Updated 2/25/2022

Topic 2: Complexity of Query Evaluation Unit 1: Conjunctive Queries (continued) Lecture 12

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

https://northeastern-datalab.github.io/cs7240/sp22/

2/25/2022

Pre-class conversations

- Recapitulation of Datalog & Query Equivalence
- Suggestion: Scribes with 2 iterations
- today:
 - Query equivalence of CQs & homomorphisms

Topic 1: Data Models and Query Languages

- Lecture 1 (Tue 1/18): Course introduction / 1 SQL / PostgreSQL setup / SQL Activities
- Lecture 2 (Fri 1/21): 1 SQL
- Lecture 3 (Tue 1/25): 1 SQL
- Lecture 4 (Fri 1/28): 1 SQL, 2 Logic & Relational Calculus
- Lecture 5 (Tue 2/1): 2 Logic & Relational Calculus
- Lecture 6 (Fri 2/4): 3 Relational Algebra & Codd's Theorem
- Lecture 7 (Tue 2/8): 3 Relational Algebra & Codd's Theorem
- Lecture 8 (Fri 2/11): 3 Relational Algebra & Codd's Theorem / 4 Datalog & Recursion
- Lecture 9 (Tue 2/15): 4 Datalog & Recursion
- Lecture 10 (Tue 2/18): 4 Datalog & Recursion
- Lecture 11 (Tue 2/22): 4 Datalog & Recursion

Pointers to relevant concepts & supplementary material:

- Unit 1. SQL: [SAMS'12], [CS 3200], [Cow'03] Ch3 & Ch5, [Complete'08] Ch6, [Silberschatz+'20] Ch3.8
- Unit 2. Logic & Relational Calculus: First-Order Logic (FOL), relational calculus (RC): [Barland+'08] 4.1.2 & 4.2.1 & 4.4, [Genesereth+] Ch6, [Halpern+'01], [Cow'03] Ch4.3 & 4.4, [Elmasri, Navathe'15] Ch8.6 & Ch8.7, [Silberschatz+'20] Ch27.1 & Ch27.2, [Alice'95] Ch3.1-3.3 & Ch4.2 & Ch4.4 & Ch5.3-5.4, [Barker-Plummer+'11] Ch11
- Unit 3. Relational Algebra & Codd's Theorem: Relational Algebra (RA), Codd's theorem: [Cow'03] Ch4.2,
 [Complete'08] Ch2.4 & Ch5.1-5.2, [Elmasri, Navathe'15] Ch8, [Silberschatz+'20] Ch2.6, [Alice'95] Ch4.4 & Ch5.4
- Unit 4. Datalog & Recursion: Datalog, recursion, Stratified Datalog with negation, Datalog evaluation strategies, Stable Model semantics, Answer Set Programming (ASP): [Complete'08] Ch5.3, [Cow'03] Ch 24, [Koutris'19] L9 & L10, [G., Suciu'10]
- Unit 5. Alternative Data Models: NoSQL: [Hellerstein, Stonebraker'05], [Sadalage, Fowler'12], [Harrison'16]

Topic 2: Complexity of Query Evaluation & Reverse Data Management

- CONTINUED Lecture 11 (Tue 2/22): 1 Conjunctive Queries
- Lecture 12 (Fri 2/25): Conjunctive Queries
- Lecture 13 (Tue 3/1): Beyond Conjunctive Queries
- Lecture 14 (Fri 3/4): Provenance
- Lecture 15 (Tue 3/8): Provenance, Reverse Data Management

Pointers to relevant concepts & supplementary material:

- Unit 1. Conjunctive Queries: Query evaluation of conjunctive queries (CQs), data vs. query complexity, homomorphisms, constraint satisfaction, query containment, query minimization, absorption: [Kolaitis, Vardi'00], [Vardi'00], [Kolaitis'16], [Koutris'19] L1 & L2
- **Unit 2. Beyond Conjunctive Queries**: unions of conjunctive queries, bag semantics, nested queries, tree pattern queries: [Kolaitis'16], [Tan+'14], [G.'11], [Martens'17]
- Unit 3. Provenance: [Buneman+02], [Green+07], [Cheney+09], [Green, Tannen'17], [Kepner+16], [Buneman, Tan'18]
- Unit 4. Reverse Data Management: update propagation, resilience: [Buneman+02], [Kimelfeld+12], [Freire+15]

Topic 3: Efficient Query Evaluation & Factorized Representations

- Lecture 16 (Fri 3/11): Acyclic Queries
- Spring break
- Lecture 17 (Tue 3/22): Acyclic Queries
- Lecture 18 (Fri 3/25): Cyclic Queries
- Lecture 19 (Tue 3/29): Cyclic Queries
- Lecture 20 (Fri 4/1): Factorized Representations
- Lecture 21 (Tue 4/5): Factorized Representations
- Lecture 22 (Fri 4/8): Top-k & Optimization Problems
- Lecture 23 (Tue 4/12): Top-k & Optimization Problems

Pointers to relevant concepts & supplementary material:

