

Inconsistent Databases

Querying Relational Databases

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	code	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

Querying Relational Databases

List the airlines that fly directly from London to Glasgow

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	code	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh



{U2}

$Q(z) :- \text{Airport}(x, \text{London}), \text{Airport}(y, \text{Glasgow}), \text{Flight}(x, y, z)$

Semantic Information About the Data

inconsistency!!!

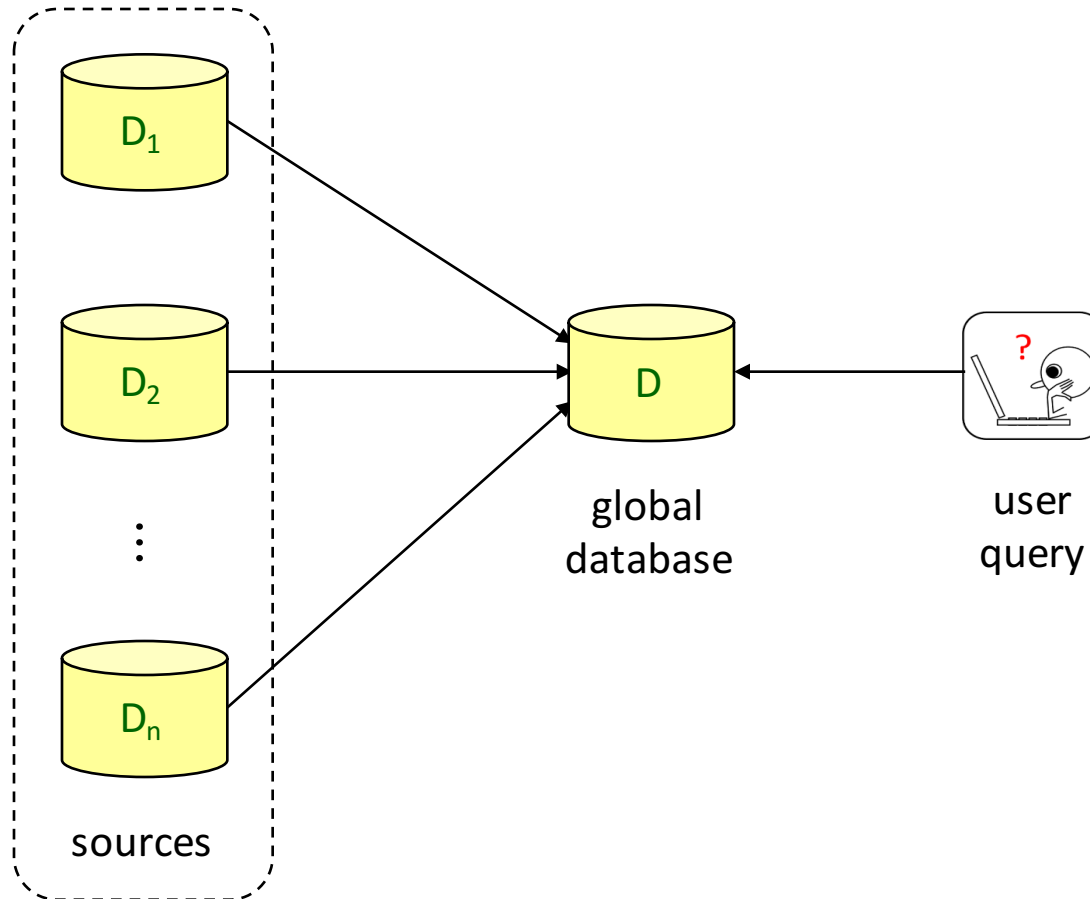
Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	code	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

The code uniquely determines the airport

Main Source of Inconsistency

data is coming from several conflicting sources



Querying Relational Databases

List the airlines that fly directly from London to Glasgow

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS



{U2} ?

Airport	code	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

$Q(z) :- \text{Airport}(x, \text{London}), \text{Airport}(y, \text{Glasgow}), \text{Flight}(x, y, z)$

Integrity Constraints

inconsistency!!!

