

Error Generation for Evaluating Data-Cleaning Algorithms

Patricia C. Arocena
University of Toronto, Canada

Boris Glavic
Illinois Inst. of Technology, US

Giansalvatore Mecca
University of Basilicata, Italy

Renée J. Miller
University of Toronto, Canada

Paolo Papotti
QCRI Doha, Qatar

Donatello Santoro
University of Basilicata, Italy

Technical Report #TR-01-2015, September 14, 2015

ABSTRACT

We address the problem of generating errors within clean databases for the purpose of benchmarking data-cleaning and data-repairing algorithms. Our goal is to provide users with the highest possible level of control over the error generation process and at the same time develop a solution that scales to large databases. This is challenging, because the error generation problem is NP-complete. The main technical contribution of this paper is to develop an efficient PTIME algorithm which sacrifices completeness in an intelligent fashion that allows us to succeed under reasonable assumptions. However, scaling to databases with millions of tuples requires additional non-trivial optimizations including exploiting a symmetry property of data quality constraints.

1. INTRODUCTION

We consider the problem of empirically evaluating data-cleaning algorithms. Currently, there are no openly-available tools for systematically generating data exhibiting different types and degrees of quality errors. This is in contrast to related fields, for example, *entity resolution*, where well known datasets and generators of duplicated data exist and can be used to assess the performance of algorithms.

We allow the user full control of the data distribution by providing a clean database (DB) into which our error generator, BART¹, introduces errors. BART supports different kinds of random errors (including typos, duplicated values, outliers and missing or bogus values). However, our main technical contributions are related to the problem of generating errors in presence of data quality rules. For this purpose, we consider the powerful language of *denial constraints* [31]. Denial constraints (DC) are a rule-language that can express many data quality rules including functional dependencies

¹BART: Benchmarking Algorithms for data Repairing and Translation

Emp	Name	Dept	Salary	Mng
t_1	John	Staff	1000	Mark
t_2	Paul	Sales	1300	Frank
t_3	Jim	Staff	1000	Carl
t_4	Frank	Mktg	1500	Jack

MD	Name	Dept	Mng
t_{m1}	John	Staff	Mark
t_{m2}	Frank	Mktg	Jack

Figure 1: Example Clean Database

(FD), conditional functional dependencies (CFD) [8], cleaning equality-generating dependencies [20], and fixing rules [35]. In addition, BART permits a user to indicate parts of a database that are *reliable* and should not be changed, and hence we can express editing rules [17] that use master data to determine a repair.

BART provides the highest possible level of control over the error-generation process, allowing users to choose, for example, the percentage of errors, whether they want a guarantee that errors are detectable using the given constraints, and even provides an estimate of how hard it will be to restore the database to its original clean state (a property we call *repairability*). This control is the main innovation of the generator and distinguishes it from previous error generators that control mainly for data size and amount of error. It permits innovative evaluations of existing cleaning algorithms that reveal new insights on their (relative) performance when executed over errors with differing degrees of repairability.

Solving this problem proves surprisingly challenging, for two main reasons. First, we introduce two different variants of the error-generation problem, and show that they are both NP-complete. To provide a scalable solution, we concentrate on a polynomial-time algorithm that is correct, but not complete. Even so, achieving the desired level of scalability remains challenging. In fact, we show that there is a duality between the problem of injecting detectable errors in clean data, and the problem of detecting violations to database constraints in dirty data, a task that is notoriously expensive from the computational viewpoint. Finding the right trade-off between control over errors and scalability is the main technical problem tackled in this paper.

Example 1: [Motivating Example] Consider a database schema composed of two relations Emp and MD and the data shown in Figure 1. Suppose we are given a single data quality rule (a functional dependency) $d_1 : \text{Emp} : \text{Name} \rightarrow \text{Dept}$.

We can introduce errors into this database in many ways.

An example error is to change the Salary value of tuple t_4 to “3000”. This change (ch_0) creates an error that is not detectable using the given data quality rule d_1 . On the contrary, a second change (ch_1) that changes the Name value of tuple t_2 to “John” introduces a *detectable error* since this change causes tuples t_1 and t_2 to agree on employee names, but not on department values. Of course, if other errors are introduced into the DB (for example, t_1 .Dept is changed to “Sales”), then ch_1 may no longer be detectable. BART permits the user to control not only the number of errors, but also how many errors are detectable and whether they are detectable by a single rule or many rules.

We not only strive to introduce detectable errors into I , we also want to estimate how hard it would be to restore I to its original, clean state. Consider the detectable error introduced by ch_1 . This change removed the value Paul from the DB. Most repair algorithms will use evidence in the DB to suggest possible ways of repairing the data. If Paul is not in the active domain of the DB, then changing t_2 .Name to Paul will not be considered as a repair. BART lets the user control the *repairability* of errors by estimating how hard it would be to repair an error to its original value. \square

Error-Generation Tasks. We formalize the problem of error-generation using the notion of an *error-generation task*, \mathbf{E} , that is composed of four key elements: (i) a database schema \mathbf{S} , (ii) a set of denial constraints (DCs), Σ , encoding data quality rules over \mathbf{S} , (iii) a DB I of \mathbf{S} that is clean with respect to Σ , and (iv) a set of configuration parameters Conf to control the error generation process. These parameters specify, among other things, which relations are immutable, how many errors should be introduced, and how many of these errors should be detectable. They also let the user control the degree of repairability of the errors.

In this paper, we concentrate on a specific update model, the one in which the database is changed by updating attribute values (cells) only, and do not consider insertions or deletions. This has several advantages. First, this covers the vast majority of algorithms that have been proposed in the recent literature [3, 6, 7, 10, 17, 20, 27, 35, 36]. Second, it suggests a simple, flexible and scalable measure to assess the quality of repairs, a fundamental goal of this work, as discussed in Section 7. While it is not difficult to introduce insertions and deletions into the error-generation process, this completely changes the nature of the quality metric: in its full generality, this requires identifying tuple homomorphisms among the two databases, a problem for which there are currently no scalable algorithms. Since scalability is a primary concern in our work, we leave this to future work.

Contributions. Our main contributions are the following.

- (i) We present the first framework for generation of data-cleaning errors with fine-grained control over error characteristics. The tool allows users to inject a fixed number of detectable errors into a clean DB, possibly with a controlled level of repairability. It is intended for data-cleaning researchers and practitioners alike, and will allow them to compare data-cleaning solutions on a leveled playing field.
- (ii) We introduce a new computational framework based on *violation-generation queries* for finding candidate cells (tuple, attribute pairs) into which detectable errors can be introduced. We study when these queries can be answered efficiently and show that determining if there are detectable errors that can be introduced can be computationally hard.

- (iii) We introduce several novel optimizations for violation-generation queries. We show that extracting tuple samples, along with computing cross-products and joins in main memory brings considerable benefits in terms of scalability. We also identify a fragment of DCs called *symmetric constraints* that considerably extend previous fragments for which scalable detection techniques have been studied. We develop new algorithms for detecting and generating errors with symmetric constraints, and show that these algorithms have significantly better performance than the ones based on joins. Finally, we discuss the benefits of these optimizations when reasoning about error detectability and repairability.

- (iv) We present a comprehensive empirical evaluation of our error generator to test its scalability. Our experiments show that the error-generation engine takes less than a few minutes to complete tasks that require the execution of dozens of queries, even on DBs of millions of tuples. In addition, we discuss the relative incidence of the various activities related to error generation, namely identifying changes, applying changes to the DB, checking detectability and repairability.

- (v) The generator provides an important service for data cleaning researchers and practitioners, enabling them to more easily do robust evaluations of algorithms, and compare solutions. To demonstrate this, we present an empirical comparison of several data-repairing algorithms over BART data. We show a number of novel insights into these algorithms, their features and relative performance that could not have been shown with existing error generators.

Our error generator is open source and publicly available.² As the tool is used, we are also building up and hosting a shared library of error-generation tasks. We believe the system and shared error-generation tasks will raise the bar for evaluation standards in data cleaning.

The problem of error generation has points of contact with other topics in database research, such as the view-update problem [2], the problems of missing answers [25] and why-not provenance [24], and database testing [5]. However, studying the error-generation problem specifically, we have developed important new algorithms and optimizations that had not previously been identified in any of these related areas, and allow our system to scale to large error-generation tasks. These innovations will certainly have applicability to these related areas.

Proofs of the theorems, pseudocode, and additional details about the system are in the appendix.

2. CONSTRAINTS AND VIOLATIONS

We assume the standard definition of a relational database schema \mathbf{S} , instance I , and tuple t . In addition, we assume the presence of *unique tuple identifiers* in an instance. That is, t_{id} denotes the tuple with identifier (id) “id” in I . A *cell* in I is a tuple element specified by a tuple id and attribute pair $\langle t_{id}, A_i \rangle$. The *value* of a cell $\langle t_{id}, A_i \rangle$ in I is the value of attribute A_i in tuple t_{id} . As an alternative notation we use $t_{id}.A_i$ for a cell $\langle t_{id}, A_i \rangle$. We call a *relational atom* over a schema \mathbf{S} a formula of the form $R(\bar{x})$ where $R \in \mathbf{S}$ and \bar{x} is a tuple of (not necessarily distinct) variables. A *comparison atom* is a formula of the form $v_1 \text{ op } v_2$, where v_1 is a variable and v_2 is either a variable or a constant, and op is one of the following comparison operators: $=, \neq, >, <, \leq, \geq$.

²<https://github.com/dbunibas/BART>

We use the language of denial constraints (DC) [31] to specify cleaning rules. We introduce a *normal form* for DCs:

Definition 1: Denial Constraint – A *denial constraint* (DC) in normal form over schema \mathbf{S} is a formula of the form

$$\forall \bar{x}_1, \dots, \bar{x}_n : \neg(R_1(\bar{x}_1), \dots, R_n(\bar{x}_n), \bigwedge_{i=1}^m (v_i \text{ op } v'_i))$$

such that (i) $R_1(\bar{x}_1), \dots, R_n(\bar{x}_n)$ are relational atoms, where variables are not reused within or between atoms, and (ii) $\bigwedge_{i=1}^m (v_i \text{ op } v'_i)$ is a conjunction of comparison atoms, each of the form $x_k \text{ op } x_l$ or $x_k \text{ op } c$, where $c \in \text{Consts}$. \square

Example 2: Continuing with Example 1, we introduce the following sample declarative constraints.

- (i) We already discussed FD $d_1 : \text{Emp} : \text{Name} \rightarrow \text{Dept}$.
- (ii) Conditional FD (CFD) $d_2 : \text{Emp} : \text{Dept}[\text{“Sales”}], \text{Name} \rightarrow \text{Mng}$ states that $\text{Name} \rightarrow \text{Mng}$ holds only for “Sales” employees.
- (iii) A *standardization rule* over table Emp states that all employees working in department “Staff” must have a salary of one-thousand dollars. We model this using a (single-tuple) CFD $d_3 : \text{Emp} : \text{Dept}[\text{“Staff”}] \rightarrow \text{Sal}[\text{“1000”}]$.
- (iv) Consider now an *editing rule* [17] d_4 prescribing how to repair inconsistencies over table Emp based on master-data values from table MD : “when Emp and MD agree on Name and Dept , change Emp.Mng to the value taken from MD.Mng ”. This rule requires that the MD table can not be changed, i.e., the MD table is *immutable*.
- (v) Finally, we consider an *ordering constraint* d_5 stating that any manager should have a salary that is not lower than the salaries of his or her employees.

Next we encode our sample constraints d_1 – d_5 as a set of DCs (universal quantifiers are omitted).

$$\begin{aligned} \text{dc}_1 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d \neq d') \\ \text{dc}_2 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d = d', \\ &\quad d = \text{“Sales”}, m \neq m') \\ \text{dc}_3 &: \neg(\text{Emp}(n, d, s, m), d = \text{“Staff”}, s \neq \text{“1000”}) \\ \text{dc}_4 &: \neg(\text{Emp}(n, d, s, m), \text{MD}(n', d', m'), n = n', d = d', m \neq m') \\ \text{dc}_5 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(m', d', s', m''), m = m', s' < s) \quad \square \end{aligned}$$

To formalize the semantics of a DC, we introduce the notion of an *assignment* to the variables of a logical formula $\phi(\bar{x})$. DCs in normal form have the nice property that each variable has a unique occurrence within a relational atom. We can therefore define an *assignment* as a mapping m of the variables \bar{x} to the cells of I . We formalize the notion of a *violation* for a constraint using assignments.