- **Unit 1. Acyclic Queries**: query hypergraph, Yannakakis algorithm, GYO reduction, dynamic programming, algebraic semirings, [Alice] Ch6.4, [Koutris'19] L4, enumeration, ranked enumeration:[Tziavelis+'20]
- Unit 2. Cyclic Queries: tree & hypertree decomposition, query widths, fractional hypertree width, AGM bound, worst-case optimal join algorithms, optimal algorithms, submodular width and multiple decompositions: [AGM'13], [NPRR'18], [KNR'17], [KNS'17]
- Unit 3. Factorized Representations: normalization, factorized databases [Olteanu, Schleich'16]
- Unit 4. Top-k & Optimization Problems: shortest paths, dynamic programming (DP), Yannakakis, semirings, rankings, top-k: [Roughgarden'10], [Ilyas+08], [Rahul, Tao'19], ranked enumeration [Tziavelis+'19]

Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
 - CQ equivalence and containment
 - Graph homomorphisms
 - Homomorphism beyond graphs
 - CQ containment
 - CQ minimization
- T2-2: Equivalence Beyond CQs
 - Union of CQs, and inequalities
 - Union of CQs equivalence under bag semantics
 - Tree pattern queries
 - Nested queries

Complexity of the Query Evaluation Problem

- The Query Evaluation Problem for Relational Calculus (RC):
 - Given a RC formula φ and a database instance D, find $\varphi^{adom}(D)$.
- Theorem: The Query Evaluation Problem for Relational Calculus is PSPACE-complete.
 - PSPACE: decision problems, can be solved using an amount of memory that is polynomial in the input length (~ in polynomial amount of space).
 - PSPACE-complete: PSPACE + every other PSPACE problem can be transformed to it in polynomial time (PSPACE-hard)
- Proof: We need to show both
 - This problem is in PSPACE.
 - This problem is PSPACE-hard. (We only focus on this task for Boolean RC queries)

Complexity of the Query Evaluation Problem

- Theorem: The Query Evaluation Problem for Boolean RC is PSPACE-hard.
- Reduction uses QBF (Quantified Boolean Formulas):
 - Given QBF $\forall x_1 \exists x_2 \dots \forall x_k \psi$, is it true or false
 - (notice every variable is <u>quantified = bound</u> at beginning of <u>sentence</u>; no free variables)
- Proof shows that QBF \preccurlyeq Query Evaluation for Relational Calculus
 - Given QBF $\forall x_1 \exists x_2 \dots \forall x_k \psi$,
 - Let V and P be two unary relations and D be the database instance with V(0), V(1), P(1)
 - Obtain ψ^* from ψ by replacing every occurrence of x_i by $P(x_i)$, and $\neg x_i$ by $\neg P(x_i)$
 - Then the following statements are equivalent:
 - $\forall x_1 \exists x_2 \dots \forall x_k \psi$ is true
 - $\forall x_1 [V(x_1) \rightarrow \exists x_2 [V(x_2) \land ... \forall x_k [V(x_k) \rightarrow \psi^*]]...]$ is true on D

Sublanguages of Relational Calculus

 Question: Are there interesting sublanguages of relational calculus for which the Query Containment Problem and the Query Evaluation Problem are "easier" than the full relational calculus?

- Answer:
 - Yes, the language of Conjunctive Queries (CQs) is such a sublanguage.
 - Moreover, conjunctive queries are the most frequently asked queries against relational databases.

Based on Phokion Kolaitis' "Logic and Databases" series at Simons Institute, 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Conjunctive Queries (CQs)

- Definition:
 - A CQ is a query expressible by a RC formula in prenex normal form built from atomic formulas $R(y_1,...,y_n)$, and ∧ and ∃ only.

{ $(x_1,...,x_k): \exists z_1 ... \exists z_m \phi(x_1,...,x_k, z_1,...,z_k)$ },

- where $\phi(x_1, ..., x_k, z_1, ..., z_k)$ is a conjunction of atomic formulas of the form $R(y_1, ..., y_m)$.
- <u>Prenex formula</u>: prefix (quantifiers & bound variables), then quantifier-free part
- Equivalently, a CQ is a query expressible by a RA expression of the form
 - $\pi_X(\sigma_{\Theta}(R_1 \times ... \times R_n))$, where
 - Θ is a conjunction of equality atomic formulas (equijoin).
- Equivalently, a CQ is a query expressible by an SQL expression of the form
 - SELECT <list of attributes>
 - FROM <list of relation names>

Based on Phokion Koraitis "Logic and Databases" series at Simons Institute, 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Conjunctive Queries (CQs)

- Definition:
 - A CQ is a query expressible by a RC formula in prenex normal form built from atomic formulas $R(y_1,...,y_n)$, and ∧ and ∃ only.

