

Extending Conditional Dependencies with Built-in Predicates

Journal:	<i>Transactions on Knowledge and Data Engineering</i>
Manuscript ID:	TKDE-2014-12-0983
Manuscript Type:	Regular
Keywords:	H.2.1.d Database models < E.0.c Data encryption < E.0 General < E Data, H.2.1.c Database integration < E.0.c Data encryption < E.0 General < E Data, E.0.b Data dependencies < E.0 General < E Data

SCHOLARONE™
Manuscripts

Review Only

Extending Conditional Dependencies with Built-in Predicates

Shuai Ma, Liang Duan, Wenfei Fan, Chunming Hu, and Wenguang Chen

Abstract—This paper proposes a natural extension of conditional functional dependencies (CFDs [23]) and conditional inclusion dependencies (CINDs [31]), denoted by CFD^p s and $CIND^p$ s, respectively, by specifying patterns of data values with $\neq, <, \leq, >$ and \geq predicates. As data quality rules, CFD^p s and $CIND^p$ s are able to capture errors that commonly arise in practice but cannot be detected by CFDs and CINDs. We establish two sets of results for central technical problems associated with CFD^p s and $CIND^p$ s. (a) One concerns the satisfiability and implication problems for CFD^p s and $CIND^p$ s, taken separately or together. These are important for, e.g., deciding whether data quality rules are dirty themselves, and for removing redundant rules. We show that despite the increased expressive power, the static analyses of CFD^p s and $CIND^p$ s retain the same complexity as their CFDs and CINDs counterparts. (b) The other concerns validation of CFD^p s and $CIND^p$ s. We show that given a set Σ of CFD^p s and $CIND^p$ s on a database D , a set of SQL queries can be automatically generated that, when evaluated against D , return all tuples in D that violate some dependencies in Σ . We also experimentally verified the efficiency and effectiveness of our SQL based error detection techniques, using real-life data. This provides commercial DBMS with an immediate capability to detect errors based on CFD^p s and $CIND^p$ s.

1 INTRODUCTION

Extensions of traditional functional dependencies (FDs) and inclusion dependencies (INDs), known as *conditional functional dependencies* (CFDs [23]) and *conditional inclusion dependencies* (CINDs [31]), respectively, have recently been proposed for improving data quality. These extensions enforce patterns of semantically related data values, and detect errors as violations of the dependencies. It has been shown that conditional dependencies are able to capture more inconsistencies than FDs and INDs [18], [22], [31].

Conditional dependencies specify constant patterns in terms of equality ($=$). In practice, however, the semantics of data often need to be specified with other predicates such as $\neq, <, \leq, >$ and \geq , as illustrated by the following example.

Example 1: An online store maintains a database of two relations: (a) item for items sold by the store, and (b) tax for the sale tax rates for the items, except artwork, in various states. The relations are specified by the following schemas:

```
item (id: string, name: string, type: string, price: float,
      shipping: float, sale: bool, state: string)
tax (state: string, rate: float)
```

where each item is specified by its id, name, type (e.g., book, CD), price, shipping fee, the state to which it is shipped, and whether it is on sale. A tax tuple specifies the sale tax rate in a state. An instance D_0 of item and tax is shown in Fig. 1.

- S. Ma, L. Duan and C. Hu are with the SKLSDE lab, School of Computer Science and Engineering, Beihang University, China.
E-mail: {mashuai, duanl, hucm}@act.buaa.edu.cn.
- W. Fan is with the RCBD center, Beihang University, China and the School of Informatics, Edinburgh University, UK.
E-mail: wenfei@inf.ed.ac.uk.
- W. Chen is with the Department of Information Management, Peking University, China.
E-mail: chenwg@pku.edu.cn.

Manuscript received XXX, 2014; revised XXX, 2014.

One wants to specify dependencies on the relations as data quality rules to detect errors in the data, such that inconsistencies emerge as violations of the dependencies. Traditional dependencies (FDs, INDs; see, e.g., [3]) and conditional dependencies (CFDs, CINDs [23], [31]) on the data include the following:

```
cf1: item (id  $\rightarrow$  name, type, price, shipping, sale)
cf2: tax (state  $\rightarrow$  rate)
cf3: item (sale = 'T'  $\rightarrow$  shipping = 0)
```

These are CFDs: (a) cf_1 assures that the id of an item uniquely determines its name, type, price, shipping and sale; (b) cf_2 states that state is a key for tax, i.e., for each state there is a unique sale tax rate; and (c) cf_3 ensures that for any item tuple t , if $t[\text{sale}] = \text{'T'}$ then $t[\text{shipping}]$ must be 0; i.e., free shipping is provided for items on sale. Here cf_3 is specified in terms of patterns of semantically related data values, namely, sale = 'T' and shipping = 0. It is to hold only on item tuples that match the pattern sale = 'T'. In contrast, cf_1 and cf_2 are traditional FDs without constant patterns, a special case of CFDs. One can verify that no sensible INDs or CINDs can be defined across item and tax.

