\}$$

Example: A Non-Trivial Bound

• $R(a, b) \wedge S(b, c) \wedge T(c, d) \wedge f_1(a, c) = d \wedge f_2(b, d) = a$ f_1, f_2 are UDFs

• $s(\mathbf{D}) = \{|R|, |S|, |T|, ac \rightarrow d, bd \rightarrow a\}$

• Constraint set:

$$\text{DC} = \{h \mid h(ab) \leq \log |R| \wedge h(bc) \leq \log |S| \wedge h(cd) \leq \log |T| \wedge h(d|ac) = h(a|bd) = 0\}$$

• Polymatroid bound:

$$\max\{h(abcd) \mid h \in \Gamma_4 \cap \text{DC}\} = \log \min\{|R| \cdot |S|, |S| \cdot |T|, |T| \cdot |R|, \sqrt{|R| \cdot |S| \cdot |T|}\}$$

Outline

The Query Optimization (and Evaluation) Problem

Cardinality Bounds and Worst-Case Optimal Joins

The Bound Hierarchy Under Degree Constraints

Research Questions

References

Main Classes of (Mostly) Open Problems

1. Is the entropic bound computable?
2. Is the polymatroid bound computable in PTIME?
3. Find classes of inputs where the polymatroid bound is computable in PTIME.
4. When is which bound asymptotically tight?
5. Approximate the bounds efficiently. Hardness of approximation.
6. Going beyond conjunctive queries?
7. Dealing with more constraints
 - Conditional independence constraints
 - Frequency moment constraints
 - Histogram constraints
 - ...

Hierarchy of Set Functions

$h : 2^{[n]} \rightarrow \mathbb{R}_+$, non-negative, monotone, $h(\emptyset) = 0$, $h(X) \leq h(Y)$ if $X \subseteq Y$

$SA_n := \{h \mid h \text{ is sub-additive}\} \quad h(X \cup Y) \leq h(X) + h(Y)$

$\Gamma_n := \{h \mid h \text{ is submodular}\} = \text{polymatroids}$

$h(X \cup Y) + h(X \cap Y) \leq h(X) + h(Y)$

$\bar{\Gamma}_n^*$: topological closure of Γ_n^* , *almost entropic*

$\Gamma_n^* = \{h : h \text{ is entropic}\}$

N_n : Normal convex-hull of step functions
(weighted coverage functions)
(non-negative multivariate mutual information)

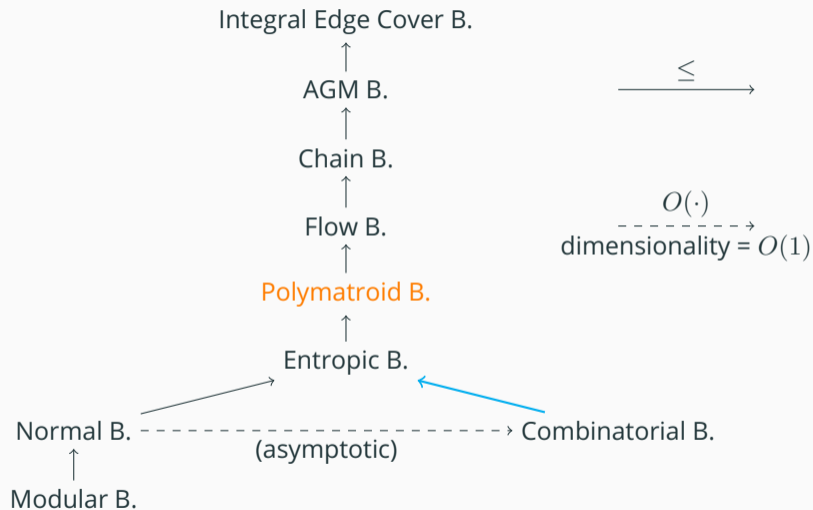
M_n : Modular $h(X) = \sum_{x \in X} h(x)$