Definition 2: Violation – Given a constraint $\text{dc} : \forall \bar{x} : \neg(\phi(\bar{x}))$ over schema \mathbf{S} , and an instance I of \mathbf{S} . We say that dc is *violated* in I if there exists an assignment m such that $I \models \phi(m(\bar{x}))$. For a set of DCs Σ , $I \models \Sigma$ if none of the constraints in Σ is violated in I . \square

Consider constraint dc_1 in Example 2. Assume two tuples are present in the instance of Emp : $t_1 = \text{Emp}(\text{John}, \text{Staff}, 1000, \text{Mark})$, $t_6 = \text{Emp}(\text{John}, \text{Sales}, 1000, \text{Mark})$. These tuples represent a violation of dc_1 according to assignment m that maps variables to cells and atoms to tuples as follows (we omit s, s', m, m'):

$$\begin{aligned} m(n) &= t_1.\text{Name} & m(d) &= t_1.\text{Dept} & m(\text{Emp}(n, d, s, m)) &= t_1 \\ m(n') &= t_6.\text{Name} & m(d') &= t_6.\text{Dept} & m(\text{Emp}(n', d', s', m')) &= t_6 \end{aligned}$$

To determine if a constraint is violated, we must examine values assigned to variables used in comparison atoms. To emphasize this, we call the *context variables* of a constraint dc in normal form those variables that appear in comparison atoms (in our example, $\{n, n', d, d'\}$). A *context* for a violation is the set of cells associated with context variables (in our example, $\{t_1.\text{Name}, t_1.\text{Dept}, t_6.\text{Name}, t_6.\text{Dept}\}$).

Definition 3: Violation Context – Given a constraint $\text{dc} : \forall \bar{x} : \neg(\phi(\bar{x}))$ over schema \mathbf{S} in normal form, and an instance I of \mathbf{S} , assume I contains a violation of dc according to assignment m . The *violation context* for dc and m , denoted by $\text{vio-context}_{\text{dc}}(m)$, is the set of cells that are the image according to m of the context variables of dc . \square

3. PROPERTIES OF ERRORS

Our primary concern is to provide users with fine-grained control over the error-generation process. To do so, we concentrate on two important properties of the changes we make to the clean instance, namely *detectability* and *repairability*.

Detectability. We define precisely when a cell change ch introduces an error that is detectable using Σ . We say that cell $t_i.A$ is *involved in a violation* with dc in instance I if there exists an assignment m such that I violates dc according to m , and $t_i.A \in \text{vio-context}_{\text{dc}}(m)$.

Definition 4: Detectable Errors – Consider a set of cell changes Ch to instance I that yield a new instance $I_d = \text{Ch}(I)$. We say that a cell change $\text{ch} = \langle t_i.A := v \rangle$ in Ch introduces a *detectable error* for constraint dc if: (i) cell $t_i.A$ is *involved in a violation* with dc in instance I_d and (ii) cell $t_i.A$ was not involved in a violation with dc in instance $I' = \text{Ch}'(I)$, where $\text{Ch}' = \text{Ch} - \{\text{ch}\}$. \square

The cell-change $\text{ch}_1 = \langle t_2.\text{Name} := \text{“John”} \rangle$ from Example 1 introduces a detectable error for dc_1 . Notice that it does not introduce a violation to any of dc_2 – dc_5 . Hence, we say that ch_1 is detectable by *exactly-one* constraint in $\Sigma = \{\text{dc}_1, \dots, \text{dc}_5\}$. This is not always the case. An example is $\text{ch}_2 = \langle t_4.\text{Mgr} = \text{“John”} \rangle$. This change is detectable using both dc_5 (which states that a manager must have a greater salary than his or her employees), and dc_4 (which states that employees who are present in the master-data table, must have the manager indicated by the master data). We say that ch_2 is *at-least-one* detectable, or simply detectable.

Reasoning about detectability is important. There are, in fact, repair algorithms [10, 27] that exploit simultaneous violations to multiple constraints to select among repairs. From the algorithmic viewpoint, given a set of constraints Σ , we want our error-generation method to control the detectability of errors, and whether they are *exactly-one* or *at-least-one* detectable.

Repairability. The notion of *repairability* provides an estimate of how hard it is to restore a DB with errors to its original, clean state. It is clear that such an estimate depends on the actual algorithm used to repair the data. To propose a practical notion, we concentrate on a specific class of repair algorithms. First, we restrict our attention to repair algorithms that use *value modification*, i.e., they fix violations by modifying one or more values in the dirty DB. Second, we assume that these algorithms rely on a *minimality principle*, according to which repairs that minimally modify the given DB are to be preferred. We notice that

the vast majority of algorithms in the recent literature [3, 6, 7, 10, 17, 20, 27, 31, 35] fall in this category.

Consider again Example 1, and assume instance I_d only has the following tuples t_{10} - t_{14} :

$t_{10} : \text{Emp}(\text{John}, \text{Staff}, \dots)$ $t_{11} : \text{Emp}(\text{John}, \text{Sales}, \dots)$
 $t_{12} : \text{Emp}(\text{John}, \text{Sales}, \dots)$ $t_{13} : \text{Emp}(\text{John}, \text{Mktg}, \dots)$
 $t_{14} : \text{Emp}(\text{John}, \text{Mktg}, \dots)$

there are eight violation contexts for constraint dc_1 , each one using two tuples with equal name and different department. However, an algorithm that is inspired by the principle of minimal repairs would be primarily interested in knowing that these 8 contexts reduce to a group of 10 cells (the name and dept cells of tuples 10-14). Variants of this notion, initially called *equivalence class* [7], have been formalized under the names of *homomorphism class* [20] and *repair context* [10]. We call this a *context equivalence class*.

Definition 5: Context Equivalence Class – Given a constraint dc over schema S in normal form, and an instance I of S , we say that two cells of I *share a context* if they share a violation context for dc and I . A *context equivalence class* \mathcal{E} for dc and I is an equivalence class of cells induced by the transitive closure of the relation “share a context”. \square

It is natural to reason about the repairability of a violation based on its context equivalence class. To formalize the notion of repairability, we start from a cell change $\text{ch} = \langle t.A := v_d \rangle$ introducing a detectable error to dc . We assume we have identified the context equivalence class \mathcal{E} for cell $t.A$ and dc . From this, we derive a bag of candidate values, denoted by *candidate-values*($t.A, \mathcal{E}, \text{dc}$). Then, we define the repairability as the probability of restoring $t.A$ to its original value by uniformly random picking a value from *candidate-values*($t.A, \mathcal{E}, \text{dc}$):

Definition 6: Repairable Error and Repairability – Given a violation $\text{ch} = \langle t.A := v_d \rangle$ to constraint dc that changes the cell $t.A$ in I from a clean value v_c to a dirty v_d , call \mathcal{E} the context equivalence class of cell $t.A$ and dc . We denote by $\mathcal{V} = \text{candidate-values}(t.A, \mathcal{E}, \text{dc})$ the bag of *candidate repair values* from \mathcal{E} for $t.A$ and dc . We say ch is a *repairable error* if $v_c \in \mathcal{V}$. The *repairability* of ch is computed by dividing the number of occurrences of v_c in \mathcal{V} by the size of \mathcal{V} . \square

The definition of function *candidate-values*($t.A, \mathcal{E}, \text{dc}$) is elaborate but quite standard in data repairing. For the sake or readability, we introduce it by means of examples:

(i) consider our sample equivalence class for tuples t_{10} - t_{14} ; for FD dc_1 in Example 2, and change $t_{11}.\text{Dept}$ from “Sales” to “Staff”, we select all values from cells of the Dept attribute in the equivalence class: {Staff, Staff, Sales, Mktg, Mktg}. The error is repairable, and the repairability is $1/5 = 0.2$;

(ii) for a CFD, like dc_3 in Example 2 (employees of the staff department need to have salaries of \$1000), and change $t_1.\text{Salary}$ from 1000 to 2000, the only candidate value is dictated by the constant, and is exactly 1000; therefore, the repairability is 1; similarly for fixing rules; editing rules like dc_4 in Example 2 use master-data tuples, i.e., immutable cells, and are treated accordingly: candidate values are taken from the master-data cells only;

(iii) finally, consider ordering constraints like dc_5 in Example 2 (managers have higher salaries than their employees); assume we change Paul’s salary to 2000 in t_2 to make it higher than the one of its manager; it is easy to see that

there are infinite real values that can be used to repair the violation; in this case, the repairability is 0.

Notice that a change may be detectable by several constraints, and thus have different repairability values. In this case, we consider the maximum of these values.

4. FORMALIZATION OF THE PROBLEM

In this section we formalize the error generation problem.

Recall that, given an instance I of S and a set Σ of DCs, to detect violations we need to find assignments that generate violation contexts. Notice this problem can be solved by running a *violation-detection query* for dc and I . Recall that id is an attribute storing the tuple id of a relation. For a constraint dc with multiple relation atoms, we abuse notation by using \vec{id} in queries to denote a vector of variables bound to all tuple ids from these atoms. Given a DC of the form $\text{dc} : \neg(\phi(\vec{id}, \vec{x}))$, we denote the context variables (i.e., variables that appear in a comparison atom) of dc by \vec{x}_c , and the rest by \vec{x}_{nc} .

Definition 7: Vio-Detection Query – The *violation-detection (vio-detection) query* for dc is a conjunctive query with free variables \vec{id}, \vec{x}_c of the form: $DQ_{\text{dc}}(\vec{id}, \vec{x}_c) = \phi(\vec{id}, \vec{x}_c, \vec{x}_{nc})$. \square

Consider our example dc_1 , its vio-detection query is the following (notice how we add the predicate $\text{id} < \text{id}'$ to avoid returning every pair of tuples twice):

$$DQ_{\text{dc}_1}(\text{id}, \text{id}', n, n', d, d') = \text{Emp}(\text{id}, n, d, s, m), \\ \text{Emp}(\text{id}', n', d', s', m'), n = n', d \neq d', \text{id} < \text{id}'$$

Injecting detectable errors into a clean DB requires the identification of cells that can be changed in order to trigger violations. To find cells with this property, we notice that each comparison in the normal form of a DC suggests a different strategy to generate errors.

Consider for example constraint dc_1 in our example: $\text{dc}_1 : \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d \neq d')$. It states that there cannot exist two employees with equal names and different departments. Given a clean DB, there are two ways to change the DB to violate this constraint:

(i) *enforce the inequality*: we may find tuples with equal names, and equal departments ($n = n', d = d'$), and change one of the departments in such a way that they become different (d becomes different from d');

(ii) *enforce the equality*: or we may find tuples with different names, and different departments ($n \neq n', d \neq d'$) and change one of the names in such a way that n becomes equal to n' .

Based on this intuition, we introduce the notion of a *violation-generating query* (vio-gen query in the following). These queries are obtained from a DC, or better, the violation-detection query for it, by changing one of the comparisons into its negation:

Definition 8: Vio-Gen Queries – Given a constraint dc in normal form, consider its vio-detection query, DQ_{dc} . A *vio-gen query* $GQ_{\text{dc}, i}$ for dc is a query obtained from DQ_{dc} by changing a comparison atom of the form $v_1^i \text{op} v_2^i$ into its negation, $\neg(v_1^i \text{op} v_2^i)$. The latter is called the *target comparison* of the query. \square

In our example, we have two vio-gen queries, as follows (target comparisons are enclosed in brackets):

$$\begin{aligned}
GQ_{dc_1,1}(id, id', n, n', d, d') &= \text{Emp}(id, n, d, s, m), \\
&\quad \text{Emp}(id', n', d', s', m'), n = n', (d = d'), id < id' \\
GQ_{dc_1,2}(id, id', n, n', d, d') &= \text{Emp}(id, n, d, s, m), \\
&\quad \text{Emp}(id', n', d', s', m'), (n \neq n'), d \neq d', id < id'
\end{aligned}$$

Definition 9: Vio-Gen Cell – Given an instance I , we call a *vio-gen cell* for vio-gen query $GQ_{dc,i}$ any cell in the result of $GQ_{dc,i}$ over I that is associated with a variable in the target comparison of the query. \square

Consider, our DB in Example 1. The two vio-gen queries for dc_1 identify cells $t_4.\text{Dept}$ and $t_1.\text{Name}$ that can be used to inject errors, as follows:

$GQ_{dc_1,1} \text{ result}$	<i>change</i>	<i>after change</i>
$t_4 : \text{Emp}(\text{Frank}, \text{Mktg}, \dots)$	$t_4.\text{Dept} := \text{"xxx"}$	$t_4 : \text{Emp}(\text{Frank}, \text{xxx}, \dots)$
$t_5 : \text{Emp}(\text{Frank}, \text{Mktg}, \dots)$		$t_5 : \text{Emp}(\text{Frank}, \text{Mktg}, \dots)$

$GQ_{dc_1,2} \text{ result}$	<i>change</i>	<i>after change</i>
$t_1 : \text{Emp}(\text{John}, \text{Staff}, \dots)$	$t_1.\text{Name} := \text{"Paul"}$	$t_1 : \text{Emp}(\text{Paul}, \text{Staff}, \dots)$
$t_2 : \text{Emp}(\text{Paul}, \text{Sales}, \dots)$		$t_2 : \text{Emp}(\text{Paul}, \text{Sales}, \dots)$

Query $GQ_{dc_1,1}$ identifies cells $t_4.\text{Dept} = t_5.\text{Dept} = \text{"Mktg"}$. By making these cells different, we are guaranteed to introduce a violation (the tuples have equal names). Similarly, query $GQ_{dc_1,2}$ captures the fact that we can introduce a detectable violation by equating cells $t_1.\text{Name} = \text{"John"}$, $t_3.\text{Name} = \text{"Paul"}$ (since they have different departments).