Note that D_0 of Fig. 1 satisfies cf_1 , cf_2 and cf_3 . That is, when these dependencies are used as data quality rules, no errors are found in D_0 .

In practice, the shipment fee of an item is typically determined by the price of the item. Moreover, when an item is on sale, the price of the item is often in a certain range. Furthermore, for any item sold by the store to a customer in a state, if the item is *not* artwork, then one expects to find the sale tax rate in the state from the tax table. These semantic relations cannot be expressed as CFDs of [23] or CINDs of [31], but can be expressed as the following dependencies:

```
pfd1: item (sale = 'F' and price  $\leq$  20  $\rightarrow$  shipping = 3)
pfd2: item (sale = 'F' and price > 20 and price  $\leq$  40
            $\rightarrow$  shipping = 6)
pfd3: item (sale = 'F' and price > 40  $\rightarrow$  shipping = 10)
```

	id	name	type	price	shipping	sale	state		state	rate
t_1 :	b1	Harry Potter	book	25.99	0	T	WA	t_5 :	PA	6
t_2 :	c1	Snow White	CD	9.99	2	F	NY	t_6 :	NY	4
t_3 :	b2	Catch-22	book	34.99	20	F	DL	t_7 :	DL	0
t_4 :	a1	Sunflowers	art	5m	500	F	DL	t_8 :	NJ	3.5

(a) An item relation

(b) A tax rate relation

Figure 1. Example instance D_0 of item and tax

pfd_4 : item (sale = 'T' \rightarrow price \geq 2.99 and price $<$ 9.99)
 pind_1 : item (state; type \neq 'art') \subseteq tax (state; nil)

Here pfd_2 states that for any item tuple, if it is not on sale and its price is in the range (20, 40], then its shipment fee must be 6; similarly for pfd_1 and pfd_3 . These dependencies extend CFDs [23] by specifying patterns of semantically related data values in terms of predicates $<$, \leq , $>$ and \geq . Similarly, pfd_4 assures that for any item tuple, if it is on sale, then its price must be in the range [2.99, 9.99). Finally, pind_1 extends CINDs [31] by specifying patterns with \neq : for any item tuple t , if $t[\text{type}]$ is not artwork, then there must exist a tax tuple t' such that $t[\text{state}] = t'[\text{state}]$, i.e., the sale tax of the item can be found from the tax relation.

Using dependencies pfd_1 – pfd_4 and pind_1 as data quality rules, we find that D_0 of Fig. 1 is not clean. Indeed, (a) t_2 violates pfd_1 : its price is less than 20, but its shipping fee is 2 rather than 3; similarly, t_3 violates pfd_2 , and t_4 violates pfd_3 . (b) Tuple t_1 violates pfd_4 : it is on sale but its price is not in the range [2.99, 9.99). (c) The database D_0 also violates pind_1 : t_1 is not artwork, but its state cannot find a match in the tax relation, i.e., no tax rate for WA is found in D_0 . \square

None of pfd_1 – pfd_4 and pind_1 can be expressed as FDs or INDs [3], which do not allow constants, or as CFDs [23] or CINDs [31], which specify patterns with equality ($=$) only. While there have been extensions of CFDs [10], [13], [30], none of these allows dependencies to be specified with patterns on data values in terms of built-in predicates \neq , $<$, \leq , $>$ or \geq . To the best of our knowledge, the earlier conference version [12] of this paper is the first to study these constraints.

These highlight the need for extending CFDs and CINDs to capture errors commonly found in real-life data. While one can consider arbitrary extensions, it is necessary to strike a balance between their expressive power and their complexity. In particular, we want to be able to reason about data quality rules expressed as extended CFDs and CINDs. Furthermore, we want to have effective algorithms to detect inconsistencies based on these extensions.

Contributions & Roadmap. To this end we introduce an extension of CFDs and CINDs, investigate the static analyses of these constraints, and develop effective SQL-based techniques for detecting errors based on these constraints.

(1) We propose two classes of dependencies, denoted by CFD^p s and CIND^p s, which respectively extend CFDs and CINDs by supporting \neq , $<$, \leq , $>$, \geq predicates (Sections 2 and 3). For example, all the dependencies we have encountered so far can be expressed as CFD^p s or CIND^p s. These dependencies are capable of capturing errors in real-world data that cannot be detected by CFDs or CINDs.

