# DATA QUALITY AWARE QUERYING\*

## Ognjen Savković Free University of Bozen-Bolzano, Italy

\* Supported by the project MAGIC, funded by the Province of Bozen-Bolzano

#### What is Data Quality?

- Data is of a high quality if it is fit for intended uses
- Typically, **Data Quality** is considered as a problem of "**dirty data**". E.g.,
  - **Unclean data**, *student('John J.','Street\*^#@0','+00399#)?:<'*), or
  - **Missing values**, *student('John J.',NULL,NULL)*
- But there are also more subtle Data Quality aspects
  - **Data Completeness** (we cannot see what is missing)
  - **Data Stability** (data that we see is updating frequently)

# QUERY COMPLETENESS REASONING UNDER CONSTRAINTS

joint work with Werner Nutt, Sergey Paramonov

#### Bolzano is in the Province of South Tyrol



- Trilingual autonomous province in the North of Italy
  - Has its own school administration

#### Ex: School Data Management in Bolzano



- School Database is decentralized
- Some data is inserted voluntarily

**Statistical Reports** 

#### Running Example: School Schema

pupil(name, ccode, sname)
 e.g., pupil(john, a1, goethe)

class(ccode, sname, level, branch)
 e.g., class(a1, gothe, 1, science)

takes(name, activity)
 e.g., takes(john, chess)

#### Query Completeness Reasoning under Constraints

- We want to reason about Query Completeness over
  - Partially complete database and
  - Constraints that hold over the database



#### Formalization: Incomplete Database

- Incomplete database *D* is a pair  $D = (D^i, D^a)$  [Motro 1989]
  - D<sup>i</sup> ideal database (complete facts that holds in the real world)
  - D<sup>a</sup> available database (actual database)



Semantics

if  $(D^i, D^a) \models Compl(Q)$  then  $Q(D^i) = Q(D^a)$ 

#### **Table Completeness Statements**

A table completeness (TC) statement for a relation R is an expression

Compl( $R(s_1, ..., s_n)$ ; G) [Halevy 96]

consisting of

• an R-atom  $R(s_1,...,s_n)$ 

• a condition G is conjunction of atoms

The TC-statement  $C = Compl(R(s_1,..., s_n); G)$  can be seen as a rule

$$r_{c} = R^{i}(s_{1},...,s_{n}), G^{i} \rightarrow R^{a}(s_{1},...,s_{n})$$

Semantics:  $(D^{i}, D^{a}) \models C$  iff  $(D^{i}, D^{a}) \models r_{C}$ 

#### **TC-QC Reasoning Problem**

• We want to decide (TC-QC entailment) "Does a set of table completeness (TC) statements C entail completeness of query Q (QC) ?"

 $\forall (D^{i}, D^{a})$  $(D^{i}, D^{a}) \models C \text{ iff } (D^{i}, D^{a}) \models Compl(Q)$ 

or shortly

 $C \models Compl(Q)$ 

## T<sub>C</sub> Operator

To  $C = Compl(R(\underline{s}); G)$  we associate the query

 $Q_{C}(\underline{s}) := R(\underline{s})$  , G

and the transformation on db instances

 $T_{C}(D) := \textcircled{R}(\underline{t}) \mid \underline{t} \in Q_{C}(D) \}$ 

For a set C of TC statements we define the transformation

 $T_{\mathcal{C}}(D) := \bigcup_{C \in \mathcal{C}} T_{C}(D)$ 

## T<sub>C</sub> Operator: Properties

- For a given (ideal) database D, operator  $T_C$  computes  $T_C(D)$  such that
- (D, T<sub>C</sub>(D)) is the least incomplete databases that satisfies C

In other words,

- 1.  $(D, T_{C}(D)) \models C$  and
- 2.  $(D^{i}, D^{a}) \models C \text{ iff } T_{C}(D^{i}) \subseteq D^{a}$

## **TC-QC: Characterization**

#### Let

- C set of TC statements
- $Q(\underline{x}) := A1,...,An$  conjunctive query
- $D_Q :=$  *frozen version* of A1,...,An with  $\theta$

#### Theorem:

 $C \models Compl(Q)$  iff  $\theta \underline{\mathbf{x}} \in Q(T_C(D_Q))^*$ 

#### **Example: Plain Reasoning**

- "Who are the pupils in class a1 in Goethe school that play chess?" Q<sub>a1g</sub>(N) :- pupil(N, a1, goethe), takes(N, chess)
- "We have all pupil from a1 Goethe school"
   C<sub>a1</sub>: Compl( pupil(N, a1, goethe) ; true)
- "We have all activity takers in Goethe school"
   C<sub>activg</sub>: Compl( takes(N, A) ; pupil(N, C, goethe))
- Is  $Q_{a1g}$  complete given the completeness assumptions?