A few remarks are in place. First, after we have identified a vio-gen cell, then we need to identify an appropriate value to generate the actual cell change and update the DB. Second, vio-gen queries also identify a context for each vio-gen cell, i.e., a set of cells that are the image of variables in comparison atoms. For example, the context of cell $t_4.\text{Dept}$ is composed by cells $\{t_4.\text{Dept}, t_4.\text{Name}, t_5.\text{Dept}, t_5.\text{Name}\}$. We need to keep track of these contexts to avoid that new changes are accidentally repairing other violations. These aspects will be discussed in the next section.

We now define the *error-generation problem* for an error-generation task $\mathbf{E} = \langle S, \Sigma, I, \text{Conf} \rangle$, as introduced in Section 1. For each constraint $dc \in \Sigma$, a parameter $\epsilon(dc)$ in Conf determines the number of detectable errors that should be introduced in I .

Definition 10: Error-Generation Problem – Given an *error-generation task* $\mathbf{E} = \langle S, \Sigma, I, \text{Conf} \rangle$, generate *at-least* (resp. *exactly*) $\epsilon(dc)$ detectable errors in I for each $dc \in \Sigma$. \square

It turns out that both variants of the error-generation problem are NP-complete.

Theorem 1: *The exactly-one and at-least-one error-generation problems for a task \mathbf{E} are NP-complete.*

Before we present algorithms and optimizations to solve the error-generation problem, we observe that the vio-gen queries for a constraint dc are nothing else than vio-detection queries for DCs that can be considered as being “*dual*” to dc . Consider our example about dc_1 . It can be seen that the two vio-gen queries correspond to the detection queries of constraints dc_1^1, dc_1^2 , as follows:

$$\begin{aligned}
dc_1^1 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d \neq d') \\
dc_1^2 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d = d') \\
dc_1^3 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n \neq n', d \neq d')
\end{aligned}$$

Since DCs are closed with respect to the negation of comparison atoms, this is a general property: any vio-gen query for constraint dc coincides with the vio-detection query of

a constraint dc' . This highlights the *duality* between the queries needed to inject new, detectable errors, and those needed to identify existing errors.

5. THE VIO-GEN ALGORITHM

In this section, we develop a *vio-gen* algorithm to solve the error-generation problem.

To efficiently generate large erroneous datasets, we aim for a solution that is in PTIME and scales to large databases. Our algorithm is a correct, but not complete, solution for the *at-least-one* error-generation problem. That is, when it succeeds, it returns a solution to the input task.

In the algorithm, we use vio-gen queries to identify cells to change; also, we avoid interactions between cell changes by enforcing that two cell changes cannot share a context. Thus, our algorithm may fail to find a solution if there is no solution without shared contexts. However, as long as there are sufficient vio-gen cells available, and a sufficient number of non-overlapping contexts, our algorithm will succeed. Intuitively, whether this is the case depends mainly on the requested error ratios (and to a lesser extent on the shape of the constraints); in practice, typical error ratios are in the range 1%-10% wrt the size of the input DB, and therefore the probability of success is very high (e.g., we succeed for all scenarios we considered in our experimental evaluation).

The main tasks are as follows.

Task 1: Finding Vio-Gen Cells and Contexts. Given dc and I , we generate the vio-gen queries for dc . Then, we execute each vio-gen query over I to identify a set of vio-gen cells (possibly all), and the associated violation contexts.

Task 2: Value Whitelists and Blacklists. Once we have identified a vio-gen cell, we need to find a set of values that would inject errors into the cell. We reason about these values in terms of a *value whitelist*, i.e., values that are candidates to be used for that cell, and a *value blacklist*, i.e., values that cannot be used. Intuitively, the set of values for a cell will be obtained by the set difference between the whitelist and the blacklist for that cell. These depend on the constraint and on the target comparison. In our previous examples:

(i) the target comparison of query $GQ_{dc_1,1}$ is an equality ($d = d'$), and we want to generate changes that falsify the equality; therefore, once we have identified cell $t_4.\text{Dept}$ and its value, “*Mktg*”, the whitelist contains any string – denoted by “***” – and the blacklist is the single value {“*Mktg*”};

(ii) the target comparison of query $GQ_{dc_1,2}$ is a difference ($n \neq n'$), and our changes should falsify it; therefore, once we have identified cell $t_1.\text{Name}$ with value “*John*”, and the value “*Paul*” in another cell $t_2.\text{Name}$ which shares a violation context with $t_1.\text{Name}$, the whitelist for cell $t_1.\text{Name}$ contains {“*Paul*”}, the blacklist contains the original value {“*John*”}.

The process of identifying values to change vio-gen cells has some additional subtleties. Consider query $GQ_{dc_2,3}$ for constraint dc_2 in our running example:

$$\begin{aligned}
GQ_{dc_2,3}(id, id', \dots) &= \text{Emp}(id, n, d, s, m), \text{Emp}(id', n', d', s', m'), \\
&\quad n = n', (d \neq d'), d = \text{"Sales"}, m \neq m'
\end{aligned}$$

Here, the value of variable d in the target comparison is constrained to be “*Sales*”. As a consequence, we cannot change its value. This general rule is necessary for detectability: any change we make to the DB for a vio-gen query must break the target comparison but preserve all other comparisons within the query. Therefore, we can only change the

value of d' and make it equal to “Sales”, that is, the whitelist of vio-gen cells for d' will only contain that value.

In the general case, our algorithm determines valid values for vio-gen cells by computing equivalence classes based on equality comparisons and use an interval algebra to deal with ordering constraints.

Task 3: Handling Interactions. The purpose of this step is to avoid possible interference among different changes to the DB. To see the problem, following our discussion of dc_1 , suppose vio-gen query $GQ_{dc_1,1}$ suggests a change to cell t_4 .Dept to make it equal to “xxx”. Assume after that query $GQ_{dc_1,2}$ suggests a change to the Name cell of the same tuple, to make it equal to “Paul”. If we apply both changes to the DB, we get the following instance:

clean DB	change	after changes
t_4 :Emp(Frank, Mktg, ...)	t_4 .Dept	t_4 :Emp(Paul, xxx, ...)
t_5 :Emp(Frank, Mktg, ...)	:= “xxx”	t_5 :Emp(Frank, Mktg, ...)
	t_4 .Name	
	:= “Paul”	

It is easy to see that the instance obtained after these two changes is not dirty wrt dc_1 , since the second change is removing the detectable error introduced by the first one.

Our algorithm, while generating changes also stores their violation contexts in main memory. We discard a candidate change that may remove violations introduced by previous changes. More precisely, a vio-gen cell is changed only if: (i) the cell itself does not belong to the context of any other previous change; (ii) none of the cells in its context have been changed by previous changes.

Task 4: Generating Random Changes. BART tries to randomly generate the desired number of changes as specified as a desired percentage of detectable errors for a constraint and an attribute in the task configuration. Recall that error-percentages are expressed wrt the size of the DB table, e.g., 1% of errors in the cells of attribute A in a table R of 100K tuples corresponds to 1000 errors.

A crucial requirement is that, during this process, vio-gen cells that may be needed later on are not discarded. In fact, we don’t know if the clean DB contains a sufficient number of non-interacting vio-gen cells. Suppose we ask for 1000 errors for an attribute, but there are less vio-gen cells. In this case, BART will generate as many changes as it is possible, and log a warning about the impossibility of meeting the configuration parameter. To handle this, we resorted to a main-memory sampling algorithm. The algorithm works as follows. Assume we are required to generate n changes for vio-gen query GQ and attribute R.A:

- (i) before executing GQ , we select a random probability p_{GQ} within a range that can be configured by the user (usually 10% to 50%), and a random offset o_{GQ} , again randomly generated in a range from 1% to 10% of the size of R; these parameters will be used to sample cells at query execution;
- (ii) then, we execute the GQ , and scan its result set; for each tuple t , we extract the vio-gen cell and vio-gen context, and check if it overlaps with the ones generated so far; in this case, we discard the tuple; otherwise, we consider the tuple as a candidate to generate a detectable change;
- (iii) a candidate tuple can be either processed or skipped; notice however that tuples that are skipped are not lost, but rather stored in a queue of candidate tuples, in case we need them later on; first, we skip the first o_{GQ} tuples that would generate valid vio-gen contexts; after this, for each candidate tuple we draw a random number $random(t)$: we process the

tuple and generate a new change whenever $random(t) < p_{GQ}$; otherwise, we store the tuple in the candidate queue;

(iv) the process stops as soon as we have generated n changes; if, however, no more tuples are available in the query result before this happens, then we go back to the queue of candidate tuples, and iterate the process.

We now state the following correctness result.

Theorem 2: *Given an error-generation task $\mathbf{E} = \langle S, \Sigma, I, Conf \rangle$, call Ch the output of the vio-gen algorithm over \mathbf{E} . If the algorithm succeeds, then Ch is an at-least-one solution for task \mathbf{E} .*

The cost of our Vio-gen algorithm is dominated by query execution. It runs in PTIME in the size of the instance and the number of constraints and is NP-hard in the maximal size of the constraint (number of atoms).

Theorem 3: *Let $\mathbf{E} = \langle S, \Sigma, I, Conf \rangle$ be an error-generation task. The Vio-gen algorithm runs in PTIME in $\|I\|$ and $\|\Sigma\|$ and is NP-hard in the maximal size (number of atoms) of constraints in $\|\Sigma\|$.*

Note that this is the worst-case complexity of the algorithm. As demonstrated by our experimental evaluation (Section 8, the algorithm perform well in practice, because using sampling we seldom have to fully evaluate a vio-gen query and the size of constraints is typically limited to a few relational atoms (i.e., the NP-hardness does not effect runtime significantly for most datasets).

Task 5: Checking Detectability and Repairability. Theorem 2 guarantees the vio-gen algorithm generates detectable changes. Our algorithm only guarantees *at-least-one* detectability, but may violate multiple constraints. If requested by a user, the system will compute how many constraints are violated by a change. This can be used to output *exactly-one* detectable changes. In addition, it may also compute repairability and filter solutions to guarantee a level of repairability if specified in the task configuration (Definition 6). This is done as follows:

- (i) after changes have been generated and applied to yield the dirty DB I_d , we run the detection query on I_d for each constraint dc in Σ to find all violation contexts for dc and I_d ; recall that this task is essentially the same as Task 1 (vio-gen queries and vio-detection queries are structurally identical); we keep counters for each cell change to know how many violation contexts the cell change occurs in, and how many constraints it violates;
- (ii) we compute repairability based on the context equivalence classes (optimizations are discussed in Section 6.2).

6. OPTIMIZATIONS

We now discuss the scalability of the vio-gen algorithm. The algorithm is designed to solve all subtasks of the error-generation problem, i.e., find n changes with at-least-one detectability, and compute detectability and repairability measures. The main cost factor is related to executing vio-gen queries (in turn, vio-detection queries when measuring detectability and repairability). These techniques do not scale to large DBs, for two main reasons.

Problem 1: Computing Cross-Products. Some constraints may result in vio-generation queries with no equalities. Such queries would end up to be executed as cross-products, and therefore have inherent quadratic cost. As an

example, consider query $GQ_{dc_5,1}$ for constraint dc_5 :

$$GQ_{dc_5,1}(id, id', \dots) = \text{Emp}(id, n, d, s, m), \text{Emp}(id', n', d', s', m'), \\ (m \neq n'), s' < s$$

This query has a very common pattern, one we have already observed in the vio-gen queries of constraint dc_1 .

Problem 2: Computing Expensive Joins. Even when equalities are present, and joins are performed, this may be slow. In fact, the distribution of values in the DB might lead to join results of very large cardinality. Consider an FD $d : \text{Emp} : \text{Dept} \rightarrow \text{Mngr}$, stating that department names imply manager names. One of the vio-gen queries for d would be:

$$GQ_{d,1}(id, id', \dots) = \text{Emp}(id, n, d, s, m), \text{Emp}(id', n', d', s', m'), \\ (d = d'), (m = m'), id < id'$$

As we execute this query on a clean instance I , where equal departments correspond to equal managers, the result of this join may be much larger than the size of table Emp (e.g., when there are few departments).