(2) We establish the complexity bounds for the satisfiability and implication problems for CFD^p s and CIND^p s, taken sepa-

rately or together (Section 4). The satisfiability problem is to determine whether a set Σ of dependencies has a nonempty model, i.e., whether the rules in Σ are consistent themselves. The implication problem is to decide whether a set Σ of dependencies entails another dependency φ , i.e., whether the rule φ is redundant in the presence of the rules in Σ . These are the central technical problems associated with any dependency language.

We show that despite the increased expressive power, CFD^p s and CIND^p s do not increase the complexity for reasoning about them. In particular, we show that the satisfiability and implication problems remain (a) NP-complete and CONP-complete for CFD^p s, respectively, (b) in $O(1)$ -time (constant-time) and EXPTIME-complete for CIND^p s, respectively, and (c) are undecidable when CFD^p s and CIND^p s are taken together. These are the same as their CFDs and CINDs counterparts [31]. In contrast, data with linearly ordered domains often makes our lives harder [35].

(3) We provide SQL-based techniques to detect errors based on CFD^p s and CIND^p s (Section 5). Given a set Σ of CFD^p s and CIND^p s on a database D , we automatically generate a set of SQL queries that, when evaluated on D , find all tuples in D that violate some dependencies in Σ . Further, the SQL queries are independent of the size and cardinality of Σ . These provide the capability of detecting errors in a single relation (CFD^p s) and across different relations (CIND^p s) within the immediate reach of commercial DBMS.

(4) Using real-life data (HOSP and DBLP), we finally conduct an extensive experimental study (Section 6). We show that (a) the running time of CFD^p s and CIND^p s is comparable to their CFDs and CINDs counterparts, which is consistent with the static analyses in Section 4, and (b) CFD^p s and CIND^p s are able to capture more errors than their CFDs and CINDs counterparts (22% on HOSP and 75% on DBLP), due to the increased expressive power.

Related work. This paper is an extension of our earlier work [12] by adding (a) the proofs for the complexity bounds for the satisfiability and implication analyses of CFD^p s and CIND^p s, separately and taken together (Section 4), and (b) an extensive experimental study of CFD^p s and CIND^p s (Section 6), which was not investigated in [12].

Recently, data dependencies have generated renewed interests for improving data quality [5], [10], [14], [15], [22], [23], [30], [31], [34], [36]. Constraint-based data cleaning was introduced in [4], which proposed to use dependencies, e.g., FDs, INDs and denial constraints, to detect and repair errors in real-life data (see, e.g., [3], [15], [22], [34] for details). Data dependencies have been studied for relational databases since the introduction of FDs by Codd [17] in 1972 (see, e.g., [3] for details), and the theory of INDs was established in [11], which developed a sound and complete

$(1) \varphi_1 = \text{tax}(\text{state} \rightarrow \text{rate}, T_1)$ $T_1: \begin{array}{c c} \text{state} & \text{rate} \\ \hline - & - \end{array}$	$(2) \varphi_2 = \text{item}(\text{sale} \rightarrow \text{shipping}, T_2)$ $T_2: \begin{array}{c c} \text{sale} & \text{shipping} \\ \hline = T & = 0 \end{array}$
$(3) \varphi_3 = \text{item}(\text{sale}, \text{price} \rightarrow \text{shipping}, T_3)$ $T_3: \begin{array}{c c c} \text{sale} & \text{price} & \text{shipping} \\ \hline = F & > 20 & = 6 \\ = F & \leq 40 & = 6 \end{array}$	$(4) \text{CFD}^p \varphi_4 = \text{item}(\text{sale} \rightarrow \text{price}, T_4)$ $T_4: \begin{array}{c c} \text{sale} & \text{price} \\ \hline = T & \geq 2.99 \\ = T & < 9.99 \end{array}$

Figure 2. Example CFD^p s

inference system and the PSPACE-completeness for the implication analysis of INDs. As an extension of traditional FDs, CFDs were developed in [23], for improving the quality of data. It was shown in [23] that the satisfiability and implication problems for CFDs are NP-complete and coNP-complete, respectively. Along the same lines, CINDs [31] were proposed to extend INDs, and it was shown [31] that the satisfiability and implication problems for CINDs are in constant time and EXPTIME-complete, respectively. SQL techniques were developed in [23] to detect errors by using CFDs, but have not been studied for CINDs. This work extends the static analyses of conditional dependencies of [23], [31], and has established several new complexity results, notably in the absence of finite-domain attributes (e.g., Theorems 2, 6, 8). In addition, it is the first work to develop SQL techniques for checking violations of CINDs and violations of CFD^p s and CIND^p s taken together.