A Collection of Optimization Problems

Many bounds can be formulated with two parameters

$$\sup\{h(V) \mid h \in P \cap \text{DC}\}$$

- DC = set of constraints $h(Y|X) \leq \log N$, one for degree constraint (X, Y, N) .
- P is a member in the aforementioned hierarchy of set functions
 - $P = M_n$: *modular bound*
 - $P = N_n$: *normal bound*
 - $P = \bar{\Gamma}_n^*$: *entropic bound*
 - $P = \Gamma_n$: *polymatroid bound*



Define $\mathbf{c} = (c_X)_{X \subseteq V}$, where $c_X = -1_{X=V}$, then – because C is linear –

$$\sup\{h(V) \mid h \in C \cap \bar{\Gamma}_n^*\} = \inf\{\langle \mathbf{c}, \mathbf{h} \rangle \mid \mathbf{A}\mathbf{h} \leq \mathbf{b} \wedge \mathbf{h} \in \bar{\Gamma}_n^*\}$$

Lagrangian: $\mathcal{L}(\boldsymbol{\delta}) = \inf_{\mathbf{h} \in \bar{\Gamma}_n^*} \langle \mathbf{c}, \mathbf{h} \rangle + \langle \mathbf{A}\mathbf{h} - \mathbf{b}, \boldsymbol{\delta} \rangle = -\langle \mathbf{b}, \boldsymbol{\delta} \rangle + \inf_{\mathbf{h} \in \bar{\Gamma}_n^*} \langle \mathbf{c} + \mathbf{A}^\top \boldsymbol{\delta}, \mathbf{h} \rangle$

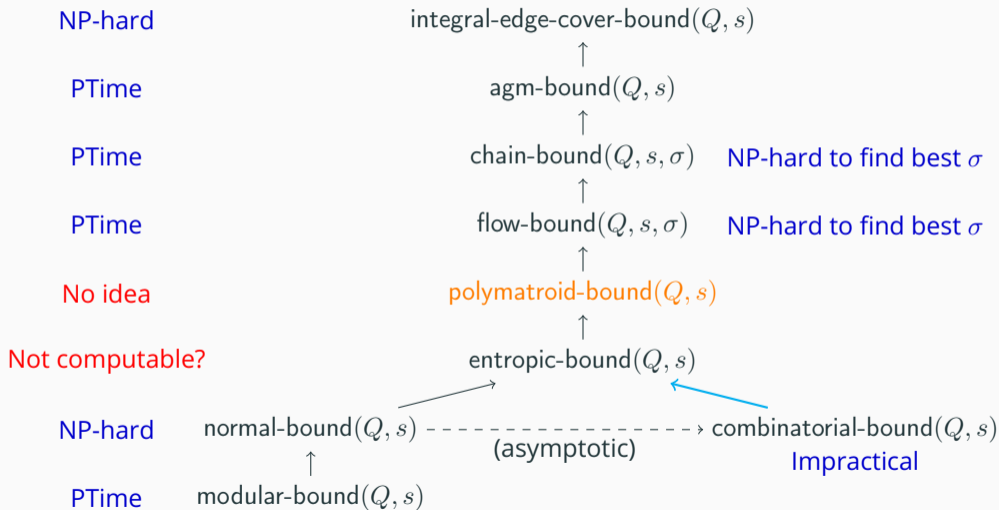
Lagrangian dual problem $(\bar{\Gamma}_n^*)^*$ denotes the dual cone of $\bar{\Gamma}_n^*$

$$\sup\{\mathcal{L}(\boldsymbol{\delta}) \mid \boldsymbol{\delta} \geq \mathbf{0}\} = \inf\{\langle \mathbf{b}, \boldsymbol{\delta} \rangle \mid \boldsymbol{\delta} \geq \mathbf{0} \wedge \mathbf{c} + \mathbf{A}^\top \boldsymbol{\delta} \in (\bar{\Gamma}_n^*)^*\}$$