These problems illustrate that the evaluation of vio-gen queries to generate errors may not scale well to large DBs. We introduce two important optimizations to greatly improve performance. We first exploit the fact that we only have to produce a certain percentage of errors by avoiding the full execution of vio-gen queries involving cross-products. Second, we identify a class of DCs, called symmetric constraints, for which the vio-gen queries can be significantly simplified.

6.1 Sampling Cross-Products

For testing data cleaning algorithms, often we want to generate a set of errors that is (much) smaller than all errors that could potentially be introduced. Thus, we consider sampling tables, and then computing cross-products over these samples in main memory as an alternative to computing the cross-product using a declarative query. To understand our intuition, consider the typical vio-gen query with inequalities for an FD, for example dc_1 :

$$GQ_{dc_1,2}(id, id', \dots) = \text{Emp}(id, n, d, s, m), \text{Emp}(id', n', d', s', m'), \\ (n \neq n'), d \neq d', id < id'$$

We search for tuples with different names and departments. In any DB with a sufficiently large number of tuples, we should have a good probability of finding tuples that differ from each other. Whenever we need to generate n changes:

(i) we scan the tables in the query to extract a sample of $c \cdot n$ tuples (c being a configuration parameter), and materialize them in memory;

(ii) we compute the cross-product in memory, filter results according to the comparisons ($(n \neq n'), d \neq d'$), and use the results for identifying vio-gen cells and their contexts; we stop as soon as n contexts have been found;

(iii) if we are not able to find n non-overlapping contexts, then we repeat the process, i.e., we choose a random offset, re-sample the tables and iterate; an iteration limit is used to stop the process before it becomes too expensive.

This strategy has a worst-case cost that is comparable (or even worse) to that of computing the whole cross-product. To avoid this cost, we limit the number of iterations. This may only happen when the cross-product is empty, or has very few results, both very unlikely cases. Our experiments confirm that typically, our optimized version runs orders of magnitude faster than computing the cross-product using

the DBMS, even when using the LIMIT clause to reduce the number of results.

6.2 Symmetric Constraints

It is known that it is possible to find tuples that violate CFDs by running join-free queries using the group-by operator [6, 15]. In this section, we build on this result, and extend it in several ways:

(i) we generalize Bohannon et al.'s [6] treatment of CFDs by formalizing the notion of *symmetric denial constraints*; this significantly enlarges the class of constraints for which this optimization can be adopted;

(ii) we develop a general algorithm to optimize the execution of symmetric queries with at least one equality and arbitrary inequalities, and show how we can efficiently compute both violation contexts and context equivalence classes;

(iii) we conduct an extensive experimental evaluation to show the benefits of this optimization.

We now blur the distinction between DCs, vio-gen queries, and detection queries, and speak simply about their logical formulas. Consider a DC $d : \neg(\phi(\bar{x}))$ in normal form; assume that $\phi(\bar{x})$ has no ordering comparisons ($>$, $<$, \leq , \geq), and contains at least one equality. Intuitively, d is *symmetric* if it can be “broken up” in two isomorphic subformulas. To formalize this idea, we make use of following definition.

Definition 11: Formula Graph – Given $\neg(\phi(\bar{x}))$ in normal form, its *formula graph*, $\mathcal{G}(\phi(\bar{x}))$, is defined as follows:

(i) it has a node for every relational atom, comparison atom, variable and constant in $\phi(\bar{x})$;

(ii) a node for relational atom $R_i(\bar{x}_i)$ has label R_i ; a node for comparison $v \text{ op } v'$ has label op ; a node for variable x (constant c) has label x (c , respectively);

(iii) it has an edge between the node for atom $R_i(\bar{x}_i)$ and each node for a variable $x \in \bar{x}_i$; if x appears within attribute A , the edge is labeled by $R_i.A$;

(iv) there is an edge between the node for atom $v \text{ op } v'$ (v, v' either variable or constant) and the nodes for v, v' . \square

We are interested in subformulas that are isomorphic to each other according to some mapping of the variables, h . Thus, given h , we use it to break up the graph, and look for isomorphic connected components. Given a mapping of the variables $h : \bar{x} \rightarrow \bar{x}$, we define the *reduced formula graph* for h as the graph obtained from $\mathcal{G}(\phi(\bar{x}))$ by removing all nodes corresponding to comparisons of the form $x \text{ op } h(x)$.

Definition 12: Symmetric Formula – A formula in normal form $\neg(\phi(\bar{x}))$ with no ordering comparisons and at least one equality is *symmetric* wrt mapping h if its reduced connection graph contains exactly two connected components, and these are isomorphic to each other according to h , i.e., (i) a variable label v can be mapped to $h(v)$, (ii) every other label is preserved. \square

In our example, both dc_1 and dc_2 are symmetric constraints (with mapping $n \rightarrow n', d \rightarrow d', s \rightarrow s', m \rightarrow m'$):

$$dc_1 : \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d \neq d') \\ dc_2 : \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d = d', \\ d = \text{“Sales”}, d' = \text{“Sales”}, m \neq m')$$

The formula graph for dc_2 is depicted in Figure 2. Notice how, for the purpose of testing symmetry, we add the implied comparison atom $d' = \text{“Sales”}$.

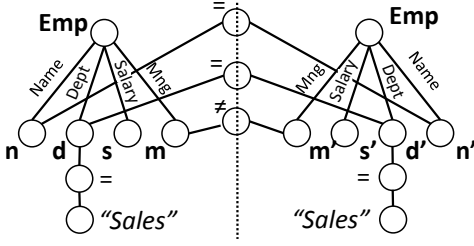


Figure 2: Formula Graph for dc_2

The vio-detection query for a symmetric constraint is always symmetric. However, not all vio-gen queries for a symmetric constraint need to be symmetric as some may not have equalities. For example, we know dc_1 is symmetric ($n = n', d \neq d'$), and among its vio-gen queries, one is symmetric ($n = n', d = d'$), the second one is not ($n \neq n', d \neq d'$).

6.3 Reduced Formulas

Symmetric constraints significantly reduce the number of joins required to execute the corresponding queries. The intuition is that we can only consider one of the isomorphic subqueries, and avoid redundant work. We represent these subcomponents of the original queries by means of relational atoms with *adornments*. An *adornment* for a variable $x \in \bar{x}_i$ is a superscript of the form $=$ or \neq , denoted by $x^=$ or x^{\neq} . We use adornments to keep track of the constraints that were imposed over variables within the original symmetric query.

Given a symmetric formula $\phi(\bar{x})$ in normal form, to optimize its execution we generate a *reduced formula with adornments*, denoted by $reduce(\phi(\bar{x}))$. Following are the reduced formulas for the symmetric vio-gen queries of dc_1, dc_2 (we use boldface to mark the variables that were involved in the target comparison of the original query):

$$\begin{aligned} reduce(GQ_{dc_1,1}) &= Emp(id, n^=, d^=, s, m) \\ reduce(GQ_{dc_2,1}) &= Emp(id, n^=, d^=, s, m^=), d = \text{"Sales"} \\ reduce(GQ_{dc_2,2}) &= Emp(id, n^{\neq}, d^=, s, m^{\neq}), d = \text{"Sales"} \\ reduce(GQ_{dc_2,3}) &= Emp(id, n^=, d^=, s, m^{\neq}), d \neq \text{"Sales"} \end{aligned}$$

Intuitively, we use $reduce(GQ_{dc_1,1})$ to derive a (join-free) SQL query that will give us all tuples t in Emp such that there exists a tuple t' with a different tuple id , equal name and equal department. Similarly for $reduce(GQ_{dc_2,1})$ (equal names, equal departments both "Sales", equal managers). Reduced formulas are constructed using the following algorithm (ids are treated separately):

- (i) start with $\phi(\bar{x})$ and variable mapping h ; consider the formula $\phi'(\bar{x}')$ corresponding to one of the connected components of the reduced formula graph for $\phi(\bar{x})$ and h ;
- (ii) for each comparison $x \text{ op } h(x)$ in $\phi(\bar{x})$, add to variable x in $\phi'(\bar{x}')$ adornment op .

Following is a more complex constraint that does not fall within classical ones. It states that a functional dependency $dc_6 : Emp : Name \rightarrow Manager$ holds for those tuples of Emp that correspond to master-data tuples on Name and Dept:

$$\begin{aligned} dc_6 : & \neg(Emp(n, d, s, m), MD(n', d', m'), n = n', d = d', \\ & Emp(n'', d'', s'', m''), MD(n''', d''', m'''), n'' = n''', d'' = d''', \\ & n = n'', m \neq m'') \end{aligned}$$

The constraint is symmetric and this is the reduced formula of one of its symmetric vio-gen queries (as usual, we explicitly mention id attributes in the formula) :

$$reduce(GQ_{dc_6,1}) = Emp(id, n^{\neq}, d, s, m^{\neq}), MD(id', n', d', s', m'), n = n', d = d'$$

6.4 The Benefits of Symmetry

Reduced formulas suggest an alternative query execution strategy that is based on a limited use of joins and favors group-by clauses. This strategy was previously used to find tuples involved in violations of (multi-tuple) CFDs [6]. We extend those algorithms to a larger class of constraints. Our work can handle multiple inequality adornments, while the original technique only considered one inequality at a time.

We first summarize the main ideas behind the use of group-by clauses for formulas with at most one inequality, then discuss the general case.

Simple Case: At Most One Inequality Adornment. Context equivalence classes for symmetric queries with at most one inequality can be computed by grouping tuples. Consider constraint dc_1 and the reduced formula for its symmetric vio-gen query:

$$reduce(GQ_{dc_2,1}) = Emp(id, n^=, d^=, s, m^{\neq}), d = \text{"Sales"}$$

In this case, a vio-gen context is made of two tuples with equal names and departments (equal to "Sales"), and different managers. A context equivalence class is composed by all tuples that have equal values of names and departments ("Sales"), such that there exist at least two different managers associated with them. Notice that this is a general property. To formalize this, we call the *equality variables* of the reduced formula the set of variables with equality adornments (these correspond to variables of the original formula involved in equalities).

Property 1: Given a symmetric constraint dc with at most one inequality, and instance I , two cells c_1, c_2 of I belong to the same context equivalence class for dc and I if and only if: (i) they belong to contexts vc_1, vc_2 for dc and I ; (ii) vc_1, vc_2 have equal values for the equality variables of dc . \square

To construct contexts and equivalence classes, we generate an SQL query with aggregates, denoted by $sym\text{-sql}(reduce(Q))$, from the reduced formula. In our example:

```
SELECT id, name, dept, mngr FROM emp WHERE name, dept IN
(SELECT name, dept FROM emp WHERE dept = 'Sales'
GROUP BY name, dept HAVING COUNT(DISTINCT mngr) > 1)
ORDER BY name, dept
```

Contexts and equivalence classes are built as follows:

- (i) we run query $sym\text{-sql}(reduce(Q))$;
- (ii) we group tuples in the result with equal name and dept attribute values, yielding all context equivalence classes;
- (iii) to construct the actual contexts, tuples within an equivalence class are combined in all possible ways in memory.

The General Case. Let us now consider a generic formula with adornments, with at least one equality and multiple inequalities, like the following:

$$reduce(GQ_{dc_2,2}) = Emp(id, n^{\neq}, d^=, s, m^{\neq}), d = \text{"Sales"}$$

Here, we are looking for pairs of tuples that have equal departments (with value "Sales"), and *both* different names, *and* different managers. Handling both inequalities makes the construction of contexts more complex. The intuition behind the algorithm is to execute multiple aggregates within our SQL query, one for each inequality, to identify cells that are candidates to generate violation contexts. In our example, $sym\text{-sql}(reduce(GQ_{dc_2,2}))$ would be the following query:

```
SELECT id, name, dept, mngr FROM emp
WHERE dept, name IN
(SELECT dept, name FROM emp
```



```

GROUP BY dept, name HAVING COUNT(DISTINCT name) > 1)
AND dept, mngr IN
(SELECT dept, mngr FROM emp
GROUP BY dept, mngr HAVING COUNT(DISTINCT mngr) > 1)
ORDER BY dept

```

Notice, however, that Property 1 does not hold in this case. Indeed, belonging to the result of this query is a necessary but not sufficient condition for a cell to belong to a violation context for $GQ_{dc_2,2}$. In fact, we have no guarantee that two tuples from the result of the SQL query above satisfy all inequalities. To select the cells that actually belong to contexts, we need to build the actual contexts, and keep only those in which all inequalities are satisfied.