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

$$\text{Key}(\text{Airport}) = \{1\} \quad \equiv \quad \forall x \forall y (\text{Airport}(x,y) \wedge \text{Airport}(x,z) \rightarrow y = z)$$

Primary Keys

at most one key per relation

inconsistency!!!

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

$$\text{Key}(\text{Airport}) = \{1\} \quad \equiv \quad \forall x \forall y (\text{Airport}(x,y) \wedge \text{Airport}(x,z) \rightarrow y = z)$$

Primary Keys

at most one key per relation

- Consider a database D , and a primary key $\sigma : \text{Key}(R) = \{i_1, \dots, i_n\}$

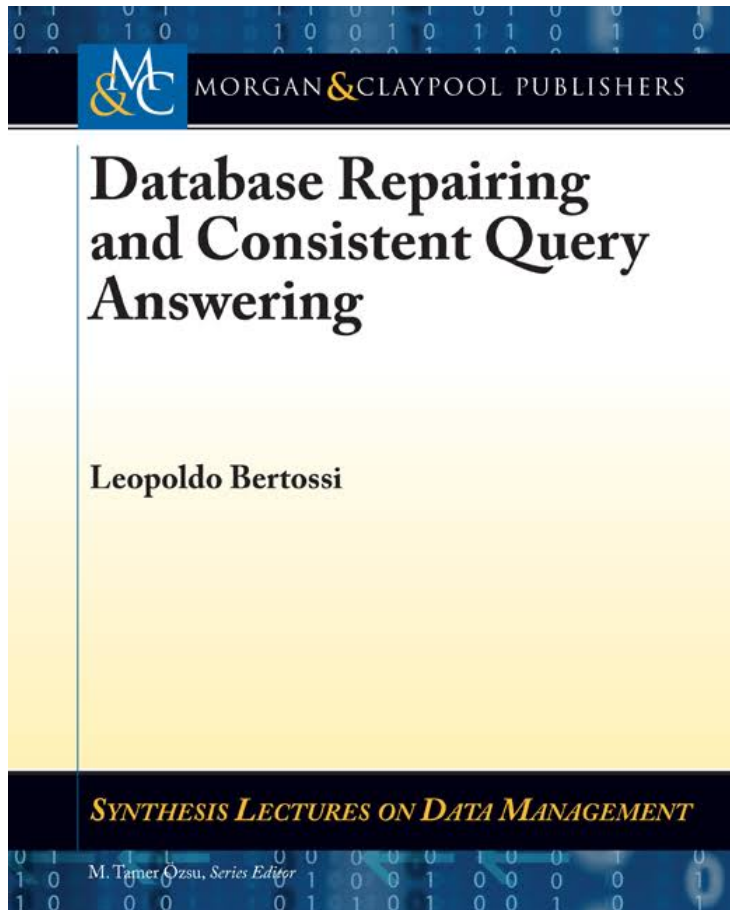
We say that D satisfies σ if, for every two atoms $R(a_1, \dots, a_m)$ and $R(b_1, \dots, b_m)$ in D

such that $a_{i_1}, \dots, a_{i_n} = b_{i_1}, \dots, b_{i_n}$, it holds that $a_1, \dots, a_m = b_1, \dots, b_m$

- D satisfies a set of primary keys Σ , denoted $D \models \Sigma$, if D satisfies every key in Σ

In this case we say that D is **consistent** w.r.t. Σ ; otherwise, D is **inconsistent** w.r.t. Σ

Consistent Query Answering (CQA)