Extensions of CFDs have been proposed to support disjunction and negation [10], cardinality constraints and synonym rules [13], and to specify patterns in terms of value ranges [30]. While CFD^p s are more powerful than the extension of [30], they cannot express disjunctions [10], cardinality constraints and synonym rules [13]. To our knowledge no extensions of CINDs have been studied. This work is the first full treatment of extensions of CFDs and CINDs by incorporating built-in predicates ($\neq, <, \leq, >, \geq$), from static analyses to error detection.

Methods have been developed for discovering CFDs [14], [30], CFD^p s [36] and CINDs [5] and for repairing data based on either CFDs [18], traditional FDs and INDs taken together [8], CFDs and CINDs taken together [20], denial constraints [7], [16], aggregate constraints [28], matching dependencies [21], matching dependencies and CFDs [27], or editing rules and master data [25]. We defer the treatment of these topics for CFD^p s and CIND^p s to future work.

A variety of extensions of FDs and INDs have been studied for specifying constraint databases and constraint logic programs [6], [9], [32], [33]. While the languages of [6], [32] cannot express CFDs, constraint-generating dependencies (CGDs) of [6] and constrained tuple-generating dependencies (CTGDs) of [33] can express CFD^p s, and CTGDs can also express CIND^p s. The increased expressive power of CTGDs comes at the price of a higher complexity: both their satisfiability and implication problems are undecidable. Built-in predicates and arbitrary constraints are supported by CGDs, for which it is not clear whether effective SQL queries can be developed to detect errors. It is worth mentioning that Theorems 2 and 6 of this work provide lower bounds for the consistency and implication analyses of CGDs, by using patterns with built-in predicates only.

2 EXTENDING CFDs WITH PREDICATES

We now define *conditional functional dependencies with predicates*, denoted by CFD^p s, by extending CFDs [23] with built-in predicates ($\neq, <, \leq, >, \geq$) in addition to equality ($=$).

Consider a relational schema R defined over a finite set of attributes, denoted by $\text{attr}(R)$. For each attribute $A \in \text{attr}(R)$, its domain is specified in R , denoted as $\text{dom}(A)$, which is either finite (e.g., bool) or infinite (e.g., string). We assume w.l.o.g. that a domain on which $<, \leq, >$ or \geq is defined is totally ordered.

Syntax. A $\text{CFD}^p \varphi$ on R is a pair $R(X \rightarrow Y, T_p)$, where (1) X, Y are sets of attributes in $\text{attr}(R)$; (2) $X \rightarrow Y$ is a standard FD, referred to as the FD *embedded in* φ ; and (3) T_p is a tableau with attributes in X and Y , referred to as the *pattern tableau* of φ , where for each A in $X \cup Y$ and each tuple $t_p \in T_p$, $t_p[A]$ is either an unnamed variable ‘ $_$ ’ that draws values from $\text{dom}(A)$, or ‘ $\text{op } a$ ’, where op is one of $\{=, \neq, <, \leq, >, \geq\}$, and ‘ a ’ is a constant in $\text{dom}(A)$.

If attribute A occurs in both X and Y , we use A_L and A_R to indicate the occurrence of A in X and Y , respectively, and we separate the X and Y attributes in a pattern tuple with ‘ $\|$ ’. We simply write φ as $(X \rightarrow Y, T_p)$ when R is clear from the context, and denote X as $\text{LHS}(\varphi)$ and Y as $\text{RHS}(\varphi)$, respectively.

Example 2: The dependencies cfd_1 – cfd_3 and pfd_1 – pfd_4 that we have seen in Example 1 can all be expressed as CFD^p s. Some of these CFD^p s are illustrated in Fig. 2, in which φ_1 is for FD cfd_2 , φ_2 is for CFD cfd_3 , φ_3 is for pfd_2 , and φ_4 is for pfd_4 , respectively. \square

Semantics. Consider $\text{CFD}^p \varphi = R(X \rightarrow Y, T_p)$, where $T_p = \{t_{p_1}, \dots, t_{p_k}\}$.

A data tuple t of R is said to *match* $\text{LHS}(\varphi)$, denoted by $t[X] \asymp T_p[X]$, if for each tuple t_{p_i} ($i \in [1, k]$) in T_p and each attribute A in X , either (a) $t_{p_i}[A]$ is the wildcard ‘ $_$ ’ (which matches any value in $\text{dom}(A)$), or (b) $t[A] \text{ op } a$ if $t_{p_i}[A]$ is ‘ $\text{op } a$ ’, where the operator op ($=, \neq, <, \leq, >$ or \geq) is interpreted by its standard semantics. Similarly, the notion that t *matches* $\text{RHS}(\varphi)$ is defined, denoted by $t[Y] \asymp T_p[Y]$.