## Example: Plain Reasoning(2)

- 1. Assume  $Q_{a1g}$  returns n' over  $D^i$
- 1. Then, pupil(n', a1, goethe), takes(N, chess)  $\in D^i$
- 1. Now, according  $C_{a1g}$  followspupil(n', a1, goethe)  $\in D^a$ and, according  $C_{activg}$  followstakes(n', chess)  $\in D^a$
- 2. Therefore,  $n' \in Q_{a1g}(D^a)$ ,  $\rightarrow Q_{a1g}$  is complete!

What if have schema constraints?

#### Schema Constraints

- We consider 2 kinds of constraints
- Foreign Keys (FKs)
  - "For every pupil record exists corresponding class record"
  - pupil[ccode, sname] ⊆ class[ccode, sname]
- Conditional Finite Doman Constraints (CFDCs)
  - "Science classes in Goethe school can be either a1 or b2"
  - class[sname=Goethe, branch=science][ccode] = {a1,b2}

#### **TC-QC under Constraints**

- $\mathcal{K}$  is a set of PKs and FKs
- $\mathcal{F}$  is a set of CFDCs
- $C \models_{\mathcal{K},\mathcal{F}} \text{Compl}(Q)$  (non-enforced FKs)  $\forall (D^i, D^a)$ : if  $D^i \models \mathcal{F}$ ,  $D^a \models \mathcal{F}$  , and  $D^i \models \mathcal{K}$ then  $(D^i, D^a) \models C$  iff  $(D^i, D^a) \models \text{Compl}(Q)$
- $C \models {}^{e}_{\mathcal{K},\mathcal{F}} \operatorname{Compl}(Q)$  (enforced FKs)  $\forall (D^{i}, D^{a})$ : if  $D^{i} \models \mathcal{F}$ ,  $D^{a} \models \mathcal{F}$  and  $D^{i} \models \mathcal{K}$  and  $D^{a} \models \mathcal{K}$ then  $(D^{i}, D^{a}) \models C$  iff  $(D^{i}, D^{a}) \models \operatorname{Compl}(Q)$

# Example: Reasoning under non-enforced FKs

- "Who are the pupils in Goethe school?"  $Q_g(N) := pupil(N, C, goethe)$
- "We have all pupil from a1 and b2 Goethe school"
   C<sub>a1</sub>: Compl( pupil(N, a1, goethe) ; class(a1, goethe, L, B) )
   C<sub>b2</sub>: Compl( pupil(N, b2, goethe) ; class(b2, goethe, L, B) )
- Can we conclude completness for Q<sub>gs</sub>?

## Example: Reasoning under non-enforced FKs (2)

Assume constraints

fk<sub>1</sub>: pupil[ccode, sname] I class[ccode, sname]
cfdc<sub>1</sub>: class[sname=goethe][ccode] = {a1,b2}

- 1.  $n' \in Q_g(D^i)$ , then pupil(n', c', goethe)  $\in D^i$
- 1. From fk<sub>1</sub>, class(c', goethe,  $l_{c',goethe}$ ,  $b_{c',goethe}$ )  $\in D^i$
- 2. From  $cfdc_1$  follows that c'=a1 or c'=b2 in class-record
- 3. In either case, acc.  $C_{a1}$  and  $C_{b2}$ , follows pupil(n', c', goethe)  $\in D^a$
- 1. Therefore,  $n' \in Q_g(D^a)$ ,  $\rightarrow Q_g$  is complete!

#### **TC-QC under CSTRs: Characterization**

- C and  $Q(\underline{x}) := A1,...,An$  as before
- $\mathbf{\mathcal{K}}$  set of PKs and acyclic FKs
  - $Chase_{\kappa}$  oblivious chase
- **F** set of CFDCs
  - $\Gamma$  set of all "maximal relevant" instantiations of  $Chase_{\mathcal{K}}(D_Q)$  according  $\mathcal{F}$

#### Theorem: The following are equivalent

- a)  $C \models_{\mathcal{K},\mathcal{F}} \text{Compl}(Q)$
- b)  $\gamma \theta \underline{\mathbf{x}}' \in Q(\gamma(T_C(\gamma(Chase_{\mathcal{K}}(D_Q))))) \text{ for all } \gamma \in \Gamma$
- \* b) is a  $\Pi^{P_2}$  problem