A crucial optimization, here, is to select a relatively small set of candidate tuples for each context equivalence class. To do this, we group the tuples in the result of the query on the values of the equality variable (`dept` in our example). Then, we combine the candidate tuples within each set to generate the actual violation contexts.

Once the contexts have been generated, building the context equivalence classes is straightforward: it suffices to hash contexts based on the values of the equality variables, and compute equivalence classes from the cells in each bucket. In the next section we show experimentally that this strategy performs very well even for large DBs.

7. USE CASES OF THE TOOL

We now summarize the main use-cases of the tool.

Use-Case 1: Generate Detectable Errors. The main purpose of BART is to generate *constraint-induced errors*. We control an error-generation task \mathbf{E} by a set of configuration parameters `Conf`. Here are the main ones:

- (i) *authoritative sources*: names of DB relations that are to be considered *immutable* (empty by default).
- (ii) *error percentages*: desired degree of detectable errors for each vio-gen query. We specify percentages with respect to the number of tuples in a table (e.g., 1% errors in a table of 100K tuples requires to make dirty 1000 cells).
- (iii) *repairability range*: in addition, users may also specify a range of repairability values for each vio-gen query; BART will estimate the repairability of changes, and only generate errors with estimated repairability within that range. This simplifies the generation of error configurations with controlled repairability (e.g., low, medium, high repairability, as shown in our experiments).

Use-Case 2: Generate Random Errors. In addition to detectable errors, BART may also generate random errors of several kinds: *typos* (e.g., ‘database’), *duplicated values*, *bogus* or *null values* (e.g., ‘999’, ‘***’), *outliers* in numerical attributes, following the standard definition [28]. Random errors may be freely mixed with constraint-induced ones.

Use-Case 3: Computing Repairability. Given a set of constraints, Σ , and a dirty DB, I_d , BART can be used to compute the repairability of the violations in I_d wrt Σ , according to Definition 6. This can be done also in the case in which I_d was not generated by the tool itself.

Use-Case 4: Measuring Repair Quality. Our ultimate goal is to benchmark data-repairing algorithms over BART data. Consider we run some algorithm A over a dirty instance I_d as generated by BART, and we obtain a repaired instance $I_{rep,A}$ by a set Ch_A of cell changes, i.e.,

$I_{rep,A} = Ch_A(I_d)$. The tool adopts a natural strategy to measure the performance of A over a task \mathbf{E} .

We call Ch^{-1} the set of cell changes that are needed to bring I_d back to its original state, I . Since we assume that I , I_d and $I_{rep,A}$ all have the same set of tuple ids, we define the *quality* of A over \mathbf{E} as the *F-Measure* of the set Ch_A , measured with respect to Ch^{-1} . That is, we compute the precision and recall of A in fixing the errors introduced in the original clean instance I . The higher the *F-measure*, the closer $I_{rep,A}$ is to the original clean instance I .

Since data-repairing algorithms have used different metrics to measure the quality of repairs, BART has been designed to be flexible in this respect. Precision and recall may be computed using the following measures:

- (i) **Value**: we count the number of cells that have been restored to their original values;
- (ii) **Cell-Var**: in addition to cells restored to their original values, we count (with 0.5 score) the cells that have been correctly identified as erroneous, and changed to a variable;
- (iii) **Cell**: we count the cells that have been identified as erroneous, regardless of the value assigned by the algorithm.

8. EXPERIMENTAL RESULTS

We describe a detailed evaluation of the BART Java prototype over synthetic and real-world datasets. We ran experiments on a machine with 8GB RAM, 2.6 GHz Intel Core i7, MacOS 10.10, and PostgreSQL 9.3.

Tasks. We tested five tasks, based on synthetic and real datasets: (i) *Employees* is the running example used in the paper; (ii) *Customers* is a synthetic scenario from Geerts et al. [19]; (iii) *Tax* is a synthetic scenario from Fan et al. [16]; (iv) *Bus* is a real-world scenario from Dallachiesa et al. [13]; and (v) *Hospital* is a real-world scenario used in several data repairing papers (e.g., [13, 17, 18]).

Note that all datasets have different characteristics. *Hospital* and *Bus* have higher redundancy in their data. *Tax* and *Employees* are the only datasets with constraints containing ordering ($<$, $>$) comparisons. Some datasets have master-data and CFDs, while others have only FDs. All these differences help to validate our techniques.

Settings. To measure the scalability of error generation, we focus on execution times. Each task has been executed five times and we report the average execution time. Because our focus is on error generation for testing data-cleaning algorithms, we report our study on different % of required errors in a range (1% – 10%) in which most cleaning algorithms can perform reasonably well. To show the effectiveness of our optimization techniques for symmetric constraints, we compare them against the standard execution.

Scalability. We discuss BART’s execution times on synthetic data sets. Figures 3.a-c report the execution times over an increasing number of tuples for 1%, 5%, and 10% injected errors, respectively. As expected, execution times increase both with the size of the input (number of tuples) and the number of changes (% of required errors). Our system is very fast, taking at most 6.6 minutes for one million tuples in the worst scenario (10% errors). The same observations apply for real-world data, as we show in Figure 3.d. Notice that, without our optimizations, the execution of all scenarios exceeds the time threshold we set (30 minutes). For this reason we do not report the data.

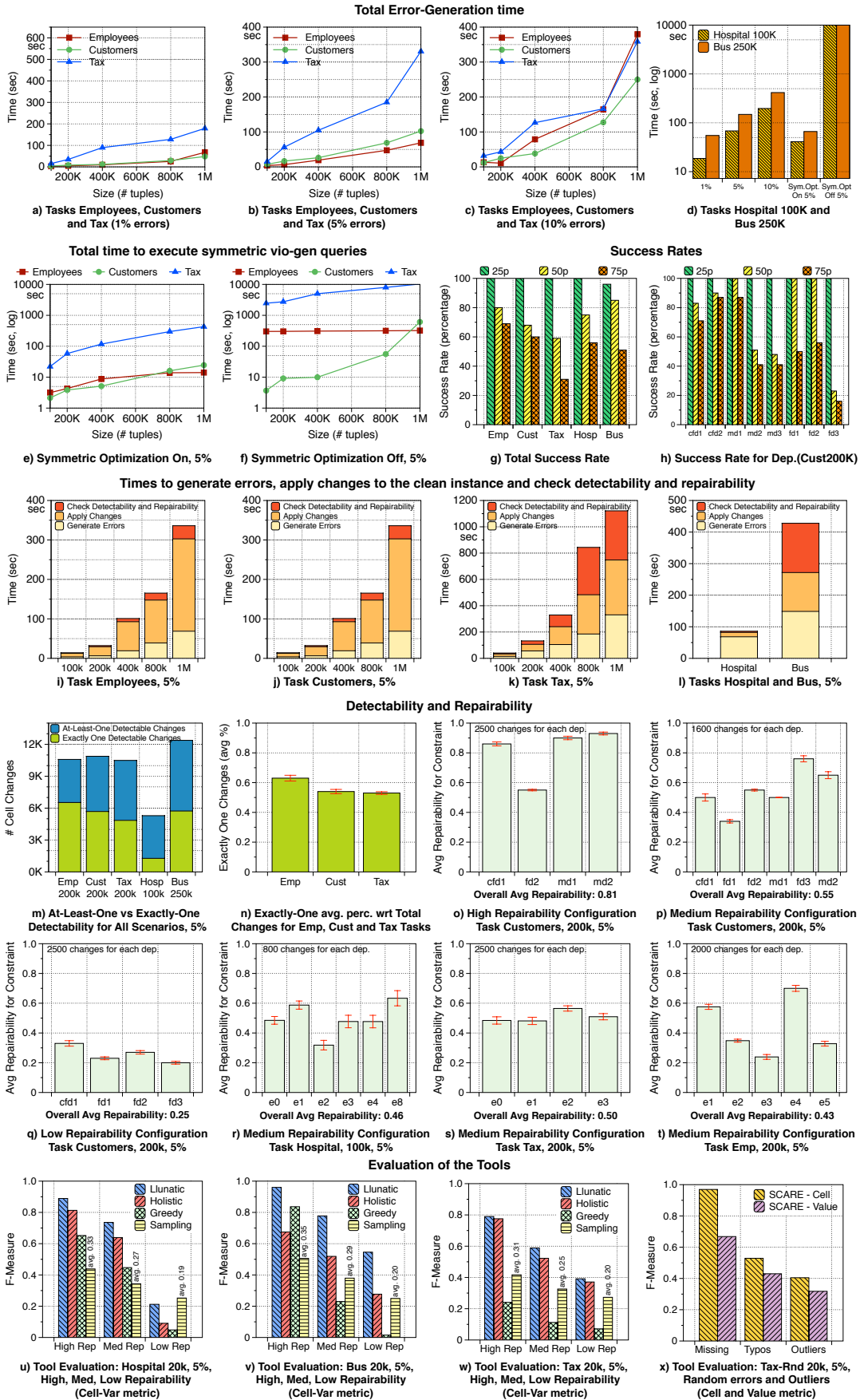


Figure 3: Experimental results

Datasets	Type and Size	#Tables	#Attributes	#Constraints	#Vio-Gen Queries	#Symmetric Queries	#Non Symmetric Queries	#Queries with No Equalities
Employees	Synthetic, 100K to 1M	2	7	5	12	4	5	3
Customers	Synthetic, 100K to 1M	3	16	8	22	7	12	3
Tax	Synthetic, 100K to 1M	1	13	5	18	12	2	4
Bus	Real 250K	1	5	12	24	10	2	12
Hospital	Real, 100K	1	19	7	16	11	0	5

Figure 4: Datasets and Tasks

Symmetric Queries. For each scenario, we extracted all symmetric vio-gen queries and executed them with and without our optimizations (we report the number of symmetric queries per scenario in Figure 4). Figure 3.e and f report the execution times over an increasing number of tuples and 5% injected errors for the synthetic datasets, with and without the symmetric optimization, respectively. Figure 3.d reports results for real datasets. For these experiments, we fixed a timeout of 10 minutes for each vio-gen query, after which we stopped the computation. Note how the symmetric strategy scales linearly with the size of the data, while the same observation does not hold for the standard join-based strategy (the DBMS adopts a nested-loop execution plan). Our symmetric optimization runs up to two orders of magnitude faster on these datasets.

Success Rate. Our algorithm trades completeness for scalability. It may fail to return the required number of changes, even if these exists, because it non-optimally uses violation contexts. To gain more insight on this aspect, we studied the behavior of the algorithm when 25%, 50%, and 75% of the size of the db is required to be made dirty. In Figure 3.g we report the fraction of changes that have been generated. Figure 3.h reports the breakdown of the success rate for every constraint of Customers. To limit the incidence of the distribution of data on the experiment, for each constraint we first computed an estimate of the maximal number of detectable changes that can be generated. Then, when running the actual experiments, we made sure to never ask for a number of changes per constraint higher than this estimate. It is interesting to note that in all scenarios, the algorithm generates 100% or the required changes for scenarios with 25% error rate. This can be considered satisfactory for most error-generation applications. The percentage reduces progressively for 50% and 75% errors.

Total Execution Time. Figures 3.i-k report total execution times over an increasing number of tuples and 5% injected errors for datasets Employees, Customers, and Tax, respectively. Not surprisingly, most of the time is consumed by the updates to the DB (there are thousands of updates, which are known to be slow). Error generation does not dominate the total time and is stable over the different scenarios, while the computation of detectability and repairability vary significantly depending on the characteristics of the constraints. Again, Tax is the scenario with the most complex constraints, and this is reflected in the execution time for the detection (note the different scale in the y axis). Real-data in Figure 3.l is also consistent with synthetic datasets of the corresponding size in the distribution of the time for the different tasks. Hospital requires less time because of the size and the smaller set of constraints.

Detectability and Repairability. Figure 3.m shows how the ratio between *at-least-one* and *exactly-one* detectable changes depends on the constraints and the data. For example, Hospital has redundant data and overlapping rules, thus it is less likely that the errors we find are *exactly-one* detectable changes.

To further study this relationship, we run 10 experiments with different db sizes and error percentage for the three synthetic datasets. We report in Figure 3.n the average percentage of exactly-one detectable changes, with 95% confidence intervals. The figure shows that, give an error-generation task, it is fairly easy to estimate the number of at-least-one detectable changes to request in order to obtain a given number of exactly-one detectable changes.

Figures 3.o-q report the average repairability of errors in Customers per constraint, using three configurations (i.e., high, medium, and low repairability). Every constraint has the same required amount of errors (5%), and all configurations used the same clean instance of 200k tuples. Intuitively, a high repairability configuration involves mostly rules with master data and CFDs, while a low repairability one involves FDs. Results for medium-repairability configurations on the remaining scenarios are in Figures 3.r-t.