{ $(x_1,...,x_k): \exists z_1 ... \exists z_m \phi(x_1,...,x_k, z_1,...,z_k)$ },

- where $\phi(x_1, ..., x_k, z_1, ..., z_k)$ is a conjunction of atomic formulas of the form $R(y_1, ..., y_m)$.
- Equivalently, a CQ can be written as a logic-programming rule:

 $Q(x_1,...,x_k) := R_1(u_1), ..., R_n(u_n)$, where

- Each variable x_i occurs in the right-hand side of the rule.
- Each **u**_i is a tuple of variables (not necessarily distinct)
- The variables occurring in the right-hand side (the body), but not in the left-hand side (the head) of the rule are existentially quantified (but the quantifiers are not displayed).

Conjunctive Queries (CQs)

- Every natural join is a conjunctive query with no existentially quantified variables
- Example: Given P(A,B,C), R(B,C,D)
 - P \bowtie R = {(x,y,z,w): P(x,y,z) \land R(y,z,w)}
 - q(x,y,z,w) := P(x,y,z), R(y,z,w)

(no variables are existentially quantified)

- SELECT P.A, P.B, P.C, R.D
 FROM P, R
 WHERE P.B = R.B AND P.C = R.C
- Conjunctive queries are also known as SPJ-queries (SELECT-PROJECT-JOIN queries)





• Return paths of Length 2: (binary output)

RC: ? RA: ? Datalog: ?

E(from, to)

Is there a path

of length 2

Ε

2

×

 $\times \neq \gamma$

1



Return paths of Length 2: (binary output)

RC: $\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$

RA:

Datalog:







- Return paths of Length 2: (binary output)
 - RC: $\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$
 - RA: $\pi_{\$1,\$4}(\sigma_{\$2=\$3}(E \times E))$ unnamed perspective Datalog: ?





• Return paths of Length 2: (binary output)

RC:
$$\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$$

RA:
$$\pi_{\$1,\$4}(\sigma_{\$2=\$3}(E \times E))$$
 unnamed perspective

Datalog: Q(x,y) := E(x,z), E(z,y)

• Is there a cycle of Length 3: (Boolean query)

RC: Datalog:

E(from, to)



• Return paths of Length 2: (binary output)

RC:
$$\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$$

RA: $\pi_{\$1,\$4}(\sigma_{\$2=\$3}(E \times E))$ unnamed perspective

Datalog: Q(x,y) := E(x,z), E(z,y)

• Is there a cycle of Length 3: (Boolean query)

RC:
$$\exists x \exists y \exists z [E(x, y) \land E(y, z) \land E(z, x)]$$

Datalog:







• Return paths of Length 2: (binary output)

RC:
$$\{(x, y) \mid \exists z [E(x, z) \land E(z, y)]\}$$

RA:
$$\pi_{\$1,\$4}(\sigma_{\$2=\$3}(E \times E))$$
 unnamed perspective

Datalog:
$$Q(x,y) := E(x,z), E(z,y)$$

• Is there a cycle of Length 3: (Boolean query)

RC:
$$\exists x \exists y \exists z [E(x,y) \land E(y,z) \land E(z,x)]$$

Datalog: Q := E(x, y), E(y, z), E(z, x)

Vardi's Taxonomy of the Query Evaluation Problem

M.Y Vardi, "The Complexity of Relational Query Languages", 1982

- Definition: Let L be a database query language.
 - The combined complexity of L is the decision problem:
 - given an L-sentence and a database instance D, is φ true on D?
 - In symbols, does D ⊧ φ (does D satisfy φ)?



- given a database instance D, does $D \models \phi$?
- The query complexity of L is the family of the following decision problems P_D, where D is a database instance:
 - given an L-sentence φ , does D $\models \varphi$?

Vardi's Taxonomy of the Query Evaluation Problem

Vardi's "empirical" discovery:

- For most query languages L:
 - The data complexity of L is of lower complexity than both the combined complexity of L and the query complexity of L.
 - The query complexity of L can be as hard as the combined complexity of L.

Taxonomy of the Query Evaluation Problem for Relational Calculus

Complexity Classes



The Query Evaluation Problem for Relational Calculus

Problem	Complexity
Combined Complexity	PSPACE-complete
Query Complexity	• in PSPACE
	 can be PSPACE- complete
Data Complexity	In Logspace

Summary

- Relational Algebra and Relational Calculus have "essentially" the same expressive power.
- The Query Equivalence Problem for Relational Calculus is undecidable.
 - Therefore also the Query Containment Problem
- The Query Evaluation Problem for Relational Calculus:
 - Data Complexity is in LOGSPACE
 - Combined Complexity is PSPACE-complete
 - Query Complexity is PSPACE-complete.