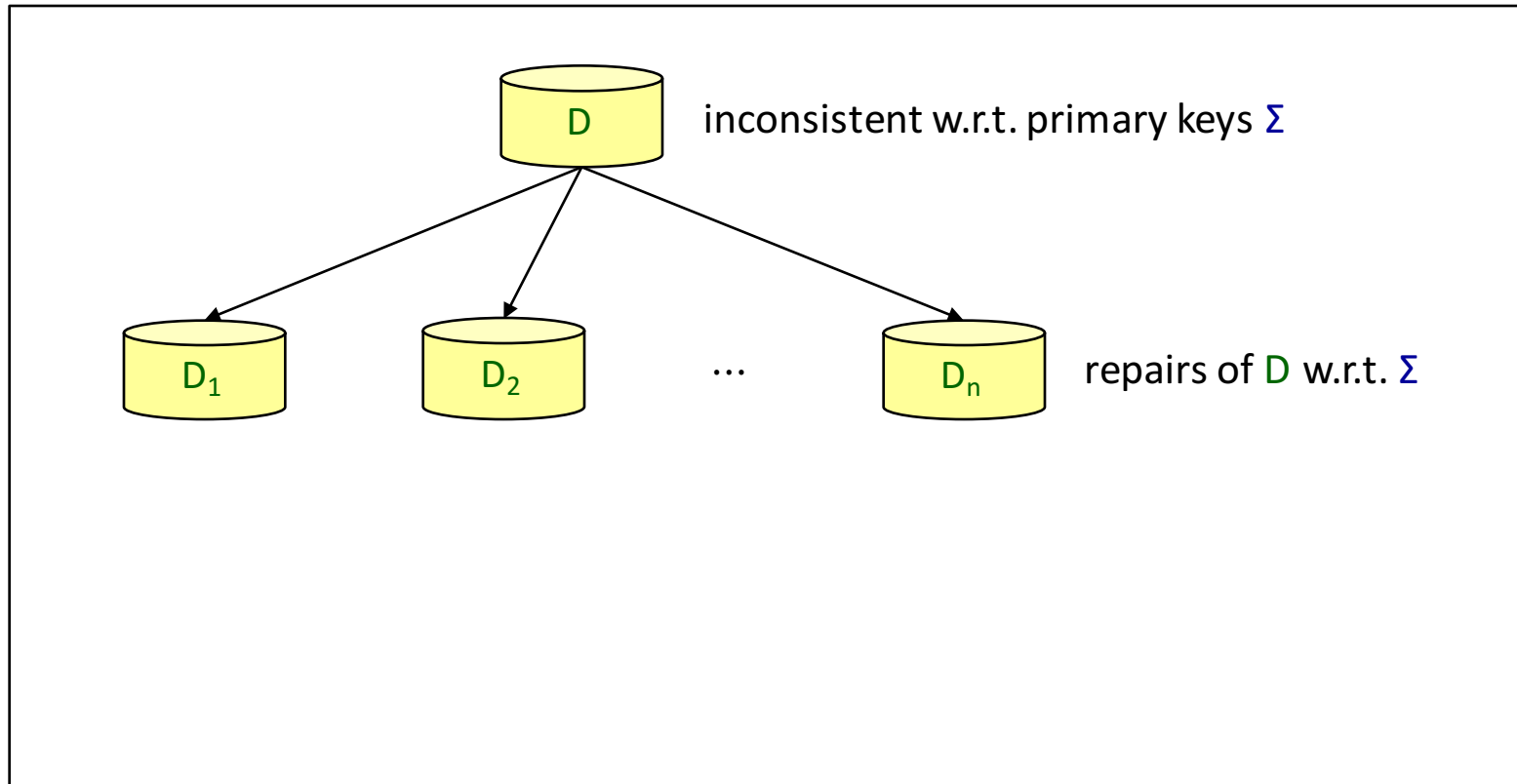


find meaningful answers to queries

when databases are inconsistent

Key Elements of CQA

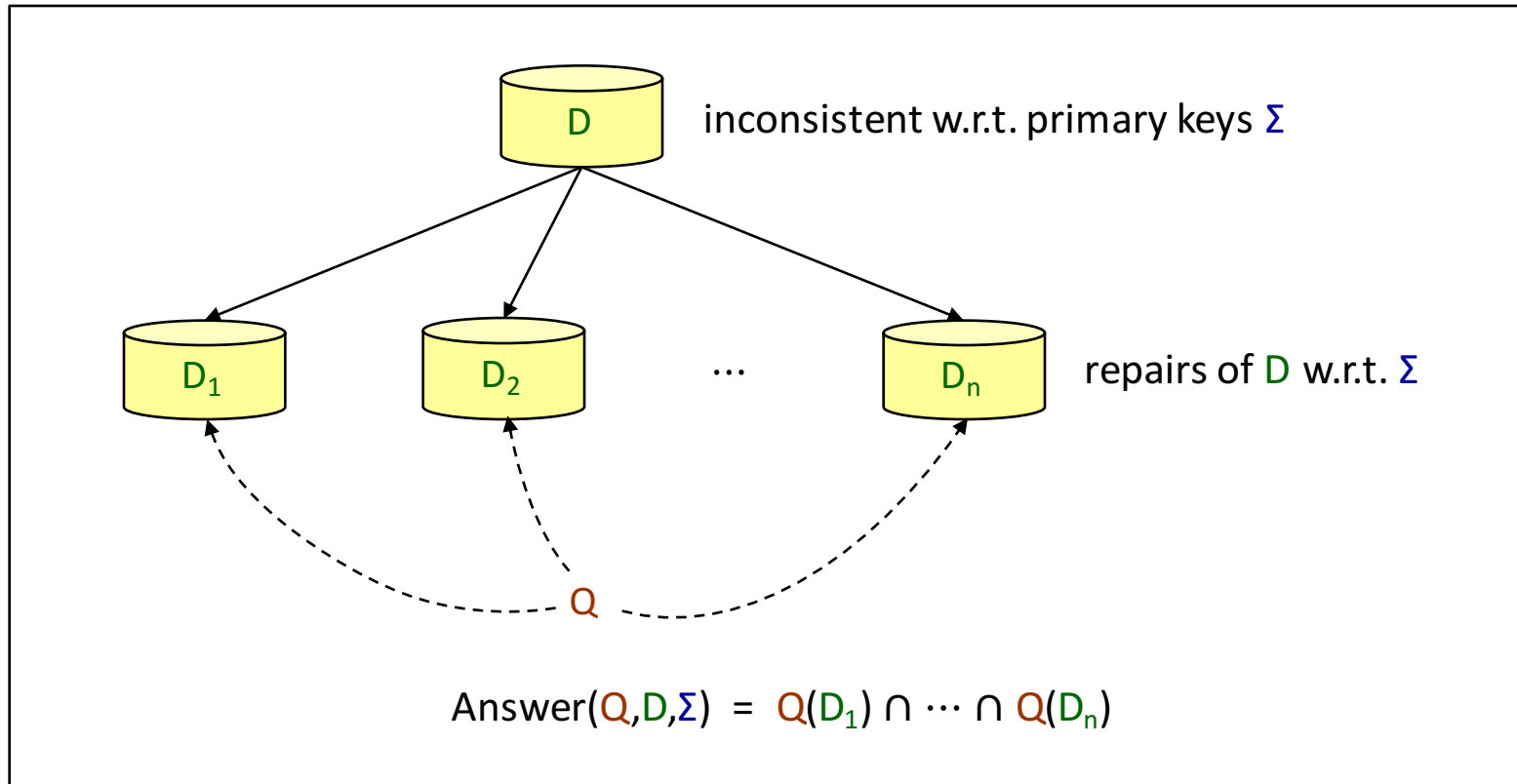
Repairs - consistent databases whose difference with D is “minimal”



Key Elements of CQA

Repairs - consistent databases whose difference with D is “minimal”

Consistent answers - answers that are true in all repairs



Key Elements of CQA

Consider a database D , and a set Σ of primary keys

A **repair** of D w.r.t. Σ is a database $D' \subseteq D$ such that the following conditions hold:

1. $D' \models \Sigma$
2. There is no $D'' \subseteq D$ such that $D'' \models \Sigma$ and $D' \subset D''$

$$\text{Answer}(Q, D, \Sigma) = \bigcap_{R \in \text{repairs}(D, \Sigma)} Q(R)$$

the set of repairs of D w.r.t. Σ

Repairs

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

Repair 1

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	GLA	Glasgow
	EDI	Edinburgh

Repair 2

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

Consistent Answers

List the airlines that fly directly from London to Glasgow

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

Consistent Answers

List the airlines that fly directly from London to Glasgow

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

Repair 1: {U2}

Consistent Answers

List the airlines that fly directly from London to Glasgow

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

Repair 1: {U2}

Repair 2: {}

Answer is empty

Consistent Query Answering

CQA(L)