Checking whether $\boldsymbol{\delta} \geq \mathbf{0}$ is dual feasible is equivalent to verifying whether

$$h(V) \leq \sum_{(X,Y) \in \text{DC}} \delta_{Y|X} h(Y|X) \quad \forall h \in \bar{\Gamma}_n^*$$

2. Computational Complexity of Polymatroid Bound



3. Parameterized Complexity of Polymatroid Bound

The polymatroid bound is computable in PTIME for some classes of inputs:

- Acyclic degree constraints
- Simple degree constraints
- Degree constraints with bounded SCCs

4. Tightness of Various Bounds

- For which class of queries does adding **conditional independence** constraints improves the bound?
- How close can we get to the entropic bound? (Ignore the combinatorial bound)
- How close we can get to combinatorial bound, under which condition?

Theorem

Except for the combinatorial \rightarrow entropic edge, there is an asymptotic gap in between every adjacent bound (connected by \rightarrow in the hierarchy), even if the number of degree constraints is fixed.

4. Tightness of Various Bounds

The following is *unsatisfactory*:

Proposition (ANS 2017)

For any $\epsilon > 0$, there exists a scale-factor k such that, if all DCs (X, Y, N) are scaled up into (X, Y, N^k) then

$$\text{entropic-bound} = (1 - \epsilon) \log \text{combinatorial-bound}$$

(Made use of **group-characterizable entropic functions**)

Outline

The Query Optimization (and Evaluation) Problem

Cardinality Bounds and Worst-Case Optimal Joins

The Bound Hierarchy Under Degree Constraints

Research Questions

References

References (and references thereof)

- N. Alon, On the number of subgraphs of prescribed type of graphs with a given number of edges, *Israel Journal of Mathematics* 38 (1981), 116–130.
- CGFS 1986 Fan R. K. Chung, Ronald L. Graham, Peter Frankl, James B. Shearer: Some intersection theorems for ordered sets and graphs. *J. Comb. Theory, Ser. A* 43(1): 23-37 (1986)
- Ehud Friedgut, Jeff Kahn, On the number of copies of one hypergraph in another. *Israel J. Math.* 105 (1998), 251–256.
- Ehud Friedgut: Hypergraphs, Entropy, and Inequalities. *Am. Math. Mon.* 111(9): 749-760 (2004)
- AGM 08 Albert Atserias, Martin Grohe, Dániel Marx: Size Bounds and Query Plans for Relational Joins. *FOCS 2008*: 739-748
- GLVV 12 Georg Gottlob, Stephanie Tien Lee, Gregory Valiant, Paul Valiant: Size and Treewidth Bounds for Conjunctive Queries. *J. ACM* 59(3): 16:1-16:35 (2012)
- ANS 16 Mahmoud Abo Khamis, Hung Q. Ngo, Dan Suciu: Computing Join Queries with Functional Dependencies. *PODS 2016*: 327-342
- ANS 17 Mahmoud Abo Khamis, Hung Q. Ngo, Dan Suciu: What Do Shannon-type Inequalities, Submodular Width, and Disjunctive Datalog Have to Do with One Another? *PODS 2017*: 429-444
- AKNS1 20 Mahmoud Abo Khamis, Phokion G. Kolaitis, Hung Q. Ngo, Dan Suciu: Decision Problems in Information Theory. *ICALP 2020*: 106:1-106:20
- AKNS2 20 Mahmoud Abo Khamis, Phokion G. Kolaitis, Hung Q. Ngo, Dan Suciu: Bag Query Containment and Information Theory. *PODS 2020*: 95-112
- Hung Q. Ngo: On an Information Theoretic Approach to Cardinality Estimation (Invited Talk). *ICDT 2022*: 1:1-1:21
- Dan Suciu: Applications of Information Inequalities to Database Theory Problems. *LICS 2023*: 1-30
- IMNPS 2023 Sungjin Im, Ben Moseley, Hung Q. Ngo, Kirk Pruhs, and Alireza Samadian: New Efficiently Computable Size Bounds for Joins. 2023.

Thank You!