#### Example: Reasoning under enforced FKs

• "Who are the science pupils in Goethe school?"  $Q_{gs}(N) := pupil(N, C, goethe), class(C, goethe,L, science)$ 

- "We have all pupil from a1 and b2 Goethe school"
   C<sub>a1</sub>: Compl( pupil(N, a1, goethe) ; class(a1, goethe, L, B) )
   C<sub>b2</sub>: Compl( pupil(N, b2, goethe) ; class(b2, goethe, L, B) )
- And assume constraints

fk<sub>1</sub>: pupil[ccode, sname] I class[ccode, sname]
cfdc<sub>1</sub>: class[sname=goethe][ccode] = {a1,b2}

• Can we conclude completness for  $Q_{gs}$ ?

#### Example: Reasoning under enforced FKs (2)

- 1.  $n' \in Q_{gs}(D^i)$ , then pupil(n', c', goethe), class(c',goethe,l',science)  $\in D^i$
- 1. From  $cfdc_1$  follows that c'=a1 or c'=b2 in class-record
- 2. In either case, acc.  $C_{a1}$  and  $C_{b2}$ , follows pupil(n', c', goethe)  $\in D^a$
- 3. From fk<sub>1</sub>, class(c', goethe,  $l_{c',goethe}$ ,  $b_{c',goethe}$ )  $\in D^a$  (no!) Since  $D^a \subseteq D^i$  it must be class(c',goethe,l',science)  $\in D^a$
- 4. Therefore,  $n' \in Q_{gs}(D^a)$ ,  $\rightarrow Q_{gs}$  is complete!

#### Implementing Completeness Reasoning

- How can one implement completeness reasoning?
- Reminder: TC-QC ranges from NP to  $\Pi_{2}^{P}$
- **SMT** (SAT modulo theories) solvers?
- ASP?
  - Receipt: Follow the characterization theorems

#### Encoding TC-QC into ASP

 "Freeze" Q<sub>a1g</sub>(N) over ideal db pupil\_i(n', a1, goethe). takes\_i(n', chess).

- 2. Represent TCs as rules from ideal to available db pupil\_a(N, a1, goethe) :- pupil\_i(N, a1, goethe) takes\_a(N, A) :- takes\_i(N,A), pupil\_i(N, C, goethe)
- Test completeness with "test predicate" on available db q\_test :- pupil\_a(N, a1, goethe), takes\_i(N, chess)

C = Compl(Q) iff q\_test is **entailed** by the program

/

THE ..