Input: a database D , a set of primary keys Σ , a query $Q/k \in L$, a tuple of constants $\mathbf{t} \in \text{adom}(D)^k$

Question: $\mathbf{t} \in \text{Answer}(Q, D, \Sigma)$?

BCQA(L)

Input: a database D , a set of primary keys Σ , a Boolean query $Q \in L$

Question: is $\text{Answer}(Q, D, \Sigma)$ non-empty?

Theorem: $\text{CQA}(L) \equiv_L \text{BCQA}(L)$, where $L \in \{\text{RA}, \text{DRC}, \text{TRC}, \text{CQ}\}$

(\equiv_L means logspace-equivalent)

Data Complexity of BCQA

input D , fixed Σ and Q

$BCQA[\Sigma, Q](L)$

Input: a database D

Question: is $\text{Answer}(Q, D, \Sigma)$ non-empty?

Data Complexity of BCQA

Theorem: For $L \in \{\text{RA}, \text{DRC}, \text{TRC}, \text{CQ}\}$, $\text{BCQA}[\Sigma, Q](L)$ is coNP-complete for a fixed set of primary keys Σ , and a query $Q \in L$.

Proof:

- Guess a repair $R \in \text{repairs}(D, \Sigma)$, and check whether $Q(R)$ is empty
- Reduction from 3-Colorability to the complement of BCQA

3-Colorability

3COL

Input: an undirected graph $\mathbf{G} = (V, E)$

Question: is there a function $c : V \rightarrow \{\mathbf{R}, \mathbf{G}, \mathbf{B}\}$ such that $(v, u) \in E \Rightarrow c(v) \neq c(u)$?

coNP-hardness

Given an undirected graph $\mathbf{G} = (V,E)$

construct a database \mathbf{D} such that, for some fixed Σ and \mathbf{Q} , it holds that

\mathbf{G} is 3-colorable iff $\text{Answer}(\mathbf{Q},\mathbf{D},\Sigma)$ is empty

$$\mathbf{D} = \{\text{Edge}(u,v) : (u,v) \in E\} \cup$$
$$\{\text{Color}(v,r), \text{Color}(v,g), \text{Color}(v,b) : v \in V\} \cup$$
$$\{\text{Bad}(r,r), \text{Bad}(g,g), \text{Bad}(b,b)\}$$
$$\Sigma = \{\text{Key}(\text{Color}) = \{1\}\}$$
$$\mathbf{Q} :- \text{Edge}(x,y), \text{Color}(x,z), \text{Color}(y,w), \text{Bad}(z,w)$$

Lemma: \mathbf{G} is 3-colorable iff there is $\mathbf{R} \in \text{repairs}(\mathbf{D},\Sigma)$ such that $\mathbf{Q}(\mathbf{R})$ is empty

Data Complexity of BCQA

Theorem: For $L \in \{\text{RA}, \text{DRC}, \text{TRC}, \text{CQ}\}$, $\text{BCQA}[\Sigma, Q](L)$ is coNP-complete for a fixed set of primary keys Σ , and a query $Q \in L$.

Proof:

- Guess a repair $R \in \text{repairs}(D, \Sigma)$, and check whether $Q(R)$ is empty
- Reduction from 3-Colorability to the complement of BCQA

Tackle High Data Complexity

Two main research directions:

1. Isolate classes of queries (in fact, classes of CQs) for which the problem can be solved efficiently in data complexity
2. Provide data-efficient approximations

Consistent Answers

List the airlines that fly directly from London to Glasgow

Flight	origin	destination	airline
	VIE	LHR	BA
	LHR	EDI	BA
	LGW	GLA	U2
	EDI	VIE	OS