Outline

Appendix

Outline

Appendix

Some known results on the parameterized complexity of the polymatroid bound

Frequency Moment Constraints

Histograms

3. Parameterized Complexity of Polymatroid Bound

- A set of DCs are **simple** if $|X| \leq 1$ for all DCs (X, Y, N)
- A set of DCs are **acyclic** if the constraint dependency graph is acyclic
 - *Constraint dependency graph*: for every degree constraint (X, Y, D) , add edges $x \rightarrow y$ for all $x \in X, y \in Y$.

Proposition

There is a polymatroid-bound-preserving reduction from an arbitrary instance to an instance where the degree constraints are a union of two sets: $DC = DC_a \cup DC_s$, where DC_s is simple, and DC_a is acyclic. Furthermore, for every non-simple DC $(X, Y, N) \in DC_a$, we have $|X| = 2$ and $|Y| \leq 3$.

Proposition

There is a poly-time algorithm computing a variable ordering σ_0 so that:

- If all input DCs are simple, then $\text{flow-bound}(Q, s, \sigma_0) = \text{polymatroid-bound}(Q, s)$*
- If all input DCs are acyclic, then $\text{flow-bound}(Q, s, \sigma_0) = \text{polymatroid-bound}(Q, s)$*

Proposition

If all DCs are acyclic, then there is a σ_0 for which

$$\text{chain-bound}(Q, s, \sigma_0) = \text{modular-bound}(Q, s).$$

All bounds in between collapse, and are all $\Theta(\text{combinatorial bound})$ in data-complexity.

Proposition

If all DCs are simple, then there is a σ_0 for which

$$\text{flow-bound}(Q, s, \sigma_0) = \text{normal-bound}(Q, s).$$

All bounds in between collapse, and are all $\Theta(\text{combinatorial bound})$ in data-complexity.

Outline

Appendix

Some known results on the parameterized complexity of the polymatroid bound

Frequency Moment Constraints

Histograms

Relation $R(\text{actor}, \text{movie}, \text{role})$, imagine a **frequency vector** $\mathbf{d}_{\text{actor}}$ (billions of entries)

| actor | movie | role |
|-------|-------|------|
| alice | | |
| bob | | |
| bob | | |
| bob | | |
| bob | | |
| carol | | |
| carol | | |

$$d_{\text{actor}}(\text{alice}) = 1$$

$$d_{\text{actor}}(\text{bob}) = 4$$

$$d_{\text{actor}}(\text{carol}) = 2$$

$$d_{\text{actor}}(v) = 0 \quad v \notin \{\text{alice}, \text{bob}, \text{carol}\}$$

The (very small) profile $s(\mathbf{D})$ contains

- $\|\mathbf{d}_{\text{actor}}\|_{\infty} = 4$ (**heaviest** frequency)
- $\|\mathbf{d}_{\text{actor}}\|_0 = 3$ (**distinct** counts)
- $\|\mathbf{d}_{\text{actor}}\|_1 = 7 = |R|$

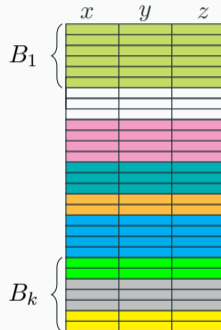
Similarly, we may have $\|\mathbf{d}_{\text{movie}}\|_p$, $\|\mathbf{d}_{\text{actor}, \text{movie}}\|_p$, $\|\mathbf{d}_{\text{role}}\|_p$, etc.