9. DATA REPAIRING TOOL EVALUATION

In this section we present an empirical comparison of several data-repairing algorithms over BART data. We show a number of novel insights into these algorithms that could not have been shown with existing error generators.

Tools and Algorithms. We used two publicly available constraint-based data-repairing tools, namely Llunatic [21] and Nadeef [13]. Using these, we were able to run four data-repairing algorithms: (i) Greedy [7, 11], as implemented within the Nadeef tool; (ii) Holistic [10], also implemented within Nadeef; (iii) Lunatic, the chase-based algorithm described in [19]; (iv) Sampling [3], implemented using Lunatic as described in [20]. In addition, we obtained a copy of the (v) SCARE [36] statistics-based tool from the authors.

Tasks. We tested four tasks, three of these constraint-based and one statistics-based (i.e., random errors). For testing constraint-based algorithms, we used (i) Hospital, (ii) Bus, and (iii) Tax. We restricted the set of DCs to FDs and CFDs as only these can be handled by the algorithms under consideration. We selected a clean instance of 20K tuples, and made it dirty with 5% errors and different repairability levels: High (approximately 0.8 rep.), Med (0.5 rep), and Low (0.25 rep.). For testing SCARE [36], we used task (iv) Tax-Rnd; we injected errors of different kinds in the Tax dataset: (a) 5% missing values, (b) 5% typos, and (c) 5% outliers over numerical attributes.

We only report results for the constraint-based algorithms over detectable errors (tasks (i)–(iii)), and for SCARE over random errors (task (iv)). In fact, the performance of the algorithms is quite poor when they are applied to errors that are outside of their scope.

Results. The purpose of this evaluation is not to assess the quality of repair algorithms, rather to show how BART can be used to uncover new insights into the data-repairing process. We measured the quality of repairs in terms of precision/recall using the Value, Cell-Var, and Cell metrics as explained in Section 7. When multiple repairs were returned, we chose the best one. We show the results in Figure 3.q–t.

(a) We notice a wide degree of variability in quality among all algorithms, from excellent repairs to low-quality ones (a trend observed in Figures 3.q–s). This variability does not clearly emerge from evaluations reported in the literature, an observation that suggests there is no definitive data-repairing algorithm yet.

(b) We observe different trends with respect to repairability: (b.1) some of the algorithms return very good repairs when sufficient information is available (i.e., high repairability); however, their quality tends to degrade quickly as repairability decreases; (b.2) on the contrary, the **Sampling** algorithm [3] – for which we return the best and average quality among a random sample of 500 repairs – is less affected by different levels of repairability.

(c) For Task (iv) (as shown in Figure 3.t), (c.1) different kinds of random errors pose challenges of different complexity, with missing values being the easiest ones to detect; (c.2) interestingly, these errors are more easily detected than fixed, as shown by the differences between the **Cell** and **Value** metrics (the first one only counts cells that have been correctly identified as erroneous, while the second one requires that they have been restored to their original value).

A key observation that emerges from this study is that repairability has a strong correlation with the quality of the repairs, thus capturing the “hardness” of the problem. Our study also shows the importance of having systematic error-generation tools for evaluating data-repairing solutions.

10. RELATED WORK

Many pieces of work in the literature have dealt with the problem of injecting errors into clean data. Typically, these approaches inject random errors by scanning database attributes involved in data-quality rules and changing attribute values with some fixed probability [3, 7, 10, 13, 18]. This strategy does not guarantee that all errors are detectable. There are cases where the authors state [11] that all injected errors are detectable, but provide little to none details on how this is done. Examples of errors are the following: dummy values (e.g., ‘xxx’) [7, 11], typos (e.g., by inserting a fixed character like ‘!’), and switched values between two cells of the same tuple or column, among others.

Missing Answers and Why-Not Provenance. Given a query Q and a DB I , the missing answer problem computes explanations for why $Q(I)$ does not include one or more tuples. *Query-based* approaches [4, 9, 22, 34] determine which parts of the query cause the answer to be missing (or how to change the query to make the answer appear), and *instance-based* approaches [24, 25] determine how the input instance can be changed to make the missing answer appear in the result. Recent work [23] also consider hybrid explanations. We view the problem of introducing errors into a clean instance as an instance-based missing answer problem, i.e., how to modify the DB to make violation query results non empty. As we strive to perform only updates of attributes values, approaches that insert or delete tuples to explain missing answers (e.g., [24]) are inapplicable to our problem. While Huang et al.’s technique [25] supports updates, it is known to produce some incorrect answers (see [24]). How-To queries as supported by Tiresias [32] could also be used to encode error generation. It is known that most reasonably complete approaches rely on

constraint solvers and moreover on possible world semantics. As such, they need to consider multiple solutions and minimize the presence of non-expected tuples in the queries (i.e., side-effects). In contrast, by applying *at-least-one* detectability we do not have to compute multiple solutions and minimize side-effects. Furthermore, we introduce a novel optimization for symmetric queries which has not been considered by missing answer approaches.

View Update. The missing answer problem and also ours, are essentially a new take on the well-studied view update problem [2]. The problem of introducing errors is equivalent to inserting tuples into views corresponding to the violation-detection queries. As before we only consider approaches that translate view insertions into updates to the base data. While an overview of this problem is beyond the scope of the present work, two examples are the work of Shu [33] and Cong et al. [12]. Shu [33] proposes to encode the problem as constraint satisfaction, similar as some of the missing answer approaches discussed before. Cong et al. [12] study the complexity of the view update problem and its relationship to annotation propagation. One major cost factor in view update is computing a translation that minimizes the side-effects on the view and/or the instance. We avoid that cost through our at-least-one detectability semantics.

Generating Test Databases. Similar to approaches used in missing answers and view update, approaches for test database generation [1, 5, 29, 30] that are applicable to larger classes of queries (e.g., [30]) encode the problem as constraint satisfaction and, thus, have to pay the high cost of running a constraint solver. Even though these approaches try to reduce the size of the constraints passed to the constraint solver, this is not always successful.

Inconsistency Measures. Many measures have been proposed to quantify the inconsistency of data, e.g., Decker has proposed metrics that count the number of constraint violations [14]. Repairability is essentially a new inconsistency measure that measures how hard it would be to restore a DB to its original clean state. An interesting area for future work is to provide control over other measures such as Hunter and Konieczny’s incompatibility ratio (informally, the number of possible repairs over the size of the data) [26].

11. REFERENCES

- [1] A. Arasu, R. Kaushik, and J. Li. Data generation using declarative constraints. In *SIGMOD*, pages 685–696, 2011.
- [2] F. Bancilhon and N. Spyratos. Update Semantics of Relational Views. *ACM TODS*, 6(4):557–575, 1981.
- [3] G. Beskales, I. F. Ilyas, and L. Golab. Sampling the repairs of functional dependency violations under hard constraints. *PVLDB*, 3:197–207, 2010.
- [4] N. Bidoit, M. Herschel, K. Tzompanaki, et al. Query-based why-not provenance with nedexplain. In *EDBT*, 2014.
- [5] C. Binnig, D. Kossmann, and E. Lo. Reverse query processing. In *ICDE*, pages 506–515, 2007.
- [6] P. Bohannon, W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for data cleaning. In *ICDE*, pages 746–755, 2007.
- [7] P. Bohannon, M. Flaster, W. Fan, and R. Rastogi. A cost-based model and effective heuristic for repairing constraints by value modification. In *SIGMOD*, pages 143–154, 2005.
- [8] L. Bravo, W. Fan, and S. Ma. Extending Dependencies with Conditions. In *VLDB*, pages 243–254, 2007.
- [9] A. Chapman and H. V. Jagadish. Why Not? In *SIGMOD*, pages 523–534, 2009.

[10] X. Chu, I. F. Ilyas, and P. Papotti. Holistic data cleaning: Putting violations into context. In *ICDE*, 2013.

[11] G. Cong, W. Fan, F. Geerts, X. Jia, and S. Ma. Improving data quality: Consistency and accuracy. In *VLDB*, pages 315–326, 2007.

[12] G. Cong, W. Fan, F. Geerts, J. Li, and J. Luo. On the Complexity of Annotation Propagation and View Update Analyses. 2012.

[13] M. Dallachiesa, A. Ebaid, A. Eldawy, A. K. Elmagarmid, I. Ilyas, M. Ouzzani, and N. Tang. Nadeef: a commodity data cleaning system. In *SIGMOD*, pages 541–552, 2013.

[14] H. Decker. Checking and Repairing the Quality of Information in Databases by Inconsistency Metrics. In *ICIQ*, pages 139–150, 2012.

[15] W. Fan and F. Geerts. *Foundations of Data Quality Management*. Morgan & Claypool, 2012.

[16] W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional Functional Dependencies for Capturing Data Inconsistencies. *ACM TODS*, 33, 2008.

[17] W. Fan, J. Li, S. Ma, N. Tang, and W. Yu. Towards certain fixes with editing rules and master data. *PVLDB*, 3(1):173–184, 2010.

[18] F. Geerts, G. Mecca, P. Papotti, and D. Santoro. The Llunatic Data-Cleaning Framework. *PVLDB*, 6(9):625–636, 2013.

[19] F. Geerts, G. Mecca, P. Papotti, and D. Santoro. Mapping and Cleaning. In *ICDE*, pages 232–243, 2014.

[20] F. Geerts, G. Mecca, P. Papotti, and D. Santoro. Mapping and Cleaning: The Llunatic Way. Technical Report TR 01-2014 - <http://www.db.unibas.it/projects/1llunatic/>, 2014.

[21] F. Geerts, G. Mecca, P. Papotti, and D. Santoro. That’s all folks! LLUNATIC goes open source. *PVLDB*, 7(13):1565–1568, 2014.

[22] Z. He and E. Lo. Answering why-not questions on top-k queries. In *ICDE*, pages 750–761, 2012.

[23] M. Herschel. Wondering why data are missing from query results?: ask conseil why-not. In *CIKM*, pages 2213–2218. ACM, 2013.

[24] M. Herschel and M. A. Hernández. Explaining missing answers to spjua queries. *PVLDB*, 3(1):185–196, 2010.

[25] J. Huang, T. Chen, A. Doan, and J. F. Naughton. On the provenance of non-answers to queries over extracted data. *PVLDB*, 1(1):736–747, 2008.

[26] A. Hunter and S. Konieczny. Approaches to measuring inconsistent information. In *Inconsistency Tolerance*, volume 3300 of *LNCS*, pages 191–236. 2005.

[27] S. Kolahi and L. V. S. Lakshmanan. On Approximating Optimum Repairs for Functional Dependency Violations. In *ICDT*, 2009.

[28] H. Liu, S. Shah, and W. Jiang. On-line outlier detection and data cleaning. *Computers & chemical engineering*, 28(9):1635–1647, 2004.

[29] E. Lo, C. Binnig, D. Kossmann, M. Tamer Özsu, and W.-K. Hon. A framework for testing dbms features. *VLDB J.*, 19(2):203–230, 2010.

[30] E. Lo, N. Cheng, and W.-K. Hon. Generating databases for query workloads. *PVLDB*, 3(1):848–859, 2010.

[31] A. Lopatenko and L. Bravo. Efficient Approximation Algorithms for Repairing Inconsistent Databases. In *ICDE*, pages 216–225, 2007.

[32] A. Meliou and D. Suciu. Tiresias: the Database Oracle for How-To Queries. In *SIGMOD*, pages 337–348, 2012.

[33] H. Shu. Using constraint satisfaction for view update. *J. of Intelligent Inf. Syst.*, 15(2):147–173, 2000.

[34] Q. T. Tran and C.-Y. Chan. How to conquer why-not questions. In *SIGMOD*, pages 15–26, 2010.

[35] J. Wang and N. Tang. Towards dependable data repairing with fixing rules. In *SIGMOD*, pages 457–468, 2014.

[36] M. Yakout, L. Berti-Équille, and A. K. Elmagarmid. Don’t be SCARED: use SCalable Automatic REpairing with

maximal likelihood and bounded changes. In *SIGMOD*, pages 553–564, 2013.

APPENDIX

A. PSEUDOCODE AND PROOFS

In this appendix we present the pseudocode for Algorithm 1, and provide proofs for our results.

Theorem 1: *The exactly-one and at-least-one error-generation problems for a task \mathbf{E} are NP-complete.*

PROOF. Both problems are clearly in NP if the number of relational atoms per DC is bound by a constant k , because after we non-deterministically pick-up a set of cell changes, we can check in polynomial time whether it is a solution by running vio-detection queries.

Exactly-one. We prove NP-completeness through a reduction from 3-colorability. Assume we are given graph $G = (V, E)$. We initially color all edges with the same color. To check 3-colorability, we construct an *exactly-one* detectable error-generation task \mathbf{E} , and show that G is 3-colorable iff there exists a solution for \mathbf{E} .