Intuitively, each pattern tuple t_{p_i} ($i \in [1, k]$) specifies a condition via $t_{p_i}[X]$, and $t[X] \asymp T_p[X]$ if $t[X]$ satisfies the *conjunction* of all these conditions. Similarly, $t[Y] \asymp T_p[Y]$ if $t[Y]$ matches all the patterns specified by $t_{p_i}[Y]$ for all pattern tuples t_{p_i} in T_p .

An instance I of R *satisfies* the $\text{CFD}^p \varphi$, denoted by $I \models \varphi$, if for each pair of tuples t_1, t_2 in I , if $t_1[X] = t_2[X] \asymp T_p[X]$, then $t_1[Y] = t_2[Y] \asymp T_p[Y]$. That is, if $t_1[X]$ and $t_2[X]$ are equal and in addition, they both match the pattern tableau

$T_p[X]$, then $t_1[Y]$ and $t_2[Y]$ must also be equal to each other and they both match the pattern tableau $T_p[Y]$.

Observe that φ is imposed only on the subset of tuples in I that match LHS(φ), rather than on the entire I . For all tuples t_1, t_2 in this subset, if $t_1[X] = t_2[X]$, then (a) $t_1[Y] = t_2[Y]$, i.e., the semantics of the embedded FDs is enforced; and (b) $t_1[Y] \prec T_p[Y]$, which assures that the *constants* in $t_1[Y]$ match the *constants* in $t_{p_i}[Y]$ for all t_{p_i} in T_p . Note that here tuples t_1 and t_2 can be the same.

An instance I of R satisfies a set Σ of CFD^ps, denoted by $I \models \Sigma$, if $I \models \varphi$ for each CFD^p φ in Σ .

Example 3: The instance D_0 of Fig. 1 satisfies φ_1 and φ_2 of Fig. 2, but neither φ_3 nor φ_4 . Indeed, tuple t_3 violates (i.e., does not satisfy) φ_3 , since $t_3[\text{sale}] = 'F'$ and $20 < t_3[\text{price}] \leq 40$, but $t_3[\text{shipping}]$ is 20 instead of 6. Note that t_3 matches LHS(φ_3) since it satisfies the condition specified by the *conjunction* of the pattern tuples in T_3 . Similarly, t_1 violates φ_4 , since $t_1[\text{sale}] = 'T'$ but $t_1[\text{price}] > 9.99$. Observe that while it takes two tuples to violate a standard FD, a single tuple may violate a CFD^p. \square

Special cases. (1) A standard FD $X \rightarrow Y$ [3] can be expressed as a CFD ($X \rightarrow Y, T_p$) in which T_p contains a single tuple consisting of $'_'$ only, without constants. (2) A CFD ($X \rightarrow Y, T_p$) [23] with $T_p = \{t_{p1}, \dots, t_{pk}\}$ can be expressed as a set $\{\varphi_1, \dots, \varphi_k\}$ of CFD^ps such that for each $i \in [1, k]$, $\varphi_i = (X \rightarrow Y, T_{p_i})$, where T_{p_i} contains the pattern tuple t_{p_i} of T_p only, defined with equality ($=$) only. For example, φ_1 and φ_2 in Fig. 2 are CFD^ps representing FD cfd₂ and CFD cfd₃ in Example 1, respectively. Note that all data quality rules in [14], [30] can be expressed as CFD^ps.

3 EXTENDING CINDS WITH PREDICATES

Similar to CFD^ps, we define *conditional inclusion dependencies with predicates*, denoted by CIND^ps, by extending CINDs [31] with built-in predicates ($\neq, <, \leq, >, \geq$) in addition to equality ($=$). Consider two relational schemas R_1 and R_2 .

Syntax. A CIND^p ψ is a pair $(R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$, where (1) X, X_p and Y, Y_p are lists of attributes in $\text{attr}(R_1)$ and $\text{attr}(R_2)$, respectively; (2) $R_1[X] \subseteq R_2[Y]$ is a standard IND, referred to as the IND *embedded* in ψ ; and (3) T_p is a tableau, called the *pattern tableau* of ψ defined over attributes $X_p \cup Y_p$, and for each A in X_p or Y_p and each pattern tuple $t_p \in T_p$, $t_p[A]$ is either an unnamed variable $'_'$ that draws values from $\text{dom}(A)$, or $'\text{op } a'$, where op is one of $=, \neq, <, \leq, >, \geq$ and $'a'$ is a constant in $\text{dom}(A)$.

We denote $X \cup X_p$ as LHS(ψ), $Y \cup Y_p$ as RHS(ψ), and separate the X_p and Y_p attributes in a pattern tuple with $'||'$. We also use nil to denote an *empty* list.