- Partition: $\text{Dom}(x) = B_1 \cup B_2 \cup \dots \cup B_k$ (*x-histogram*)
 - Typically about $k \approx 200$ buckets (e.g., MS SQL Server)
 - Top 10 heavy hitters are in their own singleton buckets
 - For the rest, equi-depth
- $\mathcal{B} := \{B_1, \dots, B_k\}$
- Give rise to per-bucket frequency vectors:

$$d_{Y|X \in B}(\mathbf{x}) := \begin{cases} |\pi_Y \sigma_{X=\mathbf{x}} R| & x \in B \\ 0 & \text{o.w.} \end{cases}$$

- HFM-Constraints $(\mathcal{B}, X, Y, c, \ell, R)$

$$F_\ell(d_{Y|X \in B}) \leq c_B \quad B \in \mathcal{B}$$



Outline

Appendix

Some known results on the parameterized complexity of the polymatroid bound

Frequency Moment Constraints

Histograms

Histograms are More Intricate

Main Question

How to turn $s(\mathbf{D})$ into constraints on polymatroids h ?

Answer sketch:

Theorem (KNN 2022, WIP)

Let Q be a conjunctive query and \mathcal{C} be a given set of simple HFM-constraints, then for any database \mathbf{D} satisfying \mathcal{C} , we have

$$\begin{aligned} \sup_{\mathbf{D} \models \mathcal{C}} \log |Q(\mathbf{D})| &\leq \max\{h(\mathbf{V}) \mid h \in \bar{\Gamma}_{n+m}^*, (h, \mathbf{P}, \mathcal{C}) \in \text{HC}\} \quad (\text{entropic bound}) \\ &\leq \max\{h(\mathbf{V}) \mid h \in \Gamma_{n+m}, (h, \mathbf{P}, \mathcal{C}) \in \text{HC}\} \quad (\text{polymatroid bound}) \end{aligned}$$

Example: $R(x, y) \wedge S(y, z)$

With a Simplifying Assumption

- Let \mathbf{f}_y and \mathbf{g}_y be the y -frequency vectors in R and S , respectively
- Assume **the same** partition $\text{dom}(y) = B_1 \cup \dots \cup B_k$ **on both** sides
- Suppose $s(\mathbf{D})$ contains the following statistics:

$$r_B := \|\mathbf{f}_{y \in B}\|_\infty \quad d_B := \|\mathbf{f}_{y \in B}\|_0 \quad s_B := \|\mathbf{g}_{y \in B}\|_\infty \quad B \in \mathcal{B}$$

- r_B = maximum number of x per y in bucket B
- d_B = number of distinct y 's in B
- s_B = maximum number of z per y in bucket B

Question: What are the constraints C on h ?

Turning $s(D)$ into Constraints on h

- Consider the uniform distribution on (X, Y, Z) chosen from $R(x, y) \wedge S(y, z)$
- Let $J \in \mathcal{B}$ be a categorical random variable, where $J = B$ iff $Y \in B$

$$p_B := \text{Prob}[J = B]$$

- Then,

$$h(J | Y) = 0 \tag{1}$$

$$h(Y | J = B) \leq \log d_B = \log \|\mathbf{f}_{y \in B}\|_0 \tag{2}$$

$$h(X | J = B) \leq \log r_B = \log \|\mathbf{f}_{y \in B}\|_\infty \tag{3}$$

$$h(Z | J = B) \leq \log s_B = \log \|\mathbf{g}_{y \in B}\|_\infty \tag{4}$$

$$h(J) = -\langle \mathbf{p}, \log \mathbf{p} \rangle \quad \|\mathbf{p}\|_1 = 1 \quad \mathbf{p} \geq \mathbf{0} \tag{5}$$

Simplifying the constraints

$$h(J | Y) = 0 \quad (6)$$

$$h(Y | J) = \sum_B h(Y | J = B) \cdot p_B \leq \langle \log \mathbf{d}, \mathbf{p} \rangle \quad (7)$$

$$h(X | J) = \sum_B h(X | J = B) \cdot p_B \leq \langle \log \mathbf{r}, \mathbf{p} \rangle \quad (8)$$

$$h(Z | J) = \sum_B h(Z | J = B) \cdot p_B \leq \langle \log \mathbf{s}, \mathbf{p} \rangle \quad (9)$$