Based on Phokion Kolaitis' "Logic and Databases" series at Simons Institute, 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
 - CQ equivalence and containment
 - Graph homomorphisms
 - Homomorphism beyond graphs
 - CQ containment
 - CQ minimization
- T2-2: Equivalence Beyond CQs
 - Union of CQs, and inequalities
 - Union of CQs equivalence under bag semantics
 - Tree pattern queries
 - Nested queries



Source: <u>https://en.wikipedia.org/wiki/Bijection, injection and surjection</u>



Source: <u>https://en.wikipedia.org/wiki/Bijection, injection and surjection</u>





maps each <u>argument</u> (element from its <u>domain</u>) to exactly one <u>image</u> (element in its <u>codomain</u>) $\forall x \in X, \exists ! y \in Y[y = f(x)]$ }

("one-to-one"): each element of the codomain is mapped to by <u>at most one</u> element of the domain (i.e. distinct elements of the domain map to distinct elements in the codomain)

$$\forall x, x' \in X. [x \neq x' \Rightarrow f(x) \neq f(x')] \\ \forall x, x' \in X. [f(x) = f(x') \Rightarrow x = x']$$

("onto"): each element of the codomain is mapped to by <u>at least one</u> element of the domain (i.e. the image and the codomain of the function are equal) $\forall y \in Y, \exists x \in X[y = f(x)]$



 $\forall y \in Y, \exists ! x \in X[y = f(x)] \}$

Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/

Source: https://en.wikipedia.org/wiki/Bijection, injection and surjection









not a mapping (or function)!

≻



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

≻€ ►



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it

injective & surjective = bijection



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it

injective & surjective = bijection

neighter



not a mapping (or function)!

injective function (or one-to-one): maps distinct elements of its domain to <u>distinct elements of its codomain</u>

surjective (or onto): every element y in the codomain Y of f has at least one element x in the domain that maps to it

injective & surjective = bijection

neighter

not even a mapping!



Bijection, Injection, and Surjection



Sources: http://mathonline.wikidot.com/injections-surjections-and-bijections,

https://www.intechopen.com/books/protein-interactions/relating-protein-structure-and-function-through-a-bijection-and-its-implications-on-protein-structur,
Bijection, Injection, and Surjection



not injective

Not a function

injective + surjective

Sources: <u>https://www.mathsisfun.com/sets/injective-surjective-bijective.html</u>, <u>https://twitter.com/jdhamkins/status/841318019397779456</u>, Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

not surjective

not surjective

We make a detour to Graph matching

• Finding a correspondence between the nodes and the edges of two graphs that satisfies some (more or less stringent) constraints



- A graph homomorphism *h* from graph $G(V_G, E_G)$ to $H(V_H, E_H)$, is a mapping from V_G to V_H such that $\{x, y\} \in E_G$ implies $\{h(x), h(y)\} \in E_H$
 - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked







- A graph homomorphism *h* from graph $G(V_G, E_G)$ to $H(V_H, E_H)$, is a mapping from V_G to V_H such that $\{x, y\} \in E_G$ implies $\{h(x), h(y)\} \in E_H$
 - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



- A graph homomorphism *h* from graph $G(V_G, E_G)$ to $H(V_H, E_H)$, is a mapping from V_G to V_H such that $\{x, y\} \in E_G$ implies $\{h(x), h(y)\} \in E_H$
 - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked



Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



- A graph homomorphism h from graph $G(V_G, E_G)$ to $H(V_H, E_H)$, is a mapping from V_G to V_H such that $\{x, y\} \in E_G$ implies $\{h(x), h(y)\} \in E_H$
 - "edge-preserving": if two nodes in G are linked by an edge, then they are mapped to two nodes in H that are also linked





- Graphs $G(V_G, E_G)$ and $H(V_H, E_H)$ are isomorphic iff there is an invertible h from V_G to V_H s.t. $\{x, y\} \in E_G$ iff $\{h(u), h(v)\} \in E_H$
 - We need to find a one-to-one correspondence



Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



- Graphs $G(V_G, E_G)$ and $H(V_H, E_H)$ are isomorphic iff there is an invertible h from V_G to V_H s.t. $\{x, y\} \in E_G$ iff $\{h(u), h(v)\} \in E_H$
 - We need to find a one-to-one correspondence





- Graphs $G(V_G, E_G)$ and $H(V_H, E_H)$ are isomorphic iff there is an invertible h from V_G to V_H s.t. $\{x, y\} \in E_G$ iff $\{h(u), h(v)\} \in E_H$
 - We need to find a one-to-one correspondence



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



- Graphs $G(V_G, E_G)$ and $H(V_H, E_H)$ are isomorphic iff there is an invertible h from V_G to V_H s.t. $\{x, y\} \in E_G$ iff $\{h(u), h(v)\} \in E_H$
 - We need to find a one-to-one correspondence





Is there an isomorphismYes: $h: \{(1,a), (2,b), (3,d), (4,c), (5,e)\}$ from G to H?bijection = surjective and injective mapping

Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
 - CQ equivalence and containment
 - Graph homomorphisms
 - Homomorphism beyond graphs
 - CQ containment
 - CQ minimization
- T2-2: Equivalence Beyond CQs
 - Union of CQs, and inequalities
 - Union of CQs equivalence under bag semantics
 - Tree pattern queries
 - Nested queries

Graph Homomorphism beyond graphs

Definition : Let G and H be graphs. A homomorphism of G to H is a function $f: V(G) \rightarrow V(H)$ such that

 \mathcal{G}

 $(x,y)\in E(G)\Rightarrow (f(x),f(y))\in E(H).$

We sometimes write $G \rightarrow H$ (G \rightarrow H) if there is a homomorphism (no homomorphism) of G to H

Definition of a homomorphism naturally extends to:

- digraphs (directed graphs)
- edge-colored graphs
- relational systems
- constraint satisfaction problems (CSPs)











Based upon an example from Rick Brewster's Graph homomorphism tutorial, 2006 Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>





Can this assignment be extended to a homomorphism?