Airport	<u>code</u>	city
	VIE	Vienna
	LHR	London
	LGW	London
	LGW	Lilongwe
	GLA	Glasgow
	EDI	Edinburgh

Repair 1: {U2}

Repair 2: {}

Answer = { (U2,50%) }

Relative Frequency of a Boolean Query

$$\text{RF}(Q, D, \Sigma) = \frac{|\{R : R \in \text{repairs}(D, \Sigma) \text{ such that } Q(R) \text{ is non-empty}\}|}{|\text{repairs}(D, \Sigma)|}$$

Consistent Query Answering Revisited

RF-BCQA(L)

Input: a database D , a set of primary keys Σ , a Boolean query $Q \in L$

Output: $RF(Q, D, \Sigma)$

we can naturally talk about the data complexity the problem

RF-BCQA[Σ, Q](L) - input D , fixed Σ and Q

Data Complexity of BCQA

Theorem: For $L \in \{\text{RA}, \text{DRC}, \text{TRC}, \text{CQ}\}$, $\text{BCQA}[\Sigma, Q](L)$ is $\text{FP}^{\#P}$ -complete for a fixed set of primary keys Σ , and a query $Q \in L$.

- This essentially means that computing the relative frequency of a Boolean query is a hard problem, even for CQs
- The goal is to *efficiently approximate* the relative frequency

Efficient Approximations

fix a set of primary keys Σ , and a Boolean CQ Q

A **fully polynomial-time randomized approximation scheme** (FPRAS) for RF-BCQA[Σ, Q](CQ)

is a randomized algorithm **Approximation** that accepts as input

a database D , and numbers $\epsilon > 0$ and $0 < \delta < 1$,

runs in polynomial time in the size of D , $1/\epsilon$ and $\log(1/\delta)$, and

produces a random variable **Approximation**(D, ϵ, δ) such that

$$\Pr(|\mathbf{Approximation}(D, \epsilon, \delta) - \text{RF}(Q, D, \Sigma)| \leq \epsilon \cdot \text{RF}(Q, D, \Sigma)) \geq 1 - \delta$$

Sampling

fix a set of primary keys Σ , and a Boolean CQ Q

Sample $[\Sigma, Q]$

Input: a database D

Output: 0 or 1

$\text{Repair} := \emptyset$

for $i = 1$ **to** n **do**

choose $P(\mathbf{t}) \in B_i$ with probability $1/|B_i|$

$\text{Repair} := \text{Repair} \cup \{P(\mathbf{t})\}$

if $\text{Repair} \models Q$ **then**

return 1

else

return 0

end

$\{B_1, B_2, \dots, B_n\}$ is a partition of D such that
each B_i collects conflicting (w.r.t. Σ) atoms of D

Efficient Approximation for RF-BCQA[Σ, Q](CQ)

fix a set of primary keys Σ , and a Boolean CQ Q

Approximation[Σ, Q]

Input: a database D , and numbers $\epsilon > 0$ and $0 < \delta < 1$

Output: random number in $[0,1]$

Experiments := $((2+\epsilon) \cdot m^k) / \epsilon^2 \cdot \ln(2/\delta)$, where k is the number of atoms in Q , and m is the size of the largest B_i

Sum := \emptyset

Counter := \emptyset

repeat

Sum := Sum + **Sample**[Σ, Q](D)

Counter := Counter + 1

until Counter = Experiments

return Sum/Experiments

Efficient Approximation for RF-BCQA[Σ, Q](CQ)

fix a set of primary keys Σ , and a Boolean CQ Q

Theorem: $\text{Approximation}[\Sigma, Q]$ is an FPRAS for $\text{BCQA}[\Sigma, Q](\text{CQ})$

Recap

- Inconsistent databases - do not conform with the integrity constraints coming with the underlying schema (such as primary keys)
- Consistent query answering (CQA) - find meaningful answers to queries when databases are inconsistent
- CQA is a hard problem, even in data complexity for CQs and primary keys
- Isolate classes of queries for which CQA is efficient in data complexity
- Provide data-efficient approximations schemes