Example 4: Figure 3 shows two example CIND^ps: ψ_1 expresses the pind₁ in Example 1, and ψ_2 refines ψ_1 by stating that for any item tuple t_1 , if its type is not art and its state is DL, then there must be a tax tuple t_2 such that its state is DL and rate is 0, i.e., ψ_2 assures that the sale tax rate in Delaware is 0. \square

Semantics. Consider CIND^p $\psi = (R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$. An instance (I_1, I_2) of (R_1, R_2) satisfies the CIND^p ψ , denoted by $(I_1, I_2) \models \psi$, iff for each tuple $t_1 \in I_1$, if $t_1[X_p] \prec$

$T_p[X_p]$, then there *exists* a tuple $t_2 \in I_2$ such that $t_1[X] = t_2[Y]$ and $t_2[Y_p] \prec T_p[Y_p]$.

That is, if $t_1[X_p]$ matches the pattern tableau $T_p[X_p]$, then ψ assures the existence of t_2 such that (1) $t_1[X] = t_2[Y]$ as needed by the standard IND embedded in ψ ; and, moreover, (2) $t_2[Y_p]$ must match the pattern tableau $T_p[Y_p]$. In other words, ψ is “conditional” since its embedded IND is applied only to the subset of tuples in I_1 that match $T_p[X_p]$, and $T_p[Y_p]$ is enforced on the tuples in I_2 that match those tuples in I_1 . As remarked in Section 2, the pattern tableau T_p specifies the *conjunction* of all the pattern tuples in T_p .

Example 5: The instance D_0 of item and tax in Fig. 1 violates CIND^p ψ_1 . Indeed, tuple t_1 in item *matches* LHS(ψ_1) since $t_1[\text{type}] \neq 'art'$, but there is no tuple t in tax such that $t[\text{state}] = t_1[\text{state}] = 'WA'$. In contrast, D_0 satisfies ψ_2 . \square

We say that a database D satisfies a set Σ of CINDs, denoted by $D \models \Sigma$, if $D \models \varphi$ for each $\varphi \in \Sigma$.

Safe CIND^ps. We say a CIND^p $(R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$ is *unsafe* if there exist pattern tuples t_p, t'_p in T_p such that either (a) there exists $B \in Y_p$, such that $t_p[B]$ and $t'_p[B]$ are not satisfiable when taken together, or (b) there exist $C \in Y, A \in X$ such that A corresponds to C in the embedded IND and $t_p[C]$ and $t'_p[A]$ are not satisfiable when taken together; e.g., $t_p[\text{price}] = 9.99$ and $t'_p[\text{price}] \geq 19.99$.

Obviously unsafe CIND^ps do not make sense: no nonempty databases satisfy unsafe CIND^ps. It takes $O(|T_p|^2)$ time in the size $|T_p|$ of T_p to decide whether a CIND^p is unsafe. Thus in the sequel we consider safe CIND^p only.

Special cases. (1) A standard CIND $(R_1[X] \subseteq R_2[Y])$ can be expressed as a CIND^p $(R_1[X; \text{nil}] \subseteq R_2[Y; \text{nil}], T_p)$ such that T_p is simply a *empty* set. (2) A CIND $(R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$ with $T_p = \{t_{p1}, \dots, t_{pk}\}$ can be expressed as a set $\{\psi_1, \dots, \psi_k\}$ of CIND^ps, where for each $i \in [1, k]$, $\psi_i = (R_1[X; X_p] \subseteq R_2[Y; Y_p], T_{p_i})$ such that T_{p_i} consists of the pattern tuple t_{p_i} of T_p , defined with equality ($=$) only.

4 REASONING ABOUT CFD^ps AND CIND^ps

The satisfiability and implication problems are the two classical questions associated with any dependency languages [3], [23], [31]. In this section we investigate these problems for CFD^ps and CIND^ps, separately and taken together.

4.1 Satisfiability Analyses

The satisfiability problem is to determine, given a set Σ of constraints, whether there exists a *nonempty* database that satisfies Σ .

The satisfiability analysis of conditional dependencies is not only of theoretical interest, but is also important in practice. Indeed, when CFD^ps and CIND^ps are used as data quality rules, this analysis helps one check whether the rules make sense themselves. The need for this is particularly evident when the rules are manually designed or discovered from various datasets [5], [14], [30].