Based upon an example from Rick Brewster's Graph homomorphism tutorial, 2006 Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u> No, this assignment requires a loop on vertex 1 (in H)





Based upon an example from Rick Brewster's Graph homomorphism tutorial, 2006 Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



Definition: Let G and H be graphs. A homom. of G to H is a function $f: V(G) \rightarrow V(H)$ s.t. that





Based upon an example from Rick Brewster's Graph homomorphism tutorial, 2006 Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

a 1



Definition: Let G and H be graphs. A homom. of G to H is a function f: $V(G) \rightarrow V(H)$ s.t. that





Basically a partitioning problem!

The quotient set of the partition (set of equivalence classes of the partition) is a subgraph of H.







Some observations





When does $G \rightarrow K_3$ hold? ($K_3 = 3$ -clique = triangle)

Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/

Some observations When does $G \rightarrow K_3$ hold? ($K_3 = 3$ -clique = triangle) iff G is 3-colorable

When does $G \rightarrow K_d$ hold? ($K_d = d$ -clique)



Some observations When does $G \rightarrow K_3$ hold? ($K_3 = 3$ -clique = triangle) iff G is 3-colorable

- When does $G \rightarrow K_d$ hold? ($K_d = d$ -clique) iff G is d-colorable
- Thus homomorphisms generalize colorings: Notation: $G \rightarrow H$ is an H-coloring of G.

What is the complexity of testing for the existence of a homomorphism (in the size of G)? ?





When does $G \rightarrow K_3$ hold? ($K_3 = 3$ -clique = triangle) iff G is 3-colorable

When does $G \rightarrow K_d$ hold? ($K_d = d$ -clique) iff G is d-colorable

Some observations

Thus homomorphisms generalize colorings: Notation: $G \rightarrow H$ is an H-coloring of G.

What is the complexity of testing for the existence of a homomorphism (in the size of G)?



109





The complexity of H-coloring

H-coloring:

Let H be a fixed graph.

Instance: A graph G.

Question: Does G admit an H-coloring?

Repeated variable names



In sentences with multiple quantifiers, <u>distinct variables do not need</u> <u>to range over distinct objects</u>! (cp. homomorphism vs. isomorphism)



Repeated variable names



In sentences with multiple quantifiers, <u>distinct variables do not need</u> <u>to range over distinct objects</u>! (cp. homomorphism vs. isomorphism)







A more abstract (general) view on homomorphisms

Homomorphisms on Binary Structures

- Definition (Binary algebraic structure): A binary algebraic structure is a set together with a binary operation on it. This is denoted by an ordered pair (*S*,*) in which *S* is a set and * is a binary operation on *S*.
- Definition (homomorphism of binary structures): Let (S,*) and (S',∘) be binary structures. A homomorphism from (S,*) to (S',∘) is a map h: S → S' that satisfies, for all x, y in S:

 $h(x \star y) = h(x) \circ h(y)$

• We can denote it by $h: (S, \star) \longrightarrow (S', \circ)$.

• Let $h(x) = e^x$. Is h a homomorphism b/w two binary structures?



- Let $h(x) = e^x$. Is h a homomorphism b/w two binary structures?
 - Yes, from the real numbers with addition (\mathbb{R} ,+) to $h(x+y) = h(x) \cdot h(y)$
 - the positive real numbers with multiplication (\mathbb{R}^+, \cdot) $h:(\mathbb{R}, +) \to (\mathbb{R}^+, \cdot)$
 - It is even an isomorphism!

The exponential map exp : $\mathbb{R} \to \mathbb{R}^+$ defined by $\exp(x) = e^x$, where *e* is the base of the natural logarithm, is an isomorphism from $(\mathbb{R}, +)$ to (\mathbb{R}^+, \times) . Exp is a bijection since it has an inverse function (namely \log_e) and exp preserves the group operations since $e^{x+y} = e^x e^y$. In this example both the elements and the operations are different yet the two groups are isomorphic, that is, as groups they have identical structures.

• Let $g(x) = e^{ix}$. Is g also a homomorphism?

7

Paragraph screenshot from p.37 in 2004 - Dummit, Foote - Abstract algebra (book, 3rd ed). <u>https://www.wiley.com/en-us/Abstract+Algebra%2C+3rd+Edition-p-9780471433347</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

- Let $h(x) = e^x$. Is h a homomorphism b/w two binary structures?
 - Yes, from the real numbers with addition (\mathbb{R} ,+) to $h(x+y) = h(x) \cdot h(y)$
 - the positive real numbers with multiplication (\mathbb{R}^+, \cdot) $h:(\mathbb{R}, +) \to (\mathbb{R}^+, \cdot)$
 - It is even an isomorphism!