$$h(J) = -\langle \mathbf{p}, \log \mathbf{p} \rangle \quad (10)$$

$$\|\mathbf{p}\|_1 = 1 \quad (11)$$

$$\mathbf{p} \geq \mathbf{0} \quad (12)$$

$$\max \quad h(XYZ) \quad (13)$$

$$\text{s.t.} \quad h(Y | J) \leq \langle \mathbf{p}, \lg \mathbf{d} \rangle \quad (14)$$

$$h(X | J) \leq \langle \mathbf{p}, \lg \mathbf{r} \rangle \quad (15)$$

$$h(Z | J) \leq \langle \mathbf{p}, \lg \mathbf{s} \rangle \quad (16)$$

$$h \in \Gamma_4 \quad \text{join distribution on } (X, Y, Z, J) \quad (17)$$

$$\mathbf{p} \geq 0, \quad (18)$$

$$h(J) = -\langle \mathbf{p}, \lg \mathbf{p} \rangle \quad (19)$$

$$h(J | Y) = 0 \quad (20)$$

$$\|\mathbf{p}\|_1 = 1. \quad (21)$$

Sanity Check: Estimator Makes Sense Combinatorially!

$$\lg |Q| = h(XYZ) \quad (22)$$

$$\text{(since } H[JY] = H[Y]) = H[XYZJ] = H[XYZ|J] + H[J] \quad (23)$$

$$\text{(since } H \in \Gamma_4) \leq H[X|J] + H[Y|J] + H[Z|J] + H[J] \quad (24)$$

$$\leq \sum_{j \in \mathcal{B}} (\lg r_B + \lg d_B + \lg s_B - \lg p_B) \cdot p_B \quad (25)$$

$$= \sum_{j \in \mathcal{B}} (\lg(r_B d_B s_B / p_B)) \cdot p_B \quad (26)$$

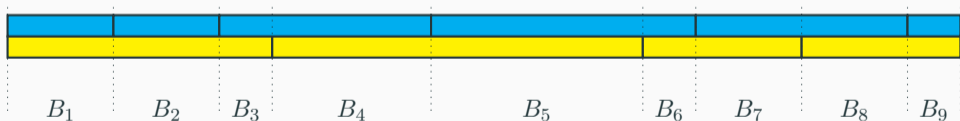
$$\text{(Jensen)} \leq \lg \left(\sum_B r_B d_B s_B \right). \quad (27)$$

$$|Q| \leq \sum_B r_B d_B s_B \quad (28)$$

Example – Removing the Simplifying Assumption

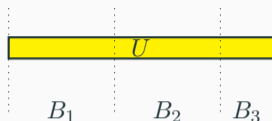
$$R(x, y) \wedge S(y, z)$$

- $\text{dom}(y) = U_1 \cup \dots \cup U_k$ on the R -side
- $\text{dom}(y) = V_1 \cup \dots \cup V_\ell$ on the S -side



Back to the “boundary aligned” model that $\text{dom}(y) = B_1 \cup B_2 \cup \dots \cup B_m$

Example – Removing the Simplifying Assumption



- For $q \in \{0, 1\}$ every constraint $\|\mathbf{f}_U^y\|_q$ is broken up into

$$\|\mathbf{f}_U^y\|_q = \|\mathbf{f}_{B_1}^y\|_q + \|\mathbf{f}_{B_2}^y\|_q + \|\mathbf{f}_{B_3}^y\|_q$$

where $\|\mathbf{f}_{B_j}^y\|_q$ are *new variables* in the optimization problem

- For $q = \infty$, set $\|\mathbf{f}_{B_i}^y\|_\infty \leq \|\mathbf{f}_{B_i}^y\|_1$ (or Hölder-type)

$$\max\{\|\mathbf{f}_{B_1}^y\|_\infty, \|\mathbf{f}_{B_2}^y\|_\infty, \|\mathbf{f}_{B_3}^y\|_\infty\} \leq \|\mathbf{f}_U^y\|_\infty$$

Under histogrammed FM constraints

Repeat the 4 (classes of) questions

- Computability
- Complexity
- Parameterized Complexity
- Tightness