~		ayık-demo.int.dhii	JZ.IL/VIGD2013/COP	istraints/index.jsp?action=reason			
				🔂 Ad	ld new query		
	ID	Descripti	Description \$		Actions	Selected Query	
	• Q	L Select the	names of all pupil	s that attend a primary school.	XL	FROM pupil AS p, school AS s	
	ୢ	2 Select the the Bolzar	names of all pupil no district and that	s that attend a primary school in learn some language.	XL	WHERE s.type='primary' AND p.sname=s.sname	
	ୢ	3 Select the the Bolza	names of all 1st le no district and that	vel pupils that attend a school in learn some language	XL		
	0 Q	4 Select all	language learners.		XL		
	ୢୣ	5 Select all	classes.		XL		
	0 Q	5 Who learn	s English?		XL		
	ୢ	3 Select all language.	pupils from school	s in Bolzano who learn a	XL		
	0 Q	Give me a	ll pupils from prim	ary schools.	XL		
- R	Compl Genez Speci	/ is not complete eteness calculat alization calcul alization(s) cal	ed in 6 ms ated in 13 ms culated in 143 m	s			<b>V</b> Ku
R	esult Quer Compl Genez Speci	/ is not complete eteness calculat alization calcul alization (s) cal plete tables ing parts of the tab	ed in 6 ms ated in 13 ms culated in 143 m ples are incomplete	s . Please collect the missing data ar Condition	nd confirm it	by adding the corresponding TC-statem	nents.
R	esult Quer CompJ Gener Speci Incom The follow	/ is not complete eteness calculat alization calcul alization (s) cal plete tables ing parts of the tak able upil(P_pname,S sn	ed in 6 ms ated in 13 ms culated in 143 m ples are incomplete ame,P_code)	<ul> <li>Please collect the missing data an</li> <li>Condition</li> <li>school(S_sname,'primary'.S di</li> </ul>	nd confirm it strict)	by adding the corresponding TC-statem	nents.
R	esult Quer CompJ Genex Speci Incom The follow	/ is not complete eteness calculat alization calcul alization (s) cal plete tables ing parts of the tab iable upil(P_pname,S_sn	ed in 6 ms ated in 13 ms culated in 143 m ples are incomplete ame,P_code)	<ul> <li>Please collect the missing data an</li> <li>Condition</li> <li>school(S_sname,'primary',S_di</li> </ul>	nd confirm it strict)	by adding the corresponding TC-statem	nents.
R	esult Compl Genes Speci Incom The follov	y is not complete eteness calculat alization calcul alization (s) cal plete tables ing parts of the tab ing parts of the tab ing parts of the tab cable upil(P_pname,S_sname) ete Query Approxim	ed in 6 ms ated in 13 ms culated in 143 m ples are incomplete ame,P_code) nation	s . Please collect the missing data an <b>Condition</b> school(S_sname,'primary',S_di	nd confirm it strict)	by adding the corresponding TC-statem	nents.
	esult CompJ Genex Speci Incom The follow	/ is not complete eteness calculat alization calcul alization (s) cal plete tables ing parts of the tab ing parts of the tab able upil(P_pname,S_sn ete Query Approxim	ed in 6 ms ated in 13 ms culated in 143 m oles are incomplete ame,P_code) nation	Please collect the missing data ar Condition school(S_sname,'primary',S_di	nd confirm it strict)	by adding the corresponding TC-statem	nents.
	esult Compl Genez Speci Incom The follow	/ is not complete eteness calculat alization calcul alization (s) cal plete tables ing parts of the tab able upil(P_pname,S_sn ete Query Approxim	ed in 6 ms ated in 13 ms culated in 143 m oles are incomplete ame,P_code) mation Query Not avail	<ul> <li>Please collect the missing data an</li> <li>Condition school(S_sname,'primary',S_di</li> </ul>	nd confirm it strict)	by adding the corresponding TC-statem	nents.
	esult Compl Genes Speci Incom The follow Comp	y is not complete eteness calculat alization calcul alization (s) cal plete tables ing parts of the tak <b>able</b> upil(P_pname,S_sn ete Query Approxin Generalization al	ed in 6 ms ated in 13 ms culated in 143 m oles are incomplete ame,P_code) mation Query Not avail SELECT I FROM pt WHERE s AND p.s	s Please collect the missing data an Condition school(S_sname,'primary',S_di able DISTINCT p.pname upil AS p, school AS s .type='primary' name=s.sname	nd confirm it strict)	by adding the corresponding TC-statem	nents.

# QUERY STABILITY IN BUSINESS PROCESSES

joint work with Werner Nutt

# **Ex: Student Registration**

Faculty of Economics	Student places 2012 / 2013	Enrollments 2011 / 2012	Enrollments 2012 / 2013	%
Bachelor in Economics and Management				20.78%
Bachelor in Tourism, Sport and Event Management				74.51%
Bachelor in Economics and Social Sciences				17.24%
Master in Entrepreneurship and Innovation				114.29%
Master in Economics and Management of the public sector				22.73%
Sum	370			41.45%

Faculty of Computer Science	Student Places 2012 / 2013	Enrollments 2011 / 2012	Enrollments 2012 / 2013	%
Bachelor in Computer Science and Engineering				30.43%
Master of Science in Computer Science				-44.44%
PhD in Computer Science				#VALUE!
Sum	195			-25.37%

## Data Stability

- Data Completeness is a rather strong requirement
  - Some information will be always incomplete or incomplete "for a short time"
- Data Stability is concerned whether data is stable in some period of time
  - "Is the number of students a stable figure during 2<sup>nd</sup> semester?"
  - "What is the longest periods for which the student numbers are stable?"



Business processes that manipulate data can tell us when the data (or query answer) is stable

#### What are Business Processes?

- Business Processes (BPs) are set of activities organized to accomplish a specific goal
  - E.g., Student registration, Car production, etc.

#### Business Processes

- Used for documenting how the companies and organizations operate their business using languages such as BPMN
- Can be executed using Business Process Engines
- Allow different types of analysis

   (simulation, verification, process mining, etc.)

# Ex: Student Registration (cont'd)

• Student Registration managed by a Business Process (in BPMN)



- Some questions:
  - How final or reliable are the figures that we see?
  - For which programs the figures become final eventually?
  - Which are the time periods in which the figures are stable?

#### Ex: Purchasing Information System @ Univ. Bolzano



- Is the available amount of the budget a stable figure?
- What is the maximal and what minimal amount in the budget that is not going to be spent? Etc.