The exponential map exp : $\mathbb{R} \to \mathbb{R}^+$ defined by $\exp(x) = e^x$, where *e* is the base of the natural logarithm, is an isomorphism from $(\mathbb{R}, +)$ to (\mathbb{R}^+, \times) . Exp is a bijection since it has an inverse function (namely \log_e) and exp preserves the group operations since $e^{x+y} = e^x e^y$. In this example both the elements and the operations are different yet the two groups are isomorphic, that is, as groups they have identical structures.

- Let $g(x) = e^{ix}$. Is g also a homomorphism?
 - Yes, from the real numbers with addition (\mathbb{R} ,+) to
 - the unit circle in the complex plane with rotation



$$\begin{array}{ll} G = \mathbb{R} \text{ under } + & f: G \to H \\ H = \{ z \in \mathbb{C} : |z| = 1 \} & z \mapsto e^{ix} \\ = \text{Group under } \times & \text{Show } f(x + y) = f(x) \times f(y) \\ e^{i(x+y)} = e^{ix} \times e^{iy} \\ e^{ix+iy} = e^{ix} \times e^{iy} \\ e^{ix+iy} = e^{ix} \times e^{iy} \\ e^{ix} \times e^{iy} = e^{ix} \times e^{iy} \\ f(0) = f(2\pi) = 1, f(2\pi\pi) = 1 \\ f \text{ is not } 1-1 \end{array}$$

Source: Socratica. Homomorphisms, 2014: <u>https://www.youtube.com/watch?v=cYzp5IWqCsg</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



Source: 3blue1brown. Euler's formula with introductory group theory, 2017: <u>https://www.youtube.com/watch?v=mvmuCPvRoWQ</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Isomorphism

- **Definition**: A homomorphism of binary structures is called an isomorphism iff the corresponding map of sets is:
 - one-to-one (injective) and
 - onto (surjective).



Some homomorphisms



CQs (Var U Constants, Relations {R_i(x,y,z), ...})

Outline: T2-1/2: Query Evaluation & Query Equivalence

- T2-1: Conjunctive Queries (CQs)
 - CQ equivalence and containment
 - Graph homomorphisms
 - Homomorphism beyond graphs
 - CQ containment
 - CQ minimization
- T2-2: Equivalence Beyond CQs
 - Union of CQs, and inequalities
 - Union of CQs equivalence under bag semantics
 - Tree pattern queries
 - Nested queries
Query Containment

Two queries q_1 , q_2 are equivalent, denoted $q_1 \equiv q_2$, if for every database instance D, we have $q_1(D) = q_2(D)$. the answer (set of tuples) - returned by one is guaranteed to be identical to the other answer

Query q_1 is contained in query q_2 , denoted $q_1 \subseteq q_2$, if for every database instance D, we have $q_1(D) \subseteq q_2(D)$

Corollary

$$q_1 \equiv q_2$$
 is equivalent to $(q_1 \subseteq q_2 \text{ and } q_1 \supseteq q_2)$

If queries are Boolean, then query containment = logical implication: $q_1 \Leftrightarrow q_2$ is equivalent to

Query Containment

Two queries q_1 , q_2 are equivalent, denoted $q_1 \equiv q_2$, if for every database instance D, we have $q_1(D) = q_2(D)$. the answer (set of tuples) - returned by one is guaranteed to be identical to the other answer

Query q_1 is contained in query q_2 , denoted $q_1 \subseteq q_2$, if for every database instance D, we have $q_1(D) \subseteq q_2(D)$

Corollary

$$q_1 \equiv q_2$$
 is equivalent to $(q_1 \subseteq q_2 \text{ and } q_1 \supseteq q_2)$

If queries are Boolean, then query containment = logical implication: $q_1 \Leftrightarrow q_2$ is equivalent to $(q_1 \Rightarrow q_2 \text{ and } q_1 \leftarrow q_2)$



A homomorphism *h* from Boolean q_2 to q_1 is a function $h: var(q_2) \rightarrow var(q_1) \cup const(q_1)$ such that:

for every atom $R(x_1, x_2, ...)$ in q_2 , there is an atom $R(h(x_1), h(x_2), ...)$ in q_1

need to be same relation!

Example

 $q_1(x) := R(x,y), R(y,y), R(y,z)$ $q_2(s) := R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)$



Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>





A homomorphism *h* from Boolean q_2 to q_1 is a function $h: var(q_2) \rightarrow var(q_1) \cup const(q_1)$ such that:

for every atom $R(x_1, x_2, ...)$ in q_2 , there is an atom $R(h(x_1), h(x_2), ...)$ in q_1

need to be same relation!

Example

 $q_1(x) := R(x,y), R(y,y), R(y,z)$ $q_2(s) := R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)$





A homomorphism *h* from Boolean q_2 to q_1 is a function $h: var(q_2) \rightarrow var(q_1) \cup const(q_1)$ such that:

for every atom $R(x_1, x_2, ...)$ in q_2 , there is an atom $R(h(x_1), h(x_2), ...)$ in q_1

need to be same relation!