To do this, the colored version of the input graph $G = (V, E)$ is encoded as relation $C(x, y, c_1, c_2)$, storing edges (x, y) with end colors c_1 and c_2 . Colors are encoded as r =red, b =blue, g =green. We add $C(x, y, r, r)$ and $C(y, x, r, r)$ to C for every $(x, y) \in E$. Now we use DCs to encode that (i) nodes may only be colored in r, g, b ;

$$\begin{aligned} \text{dc}_{chooseL} &: \neg(C(x, y, c_1, c_2), c_1 \neq r, c_1 \neq b, c_1 \neq g) \\ \text{dc}_{chooseR} &: \neg(C(x, y, c_1, c_2), c_2 \neq r, c_2 \neq b, c_2 \neq g) \end{aligned}$$

(ii) each vertex must have a unique color;

$$\begin{aligned} \text{dc}_{VColorL} &: \neg(C(x, y, c_1, c_2), C(x', y', c'_1, c'_2), x = x', c_1 \neq c'_1) \\ \text{dc}_{VColorR} &: \neg(C(x, y, c_1, c_2), C(x', y', c'_1, c'_2), y = y', c_2 \neq c'_2) \\ \text{dc}_{VColorLR} &: \neg(C(x, y, c_1, c_2), C(x', y', c'_1, c'_2), x = y', c_1 \neq c'_2) \\ \text{dc}_{VColorRL} &: \neg(C(x, y, c_1, c_2), C(x', y', c'_1, c'_2), x' = y, c_2 \neq c'_1) \end{aligned}$$

(iii) adjacent vertices all have the same color

$$\text{dc}_{noSameCE} : \neg(E(x, y, c_1, c_2), c_1 \neq c_2)$$

Note that all constraints trivially hold in the initial instance we have constructed, as all nodes are of the same color (red). To show that the graph G is 3-colorable, we shall introduce violations to constraint $\text{dc}_{noSameCE}$, thus making adjacent vertices of different colors. We require this to happen for all vertices, i.e. the constraint is violated $2 \times |E|$ times.

We create an error-generation task \mathbf{E} with $I = \{C\}$, Σ including all constraints above, and Conf such that:

$$\begin{aligned} \epsilon(\text{dc}_{chooseL}) &= 0, \epsilon(\text{dc}_{chooseR}) = 0 \\ \epsilon(\text{dc}_{noSameCE}) &= 2 \times |E| \\ \epsilon(\text{dc}_{VColorL}) &= 0, \epsilon(\text{dc}_{VColorR}) = 0, \epsilon(\text{dc}_{VColorLR}) = 0 \end{aligned}$$

Applying the cell changes of a solution to the input DB, according to the *exactly-one* variant of task \mathbf{E} , generates a 3-coloring of the input graph. If there exists no solution for \mathbf{E} , then the graph is not 3-colorable. This proves the claim. \square

At-least-one. We prove NP-completeness through a reduction from the Hitting Set problem: Given sets S_1 to S_n with $\mathcal{U} = \bigcup_{i \in \{1, n\}} S_i$, find $S \subseteq \mathcal{U}$ such that $|S| \leq k$ and for every $i \in \{1, n\}$ we have $S_i \cap S \neq \emptyset$. We construct an error-generation task such that finding a solution returns a hitting set S of size less than or equal to k for S_1 to S_n , and if there exists no solution then this implies that no such hitting set exists. We use relations $S_i(x)$, $i = 1, n$ to encode

Algorithm 1 VIO-GEN ALGORITHM($S, \Sigma, I, \text{Conf}$)

```
1:  $Ch = \emptyset$ 
2:  $changeContexts = \text{new Map} < \text{Cell}, \text{Context} >$ 
3: for all  $dc \in \Sigma$  do
4:    $\mathcal{Q}_{dc} = \text{GENERATEVIOGENQUERIES}(dc)$ 
5:   for all  $GQ_{dc,i} \in \mathcal{Q}_{dc}$  do
6:      $skipped = \emptyset$ 
7:      $changes = 0$ 
8:      $n = \text{NUMOFERRORSINCONFIG}(GQ_{dc,i}, \text{Conf})$ 
9:     if  $n = 0$  then continue end if
10:     $p = \text{CHOOSEPROBABILITY}(GQ_{dc,i}, \text{Conf})$ 
11:     $o = \text{CHOOSEOFFSET}(GQ_{dc,i}, \text{Conf})$ 
12:     $queryPlan = \text{STANDARDJOINQUERYPLAN}$ 
13:    if  $\text{HASNOEQUALITIES}(GQ_{dc,i})$  then
14:       $queryPlan = \text{CROSSPRODUCTSAMPLINGQUERYPLAN}$ 
15:    end if
16:    if  $\text{ISYMMETRIC}(GQ_{dc,i})$  then
17:       $queryPlan = \text{SYMMETRICQUERYPLAN}$ 
18:    end if
19:     $contexts = \text{EXECUTEQUERY}(queryPlan)$ 
20:    while  $\text{HASMORECONTEXTS}(contexts, skipped)$  do
21:       $context = \text{GETNEXTCONTEXT}(contexts, skipped)$ 
22:       $skipped.remove(context)$ 
23:       $cells = \text{SELECTVIOGENCELLS}(context)$ 
24:      if  $\text{IMMUTABLE}(cells)$  then continue end if
25:      if  $\text{OVERLAP}(context, cells, changeContexts)$  then
26:        continue
27:      end if
28:      if  $skipped.size() < o$  then
29:         $skipped.add(context)$ 
30:        continue
31:      end if
32:      if  $\text{random}() > p$  then
33:         $skipped.add(context)$ 
34:        continue
35:      end if
36:       $white = \text{COMPUTEWHITELIST}(context, cells)$ 
37:       $black = \text{COMPUTEBLACKLIST}(context, cells)$ 
38:       $value = \text{SELECTNEWVALUE}(white, black)$ 
39:       $changes = \text{GENERATECELLCHANGE}(value, cells)$ 
40:       $Ch.add(changes)$ 
41:       $changeContexts.putAll(cells, context)$ 
42:       $changes = changes + cells.size()$ 
43:      if  $(n \neq ALL)$  and  $(n = changes)$  then break end if
44:    end while
45:  end for
46: end for
47: return  $Ch$ 
```

the input sets. Each element $e \in S_i$ is represented as a tuple $S_i(e)$. These relations are marked as immutable (master data) for the error-generation task. The hitting set is encoded as a relation $H(y, z)$ which initially contains k tuples, $(1, e_{new})$ to (k, e_{new}) , each with a new value e_{new} not in \mathcal{U} . The initial DB instance is therefore such that S_i and H have no common values, as encoded by the following DCs, one for each $i = 1, n$:

$$dc_{hit(i)} : \neg(S(x), H(y, x'), x = x')$$

We create an error-generation task \mathbf{E} with $I = \{S_1, \dots, S_n, H\}$, $\Sigma = \{dc_{hit(i)} | i = 1, n\}$, and Conf such that:
 $\forall i \in \{1, n\} : \epsilon(dc_{hit(i)}) = 1$

A solution to this error generation problem with at-least-one detectability creates a hitting set encoded in H . Note that if a solution contains unmodified e_{new} elements we ignore these elements, and obtain a hitting set H with $|S| < k$.

Theorem 2: *Given an error-generation task $\mathbf{E} = \langle S, \Sigma, I, \text{Conf} \rangle$, call Ch the output of the vio-gen algorithm detailed in Algorithm 1 over \mathbf{E} . If the algorithm succeeds, then Ch is an at-least-one solution for task \mathbf{E} .*

PROOF. A vio-gen query $GQ_{dc,i}$ only differs from the vio-detection query DQ_{dc} for a denial constraint dc in that the i^{th} comparison atom (the target comparison) is negated. Let t and t' be two tuples from the instance that are returned by $GQ_{dc,i}$ the variable (if any) of the LHS, respective RHS, of the target comparison was assigned to t , respective t' . Let $t.A$ and $t'.A$ be the cells assigned to the target comparison. WLOG assume that the algorithm has chosen to apply change $ch : t.A := v$. We have to prove that $t.A$ was not part of any violation context in I and there exists an assignment m so that dc is violated in I_d with m and $t.A \in \text{vio-context}_{dc}(m)$. Since I is clean, we know that the first requirement is fulfilled. The second requirement is fulfilled if the vio-detection DQ_{dc} returns a result containing t . Consider the tuples in t_1 to t_m that were assigned to relational atoms of query $GQ_{dc,i}$ to derive t and t' . This assignment is unique, because id's are included in the query result. Let t'_i denote the updated version of tuple t_i in instance I_d . We know that $t_i = t'_i$ with the only exception of t , because our algorithm discards changes to tuples that are part of the violation context of another cell change. Thus, for t'_1, \dots, t'_m all comparison atoms of DQ_{dc} are guaranteed to evaluate to true except for comparison atoms involving $t.A$ including the target comparison. Our algorithms never updates cells in a way that would invalidate other comparisons than the target comparison using the white- and black-lists. Since we only introduce cell changes that make the target comparison fail, we know that in I_d the target comparison's negation in DQ_{dc} will evaluate to true. Since the target comparison is the only difference between $GQ_{dc,i}$ and DQ_{dc} , it follows that $DQ_{dc}(I_d)$ returns a violation involving $t.A$. \square

Theorem 3 *Let $\mathbf{E} = \langle S, \Sigma, I, \text{Conf} \rangle$ be an error-generation task. The Vio-gen algorithm runs in PTIME in $\|I\|$ and $\|\Sigma\|$ and is NP-hard in the maximal size (number of atoms) of constraints in $\|\Sigma\|$.*

PROOF. The main cost factors in the vio-gen algorithm are the execution of queries and bookkeeping of contexts. We first analyze the cost of query evaluation and then study the cost of client-side operations.

Query Evaluation Cost: The algorithm first creates all vio-gen queries for each constraint in Σ . For a constraint with n comparisons, there are n potential vio-gen queries. Let $maxSize$ denote the maximum number of atoms (comparisons and relational atoms) that occurs in a constraint from Σ . Recall that a vio-gen query is generated by negating one of the comparisons of a vio-detection query. Thus, the number of generated queries is bound by $\|\Sigma\| \times maxSize$. The data complexity of relational algebra queries is PTIME while the combined complexity is PSPACE-complete (and, thus, NP-hard). In worst-case all results for all generated vio-gen queries have to be retrieved. Thus, the complexity of the vio-gen algorithm in terms of query evaluation is PTIME in $\|\Sigma\|$, $\|I\|$, and NP-hard in $maxSize$.

Client-side operations: Several tasks in the algorithm are executed on the client-side. We now analyze the complexity of each of these tasks. For each $dc \in \Sigma$ we have to generate vio-gen queries (line 4). This is achieved by inverting one comparison in the vio-detection query DQ_{dc} for dc . We know that there are less than $maxSize$ comparisons. Thus, both the size and number of vio-gen queries for dc is bound by $maxSize$ resulting in $\mathcal{O}(maxSize^2 \times \|\Sigma\|)$ runtime. Checking whether a vio-gen query $GQ_{dc,i}$ has no equalities (line 13) can be done in time linear in $maxSize$ and symmetry can be tested in PTIME in $maxSize$ (line 16). The remaining tasks are in PTIME in $maxSize$: selecting vio-gen cells

from a context (line 23); determining whether a cell is immutable (line 24); computing overlap of a context vio-gen cells with contexts for previous changes (line 25); generation of white- and black-lists (lines 36 and 37); selecting a new value from these lists (line 38); and generating the actual cell change (line 39). In combination, the operations applied to each returned context are PTIME in $maxSize$.

To summarize, the runtime complexity of the approach is dominated by query execution which is PTIME in $\|I\|$ and $\|\Sigma\|$ and NP-hard in $maxSize$. \square

B. DETAILS ON THE EXPERIMENTS

We present here full details regarding the datasets and the constraints defined over them, as used in our empirical evaluation (Sections 8 and 9).