Satisfiability analysis of CFD^ps. Given any FDs, one does not need to worry about their satisfiability as any set of FDs is always satisfiable. However, as observed in [23], for a set Σ of CFDs on a relational schema R , there may not exist

$$\begin{aligned}
(1) \psi_1 &= (\text{item} [\text{state}; \text{type}] \subseteq \text{tax} [\text{state}; \text{nil}], T_1), \\
(2) \psi_2 &= (\text{item} [\text{state}; \text{type}, \text{state}] \subseteq \text{tax} [\text{state}; \text{rate}], T_2) \\
T_1: & \begin{array}{c|c} \text{type} & \text{nil} \\ \hline \neq \text{art} & \end{array} & T_2: & \begin{array}{c|c|c} \text{type} & \text{state} & \text{rate} \\ \hline \neq \text{art} & = \text{DL} & = 0 \end{array}
\end{aligned}$$

Figure 3. Example CIND^ps

a *nonempty* instance I of R such that $I \models \Sigma$. As CFDs are a special case of CFD^ps, the same problem exists when it comes to CFD^ps.

Example 6: Consider a CFD^p $\varphi = (R : A \rightarrow B, T_p)$ such that $T_p = \{(_ \parallel a), (_ \parallel \neq a)\}$. There is no *nonempty* instance I of R that satisfies φ . Indeed, for any R tuple t , φ requires that both $t[B] = a$ and $t[B] \neq a$, which is impossible. \square

This problem is already NP-complete for CFDs [23]. Below we show that it remains the same complexity for CFD^ps despite their increased expressive power.

Proposition 1: *The satisfiability problem for CFD^ps is NP-complete.* \square

Proof: The lower bound follows from the NP-hardness of their CFDs counterparts [23], since CFDs are a special case of CFD^ps. The upper bound is verified by presenting an NP algorithm that, given a set Σ of CFD^ps defined on a relational schema R , determines whether Σ is satisfiable.

We next present an NP algorithm that, given a set Σ of CFD^ps defined on a relational schema R , determines whether Σ is satisfiable or not. The satisfiability problem has the following small model property: If there is a nonempty R instance I such that $I \models \Sigma$, then for any tuple $t \in I$, instance $I_t = \{t\}$ satisfies Σ . Thus it suffices to consider single-tuple instances $I = \{t\}$ for deciding whether Σ is satisfiable.

Assume *w.l.o.g.* that the attributes $\text{attr}(R) = \{A_1, \dots, A_m\}$ and the total number of pattern tuples in all pattern tableaux T_p of CFD^ps in Σ is h . For each $i \in [1, m]$, define the active domain of A_i to be a set $\text{adom}(A_i) = C_0 \cup C_1$, where (1) C_0 consists of all constants in $T_p[A_i]$ of all pattern tableaux T_p in Σ , and if C_0 is empty, we further let $C_0 = \{a_1, a_2\}$, where $a_1, a_2 \in \text{dom}(A_i)$ and $a_1 \neq a_2$, and (2) C_1 contains the set of constants for the attributes whose domains have total orders, *i.e.*, involved with predicates $\neq, <, \leq, >$ or \geq :

- (1) Arrange all constants in C_0 in the increasing order, and assume the resulting $C_0 = \{a_1, \dots, a_k\}$ ($k \geq 1$);
- (2) Add a constant $b_{01} \in \text{dom}(A_i)$ to C_1 such that $b_{01} < a_1$ if there exists one; And also add another constant $b_{02} \in \text{dom}(A_i)$ to C_1 such that $b_{02} < a_1$ and $b_{02} \neq b_{01}$ if there exists one;
- (3) Similarly, for each $j \in [1, k-1]$, add a constant $b_{j1} \in \text{dom}(A_i)$ to C_1 such that $a_j < b_{j1} < a_{j+1}$ if there exists one; And also add another constant $b_{j2} \in \text{dom}(A_i)$ to C_1 such that $a_j < b_{j2} < a_{j+1}$ and $b_{j2} \neq b_{j1}$ if there exists one;
- (4) Finally, add a constant $b_{k1} \in \text{dom}(A_i)$ to C_1 such that $b_{k1} > a_k$ if there exists one; And also add another constant $b_{k2} \in \text{dom}(A_i)$ to C_1 such that $b_{k2} > a_k$ and $b_{k2} \neq b_{k1}$ if there exists one.

Moreover, the number of elements in $\text{adom}(A_i)$ is at most $3 * h + 2$. Then one can easily verify that Σ is satisfiable iff

there exists a mapping ρ from $t[A_i]$ to $\text{adom}(A_i)$ ($i \in [1, m]$) such that $I = \{(\rho(t[A_1]), \dots, \rho(t[A_m]))\}$ and $I \models \Sigma$.