Example

 $q_1(x) := R(x,y), R(y,y), R(y,z)$ $q_2(s) := R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)$



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/





A homomorphism *h* from Boolean q_2 to q_1 is a function $h: var(q_2) \rightarrow var(q_1) \cup const(q_1)$ such that:

for every atom $R(x_1, x_2, ...)$ in q_2 , there is an atom $R(h(x_1), h(x_2), ...)$ in q_1

need to be same relation!

Example

Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



A homomorphism *h* from Boolean q_2 to q_1 is a function $h: var(q_2) \rightarrow var(q_1) \cup const(q_1)$ such that:

for every atom $R(x_1, x_2, ...)$ in q_2 , there is an atom $R(h(x_1), h(x_2), ...)$ in q_1

need to be same relation!

Example

 $\begin{aligned} & q_1(x) := R(x,y), \, R(y,y), \, R(y,z) \\ & q_2(s) := R(s,u), \, R(u,w), \, R(s,v), \, R(v,w), \, R(u,v), \, R(v,v) \end{aligned}$



Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>







A homomorphism *h* from Boolean q_2 to q_1 is a function $h: var(q_2) \rightarrow var(q_1) \cup const(q_1)$ such that:

for every atom $R(x_1, x_2, ...)$ in q_2 , there is an atom $R(h(x_1), h(x_2), ...)$ in q_1





Canonical database



Definition Canonical database

Given a conjunctive query q, the canonical database $D_c[q]$ is the database instance where each atom in q becomes a fact in the instance.

Example $q_1(x) := R(x,y), R(y,y), R(y,z)$ $D_c[q_1] = ?$

Definition Canonical database

Given a conjunctive query q, the canonical database $D_c[q]$ is the database instance where each atom in q becomes a fact in the instance.

Example

- $q_1(x) := R(x,y), R(y,y), R(y,z)$
- $D_{c}[q_{1}] = \{R('x', 'y'), R('y', 'y'), R('y', 'z')\}$
 - $\equiv \{R(\mathsf{a},\mathsf{b}), R(\mathsf{b},\mathsf{b}), R(\mathsf{b},\mathsf{c})\}$
 - $\equiv \{R(1,2),\,R(2,2),\,R(2,3)\}$

Just treat each variable as different constant 😳



AR ()_1_v).



THEOREM (Query Containment) Given two Boolean CQs q_1 , q_2 , the following statements are equivalent: 1) $q_1 \subseteq q_2$ $(q_1 \Rightarrow q_2)$ 2) There is a homomorphism $h_{2 \rightarrow 1}$ from q_2 to q_1 3) $q_2(D_C[q_1])$ is true

We will only look at 2) \Rightarrow 1)



Chandra, Merlin. "Optimal implementation of conjunctive queries in relational data bases." STOC 1977. <u>https://doi.org/10.1145/800105.803397</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

BORFAN

 $q = v \circ h$

q(x) = v(h(x))

- We show: If there is a homomorphism $h_{2 \rightarrow 1}$, then for any D: $q_1(D) \Rightarrow q_2(D)$
- 1. For $q_1(D)$ to hold, there is a valuation v s.t. $v(q_1) \in D$
- 2. We will show that the composition $g = v \circ h$ is a valuation for q_2

 - 2a. By definition of h, for every $R(x_1, x_2, ...)$ in q_2 , $R(h(x_1), h(x_2), ...)$ in q_1 2b. By definition of v, for every $R(x_1, x_2, ...)$ in q_2 , $R(v(h(x_1)), v(h(x_2)), ...)$ in D



- We show: If there is a homomorphism $h_{2 \rightarrow 1}$, then for any D: $q_1(D) \Rightarrow q_2(D)$
- 1. For $q_1(D)$ to hold, there is a valuation v s.t. $v(q_1) \in D$
- 2. We will show that the composition $g = v \circ h$ is a valuation for q_2
 - 2a. By definition of h, for every $R(x_1, x_2, ...)$ in q_2 , $R(h(x_1), h(x_2), ...)$ in q_1
 - 2b. By definition of v, for every $R(x_1, x_2, ...)$ in q_2 , $R(v(h(x_1)), v(h(x_2)), ...)$ in D

 $q = v \circ h$

q(x) = v(h(x))

<u>Example</u>

 $q_1() := R(x,y), R(y,y), R(y,z)$ $q_2() := R(s,u), R(u,w), R(s,v), R(v,w), R(u,v)$



- We show: If there is a homomorphism $h_{2 \rightarrow 1}$, then for any D: $q_1(D) \Rightarrow q_2(D)$
- 1. For $q_1(D)$ to hold, there is a valuation v s.t. $v(q_1) \in D$
- 2. We will show that the composition $g = v \circ h$ is a valuation for q_2
 - 2a. By definition of h, for every $R(x_1, x_2, ...)$ in q_2 , $R(h(x_1), h(x_2), ...)$ in q_1
 - 2b. By definition of v, for every $R(x_1, x_2, ...)$ in q_2 , $R(v(h(x_1)), v(h(x_2)), ...)$ in D

 $q = v \circ h$

q(x) = v(h(x))