Employees is the running example used through the paper. There are two relations (one is master data) with a total of seven attributes and five constraints.

$$\begin{aligned} dc_1 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d \neq d') \\ dc_2 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d = d', \\ &\quad d = \text{“Sales”}, m \neq m') \\ dc_3 &: \neg(\text{Emp}(n, d, s, m), d = \text{“Staff”}, s \neq \text{“1000”}) \\ dc_4 &: \neg(\text{Emp}(n, d, s, m), \text{MD}(n', d', m'), n = n', d = d', m \neq m') \\ dc_5 &: \neg(\text{Emp}(n, d, s, m), \text{Emp}(m', d', s', m''), m = m', s' < s) \end{aligned}$$

Customers is a synthetic scenario from Geerts et al. [20]. There are three relations (one is master data) with a total of sixteen attributes and eight constraints.

$$\begin{aligned} dc_1 &: \neg(\text{Cust}(\text{SSN}, \text{Name}, \text{Phon}), \text{Cust}(\text{SSN}', \text{Name}', \text{Phon}'), \\ &\quad \text{SSN} = \text{SSN}', \text{Phon} = \text{Name}', \text{Phon} \neq \text{Phon}') \\ dc_2 &: \neg(\text{Cust}(\text{SSN}, \text{Name}, \text{CC}), \text{Cust}(\text{SSN}', \text{Name}', \text{CC}'), \\ &\quad \text{SSN} = \text{SSN}', \text{Phon} = \text{Name}', \text{CC} \neq \text{CC}') \\ dc_3 &: \neg(\text{Treat}(\text{Insu}, \text{Trea}), \text{Insu} = \text{“Abx”}, \text{Trea} \neq \text{“Dental”}) \\ dc_4 &: \neg(\text{Cust}(\text{SSN}, \text{Name}, \text{Phon}), \text{MD}(\text{SSN}', \text{Name}', \text{Phon}'), \\ &\quad \text{SSN} = \text{SSN}', \text{Name} \neq \text{Name}', \text{Phon} = \text{Phon}') \\ dc_5 &: \neg(\text{Cust}(\text{SSN}, \text{Phon}, \text{Str}), \text{MD}(\text{SSN}', \text{Phon}', \\ &\quad \text{Str}'), \text{SSN} = \text{SSN}', \text{Str} \neq \text{Str}', \text{Phon} = \text{Phon}') \\ dc_6 &: \neg(\text{Cust}(\text{SSN}, \text{Phon}, \text{City}), \text{MD}(\text{SSN}', \text{Phon}', \\ &\quad \text{City}'), \text{SSN} = \text{SSN}', \text{City} \neq \text{City}', \text{Phon} = \text{Phon}') \\ dc_7 &: \neg(\text{Cust}(\text{SSN}, \text{City}), \text{Treat}(\text{SSN}', \text{Insu}'), \\ &\quad \text{SSN} = \text{SSN}', \text{Insu} = \text{“Abx”}, \text{City} = \text{“SF”}), \\ dc_8 &: \neg(\text{Treat}(\text{SSN}, \text{Sal}), \text{Treat}(\text{SSN}', \text{Sal}') \text{SSN} = \text{SSN}', \text{Sal} \neq \text{Sal}') \end{aligned}$$

Tax is a synthetic scenario from Fan et al. [16]. Each record represents an individual’s address and tax information. The relation has thirteen attributes and five constraints.

$$\begin{aligned} dc_1 &: \neg(\text{Tax}(\text{ac}, \text{st}), \text{Tax}(\text{ac}', \text{st}'), \text{ac} = \text{ac}', \text{st} \neq \text{st}') \\ dc_2 &: \neg(\text{Tax}(\text{st}, \text{zip}), \text{Tax}(\text{st}', \text{zip}'), \text{st} \neq \text{st}', \text{zip} = \text{zip}') \\ dc_3 &: \neg(\text{Tax}(\text{st}, \text{ch}, \text{ctx}), \text{Tax}(\text{st}', \text{ch}', \text{ctx}'), \text{st} = \text{st}', \text{ch} = \text{ch}', \text{ctx} \neq \text{ctx}') \\ dc_4 &: \neg(\text{Tax}(\text{st}, \text{ms}, \text{stx}), \text{Tax}(\text{st}', \text{ms}', \text{stx}'), \text{st} = \text{st}', \text{ms} = \text{ms}', \text{stx} \neq \text{stx}') \\ dc_5 &: \neg(\text{Tax}(\text{st}, \text{sal}, \text{tr}), \text{Tax}(\text{st}', \text{sal}', \text{tr}'), \text{st} = \text{st}', \text{sal} > \text{sal}', \text{tr} < \text{tr}') \end{aligned}$$

Bus is a real-world scenario from Dallachiesa et al. [13]. It consists of one relation obtained by joining eight tables from the UK government public datasets³. The table has 250K

tuples with twenty-five attributes and twelve constraints.

$$\begin{aligned} dc_1 &: \neg(\text{Bus}(\text{loco}, \text{lona}), \text{Bus}(\text{loco}', \text{lona}'), \text{loco} = \text{loco}', \text{lona} \neq \text{lona}') \\ dc_2 &: \neg(\text{Bus}(\text{loco}, \text{lnla}), \text{Bus}(\text{loco}', \text{lnla}'), \text{loco} = \text{loco}', \text{lnla} \neq \text{lnla}') \\ dc_3 &: \neg(\text{Bus}(\text{loco}, \text{aaco}), \text{Bus}(\text{loco}', \text{aaco}'), \text{loco} = \text{loco}', \text{aaco} \neq \text{aaco}') \\ dc_4 &: \neg(\text{Bus}(\text{loco}, \text{ndco}), \text{Bus}(\text{loco}', \text{ndco}'), \text{loco} = \text{loco}', \text{ndco} \neq \text{ndco}') \\ dc_5 &: \neg(\text{Bus}(\text{reco}, \text{rena}), \text{Bus}(\text{reco}', \text{rena}'), \text{reco} = \text{reco}', \text{rena} \neq \text{rena}') \\ dc_6 &: \neg(\text{Bus}(\text{pzco}, \text{coun}), \text{Bus}(\text{pzco}', \text{coun}'), \text{pzco} = \text{pzco}', \text{coun} \neq \text{coun}') \\ dc_7 &: \neg(\text{Bus}(\text{dico}, \text{dina}), \text{Bus}(\text{dico}', \text{dina}'), \text{dico} = \text{dico}', \text{dina} \neq \text{dina}') \\ dc_8 &: \neg(\text{Bus}(\text{aaco}, \text{dico}), \text{Bus}(\text{aaco}', \text{dico}'), \text{dico} = \text{dico}', \text{aaco} \neq \text{aaco}') \\ dc_9 &: \neg(\text{Bus}(\text{aaco}, \text{arna}), \text{Bus}(\text{aaco}', \text{arna}'), \text{aaco} = \text{aaco}', \text{arna} \neq \text{arna}') \\ dc_{10} &: \neg(\text{Bus}(\text{aaco}, \text{reco}), \text{Bus}(\text{aaco}', \text{reco}'), \text{aaco} = \text{aaco}', \text{reco} \neq \text{reco}') \\ dc_{11} &: \neg(\text{Bus}(\text{cdti}, \text{mdti}, \text{reno}), \text{cdti} = \text{mdti}, \text{reno} \neq \text{“0”}) \\ dc_{12} &: \neg(\text{Bus}(\text{cdti}, \text{mdti}, \text{modi}), \text{cdti} = \text{mdti}, \text{modi} \neq \text{“new”}) \end{aligned}$$

Hospital is a real-world scenario used in several data-repairing papers (e.g., [13, 17, 20]). It is based on a dataset from the US Department of Health & Human Services.⁴ The database contains a single table with 100k tuples and nineteen attributes, over which seven constraints are defined.

$$\begin{aligned} dc_1 &: \neg(\text{Hosp}(\text{City}, \text{Zip}), \text{Hosp}(\text{City}', \text{Zip}'), \text{City} \neq \text{City}', \text{Zip} = \text{Zip}') \\ dc_2 &: \neg(\text{Hosp}(\text{Stat}, \text{Zip}), \text{Hosp}(\text{Stat}', \text{Zip}'), \text{Stat} \neq \text{Stat}', \text{Zip} = \text{Zip}') \\ dc_3 &: \neg(\text{Hosp}(\text{Zip}, \text{Ph}), \text{Hosp}(\text{Zip}', \text{Ph}'), \text{Zip} \neq \text{Zip}', \text{Ph} = \text{Ph}') \\ dc_4 &: \neg(\text{Hosp}(\text{City}, \text{Ph}), \text{Hosp}(\text{City}', \text{Ph}'), \text{City} \neq \text{City}', \text{Ph} = \text{Ph}') \\ dc_5 &: \neg(\text{Hosp}(\text{Stat}, \text{Ph}), \text{Hosp}(\text{Stat}', \text{Ph}'), \text{Stat} \neq \text{Stat}', \text{Ph} = \text{Ph}') \\ dc_6 &: \neg(\text{Hosp}(\text{PrNo}, \text{MC}, \text{StAv}), \text{Hosp}(\text{PrNo}', \text{MC}', \\ &\quad \text{StAv}'), \text{PrNo} = \text{PrNo}', \text{MC} = \text{MC}', \text{StAv} \neq \text{StAv}') \\ dc_7 &: \neg(\text{Hosp}(\text{Stat}, \text{MC}, \text{StAv}), \text{Hosp}(\text{Stat}', \text{MC}', \\ &\quad \text{StAv}'), \text{Stat} = \text{Stat}', \text{MC} = \text{MC}', \text{StAv} \neq \text{StAv}') \end{aligned}$$

For the synthetic scenarios, we generated up to one million tuples for the relations. When available, the master-data table contains 20% of the tuples in the other tables. Further details are provided in Figure 4.

C. USAGE DETAILS

We further describe how to use our error-generation tool. Given an error-generation task \mathbf{E} , we control the error-generation process by specifying a set of configuration parameters Conf . Due to space limitations, we only discuss the parameters that have the strongest impact over BART.

Types of Relations. We use a configuration parameter called “Authoritative Source” to specify which DB relations are to be considered *immutable*. By default, we assume we can inject errors over any relation.

Types of Errors. We support two kinds of errors: (i) *constraint-induced errors*, in which we use the input constraints to decide how the noise is added to the data; and (ii) *random errors*, in which we inject randomly different error types. For the latter, we describe below different ways in which a value in a DB may be erroneous and discuss how we can simulate these random errors using BART.⁵ (i) for constraint-induced errors, we may specify the desired error percentage for each constraint dc , denoted by $\epsilon(dc)$. Given a constraint $dc \in \Sigma$, we call R.A an attribute *covered by* dc if there exists a variable x_i that has more than one occurrence in dc , and one of the occurrences of x_i appears within a relational atom for table R associated with attribute A. The error percentage for dc and R.A is denoted by $\epsilon(dc, R.A)$, and is also specified with respect to the number of tuples in table R. In addition, we may also specify if errors should be *at-least-one* or *exactly-one* detectable;

⁴<http://www.medicare.gov/hospitalcompare/>

⁵To guarantee detectability (as presented in Section 5), we take black-lists and vio-gen contexts into account.

³<http://data.gov.uk/data>

(ii) for random errors, we may specify the percentage of various kinds of random errors, denoted by $\epsilon(rand)$, with respect to the total number of tuples of the database.

Example 3: Consider a constraint over Emp:

$$dc_1 : \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d \neq d')$$

Covered attributes: Emp.Name, Emp.Dept

By requiring that $\epsilon(dc_1, \text{Emp.Dept}) = 1\%$, $\epsilon(dc_1, \text{Emp.Name}) = 1\%$, we ask BART to make dirty both attributes with the same ratio. We may decide how many errors will be introduced for the equality comparison with respect to those for the inequality by changing the two values. \square

Error Percentages. We use a configuration parameter called “Error Percentage” to specify the desired degree of noise in task **E**. We specify percentages by reasoning with respect to the number of tuples in a table, rather than on the number of cells, which depends also on the number of attributes. For example, it is easy to see that 1% errors in a table of 100K tuples requires to make dirty 1000 cells. BART supports the following options:

(i) for constraint-induced errors, we may specify the desired error percentage for each constraint dc , denoted by $\epsilon(dc)$. Given a constraint $dc \in \Sigma$, we call R.A an attribute *covered by dc* if there exists a variable x_i that has more than one occurrence in dc , and one of the occurrences of x_i appears within a relational atom for table R associated with attribute A. The error percentage for dc and R.A is denoted by $\epsilon(dc, R.A)$, and is also specified with respect to the number of tuples in table R. In addition, we may also specify if errors should be *at-least-one* or *exactly-one* detectable;

(ii) for random errors, we may specify the percentage of various kinds of random errors, denoted by $\epsilon(rand)$, with respect to the total number of tuples of the database.

Example 4: Consider a constraint over Emp:

$$dc_1 : \neg(\text{Emp}(n, d, s, m), \text{Emp}(n', d', s', m'), n = n', d \neq d')$$

Covered attributes: Emp.Name, Emp.Dept

By requiring that $\epsilon(dc_1, \text{Emp.Dept}) = 1\%$, $\epsilon(dc_1, \text{Emp.Name}) = 1\%$, we ask BART to make dirty both attributes with the same ratio. We may decide how many errors will be introduced for the equality comparison with respect to those for the inequality by changing the two values. \square