#### Process and Data: Two sides of the same coin

• In current BP languages,

the interaction between BPs and Data is very limited

- Recently, several models have been proposed
  - Relational transducers [Abiteboul et al., 1998]
  - Data-driven Web Systems [Deutsch et al., 2004; 2007]
  - Data-Centric Dynamic Systems [Calvanese et. al., 2012; 2013]

#### **Data-aware Business Process**

- Data-aware Business Process (DABP) is defined as a tuple D=(N, L, D, 0)
  - N = (P,T) process net represented as a finite graph
  - $L = (L_e, L_w) labeling functions$  that label edges t from T
    - $L_e(t)$  execution condition, Boolean Conjunctive Query over D, O
    - $L_w(t)$  writing action, a horn rule Head<sub>t</sub> :- Body<sub>t</sub>
  - D relational database
  - O set of (business) objects, where each o has associated
    - $M(o) \in P place$  of the process net
    - $Val(o) = R(\underline{s}) value$  as a record of fixed size
- $\mathcal{P} = (N, L) process part$  (static part)
- C = (D, 0) configuration part (dynamic part)

#### **Execution in DABP**

- An execution step in DABP is either
  - a) Traversal of an edge by an object, Object o can traverse t from T iff  $D,Val(o) = L_e(t)$
  - b) Creation of a new object in the start place of a process which is "always possible"
- Result of an execution step is a new configuration C' = (D', O'), written  $C \rightarrow C'$ 
  - a) If o traverses t then  $D' = D \cup \{\alpha \text{ Head}_t \mid \alpha \text{ Body}_t \subseteq D, \text{Val}(o)\}$ M(o) points to the new place, and the rest remains the same
  - b) If o is a newly introduced object then O' = O U {o} and M(o)=start, and the rest remains the same

#### **Query Stability in DABP**

• Given DABP  $\mathcal{D}=(N, L, D, 0)$ , configuration C' = (D', 0') is **reachable**, if there exist execution steps  $C = (D, 0) \rightarrow C_1 \rightarrow ... \rightarrow C_n \rightarrow C'$ 

#### Query Stability

Conjunctive query  $Q(\underline{x}) := A1,...,An$  is stable in DABP  $\mathcal{D}=(N, L, D, 0)$ , if for every configuration C' = (D', 0') reachable from C = (D, 0)

Q(D) = Q(D')

## Ex: Student Registration (cont'd)

 "Do we have stable figures for the students of informatics?" Q(Name) ← student(Name, 'informatics')
 No! Because we can still enroll new students

 "Do we have stable figures for the students of economics?" Q(Name) ← student(Name,'economics')

 Yes! Because we cannot enroll any new student

## Ex: Student Registration (cont'd)

Faculty of Economics	Student places 2012 / 2013	Enrollments 2011 / 2012	Enrollments 2012 / 2013	%	Stable?
Bachelor in Economics and Management				20.78%	YES!
Bachelor in Tourism, Sport and Event Management				74.51%	YES!
Bachelor in Economics and Social Sciences				17.24%	YES!
Master in Entrepreneurship and Innovation				114.29%	YES!
Master in Economics and Management of the public sector				22.73%	YES!
Sum				41.45%	YES!

Faculty of Computer Science	Student Places 2012 / 2013	Enrollments 2011 / 2012	Enrollments 2012 / 2013	%	Stable?
Bachelor in Computer Science and Engineering				30.43%	YES!
Master of Science in Computer Science				-44.44%	NO!
PhD in Computer Science				#VALUE!	NO!
Si	um 🥵			-25.37%	NO!

## Checking Stability in DABP

- 3 DABP dialects:
  - **DABP**<sup>core</sup> forbids recursion in the horn writing rules
  - DABP<sup>chron</sup> assumes no objects initially
  - DABP<sup>recur</sup> no restrictions
- Computational Complexity of checking stability for CQ

complexity	DABP <sup>core</sup>	DABP <sup>chron</sup>	DABP <sup>recur</sup>	
data	$AC^0$	PTIME	CO-NP	
combined	CO-NP?	EXPTIME	co-NExpTime	
char. query lang.	CQ	DATALOG	DATALOG <sup>neg</sup>	

- data is measured in the size of the configuration
- combined is measured in the sizes of DABP and query

#### Conclusion

- 2 "subtle" aspects of data quality
  - query completeness
  - stability of queries answers
- Asses quality of query answers using metadata
  - completeness statements, constraints
  - business processes
- Characterization of query completeness under constraints
  - Basis for implementation in ASP
- Modeling query stability problem
  - Reasoning techniques for conjunctive queries

Thank you!