Example



- We show: If there is a homomorphism $h_{2 \rightarrow 1}$, then for any D: $q_1(D) \Rightarrow q_2(D)$
- 1. For $q_1(D)$ to hold, there is a valuation v s.t. $v(q_1) \in D$
- 2. We will show that the composition $g = v \circ h$ is a valuation for q_2
 - 2a. By definition of h, for every $R(x_1, x_2, ...)$ in q_2 , $R(h(x_1), h(x_2), ...)$ in q_1
 - 2b. By definition of v, for every $R(x_1, x_2, ...)$ in q_2 , $R(v(h(x_1)), v(h(x_2)), ...)$ in D

 $q = v \circ h$

q(x) = v(h(x))

<u>Example</u>



Combined complexity of CQC and CQE

Corollary:

The following problems are NP-complete (in the size of Q or Q'):

- 1) Given two (Boolean) conjunctive queries Q and Q', is $Q \subseteq Q'$?
- 2) Given a Boolean conjunctive query Q and an instance D, does $D \models Q$?

Proof:

(a) Membership in NP follows from the Homomophism Theorem: $Q \subseteq Q'$ if and only if there is a homomorphism h: $Q' \rightarrow Q$

(b) NP-hardness follows from 3-Colorability: G is 3-colorable if and only if $Q^{K_3} \subseteq Q^{G_1}$

The Complexity of Database Query Languages

	Relational	CQs
	Calculus	
Query Eval.:	In LOGSPACE	In LOGSPACE
Data Complexity	(hence, in P)	(hence, in P)
Query Eval.:	PSPACE-	NP-complete
Combined Compl.	complete	
Query Equivalence	Undecidable	NP-complete
& Containment		

Based on Phokion Kolaitis' "Logic and Databases" series at Simons Institute, 2016. <u>https://simons.berkeley.edu/talks/logic-and-databases</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Exercise: Find Homomorphisms

$q_1: \{E(x,y), E(y,z), E(z,w)\}$

Order of subgoals in the query does not matter (thus written here as sets)

q₂: {E(x,y),E(y,z),E(z,x)}

 $q_3: \{E(x,y), E(y,x)\}$

what is the containment relation between these queries ?

 $q_4: \{E(x,y), E(y,x), E(y,y)\}$ $q_5: \{E(x,x)\}$

Example by Andreas Pieris: <u>https://homepages.inf.ed.ac.uk/apieris/courses/atfd2020/index.html</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Exercise: Find the Homomorphisms



q₁: {E(x,y),E(y,z),E(z,w)} $x \longrightarrow y \longrightarrow z \longrightarrow w$

Order of subgoals in the query does not matter (thus written here as sets)

 $q_{2}: \{E(x,y), E(y,z), E(z,x)\} \qquad q_{3}: \{E(x,y), E(y,x)\} \\ \times \longrightarrow y \\ What is the containment relation \\ between these queries ? \\ q_{4}: \{E(x,y), E(y,x), E(y,y)\} \qquad q_{5}: \{E(x,y), E(y,y)\}$





Example by Andreas Pieris: <u>https://homepages.inf.ed.ac.uk/apieris/courses/atfd2020/index.html</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Exercise: Find the Homomorphisms







 $q_1(x,y) := R(x,u), R(v,u), R(v,y)$ $var(q_1) = \{x, u, v, y\}$

 $q_2(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y)$ var $(q_2) = \{x, u, v, w, t, y\}$

Are these queries equivalent ?



 $q_{1}(x,y) := R(x,u), R(v,u), R(v,y) \qquad var(q_{1}) = \{x, u, v, y\}$ $q_{2}(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y) \qquad var(q_{2}) = \{x, u, v, w, t, y\}$

$q_1 \rightarrow q_2$ Thus ?

which query contains the other?



 $q_{1}(x,y) := R(x,u), R(v,u), R(v,y) \qquad var(q_{1}) = \{x, u, v, y\}$ $q_{2}(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y) \qquad var(q_{2}) = \{x, u, v, w, t, y\}$





 $q_1(x,y) := R(x,u), R(v,u), R(v,y) \qquad var(q_1) = \{x, u, v, y\}$

 $q_2(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y)$ $var(q_2) = \{x, u, v, w, t, y\}$

Is there any homomorphism $q_2 \rightarrow q_1$ and then $q_1 \subseteq q_2$?



 $q_1(x,y) := R(x,u), R(v,u), R(v,y)$ var($q_2(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y)$ var(

$$var(q_1) = \{x, u, v, y\}$$

$$\uparrow \uparrow \uparrow \downarrow$$

$$var(q_2) = \{x, u, v, w, t, y\}$$



 $q_2 \rightarrow q_1$ Thus $q_1 \subseteq q_2$?



 $q_1(x,y) := R(x,u), R(v,u), R(v,y)$ $q_2(x,y) := R(x,u), R(v,u), R(v,w), R(t,w), R(t,y)$

$$var(q_1) = \{x, u, v, y\}$$

$$\uparrow \uparrow \uparrow \downarrow$$

$$var(q_2) = \{x, u, v, w, t, y\}$$



Thus $q_1 \subseteq q_2$