We now give an NP algorithm as follows: (1) Guess an instance, which contains a single tuple t of R such that $t[A_i] \in \text{adom}(A_i)$ for each $i \in [1, m]$. (2) Check whether $I \models \Sigma$. If so the algorithm returns ‘yes’, and otherwise it repeats steps (1) and (2). Obviously step (2) can be done in PTIME in the size of Σ . Hence the algorithm is in NP, and so is the problem. \square

It is known [23] that the satisfiability problem for CFDs is in PTIME when the CFDs considered are defined over attributes that have an infinite domain, *i.e.*, in the absence of finite domain attributes. However, this is no longer the case for CFD^ps. This tells us that the increased expressive power of CFD^ps does take a toll in this special case. It should be remarked that while the proof of Proposition 1 is an extension of its counterpart in [23], the result below is new.

Theorem 2: *In the absence of finite domain attributes, the satisfiability problem for CFD^ps remains NP-complete.* \square

Proof: The problem is in NP by Proposition 1. Its NP-hardness is shown by reduction from the 3SAT problem, which is NP-complete (cf. [29]).

We next show the reduction from the 3SAT problem. Consider an instance $\phi = C_1 \wedge \dots \wedge C_n$ of 3SAT, where all the variables in ϕ are x_1, \dots, x_m , C_j is of the form $y_{j1} \vee y_{j2} \vee y_{j3}$ such that for each $i \in [1, 3]$, y_{ji} is either $x_{p_{ji}}$ or $\overline{x_{p_{ji}}}$ for $p_{ji} \in [1, m]$. Given ϕ , we construct a relational schema R and a set Σ of CFD^ps defined on R such that ϕ is satisfiable iff Σ is satisfiable.

(1) We first define the relational schema $R(X_1, \dots, X_m, C_1, \dots, C_n, Z)$, where all attributes share a common infinite domain dom that contains constant a . Intuitively, for each R tuple t , $t[X_1, \dots, X_m]$ specifies a truth assignment ξ for variables x_1, \dots, x_m of ϕ , and $t[C_i]$ ($i \in [1, n]$) and $t[Z]$ are the truth values of clause C_i and sentence ϕ w.r.t. the assignment ξ , respectively.

(2) We then construct the set CFD^ps $\Sigma = \Sigma_0 \cup \Sigma_1 \cup \dots \cup \Sigma_n \cup \Sigma_{n+1}$, defined as follows.

- (a) Σ_0 contains $n+1$ CFD^ps, which intuitively encode the relationships of the truth values between the clauses C_1, \dots, C_n and sentence ϕ .

For each clause C_i ($i \in [1, n]$), we add to Σ_0 a CFD^p $\varphi_i = (C_1, \dots, C_n \rightarrow Z, T_{pi})$, in which $T_{pi} = \{t_{pi}\}$ such that $t_{pi}[C_i, Z] = (\neq a \parallel \neq a)$ and $t_{pi}[C_j] = _$ for any $j \neq i$ and $j \in [1, n]$. We also add to Σ_0 a CFD^p $\varphi_0 = (C_1, \dots, C_n \rightarrow Z, T_{p0})$, where $T_{p0} = \{(\neq a, \dots, \neq a \parallel = a)\}$. Intuitively, we use $\neq a$ and $= a$ to represent *false* and *true*, respectively.

- (b) For $i \in [1, n]$, Σ_i contains 8 CFD^ps, which intuitively encode the relationships of the truth values between the clause C_i and its three variables.

For clause $C_i = x_{j_1} \vee \overline{x_{j_2}} \vee x_{j_3}$ of ϕ with $1 \leq j_1, j_2, j_3 \leq m$, we define CFD^ps (a) $\varphi_{i,0} = (X_{j_1}, X_{j_2}, X_{j_3} \rightarrow C_i, T_{pi,0})$ with $T_{pi,0} = \{(\langle a, \langle a, \langle a, \langle a \rangle \rangle \rangle = a)\}$, (b) $\varphi_{i,1} = (X_{j_1}, X_{j_2}, X_{j_3} \rightarrow C_i, T_{pi,1})$ with $T_{pi,1} = \{(\langle a, \langle a, \langle a, \geq a \rangle \rangle \rangle = a)\}$, and (c) $\varphi_{i,2} = (X_{j_1}, X_{j_2}, X_{j_3} \rightarrow C_i, T_{pi,2})$ with $T_{pi,2} = \{(\langle a, \langle a, \geq a, \langle a \rangle \rangle \rangle = a)\}$. Similarly, we can define the rest 5 CFD^ps $\varphi_{i,3}, \varphi_{i,4}, \varphi_{i,5}, \varphi_{i,6}$ and $\varphi_{i,7}$. Intuitively, we further use $\langle a$ and $\geq a$ to represent *false* and *true* for a variable, respectively.

- (c) Σ_{n+1} contains a single CFD^p $\varphi_{n+1} = (Z \rightarrow Z, T_{p(n+1)})$ with $T_{p(n+1)} = \{(_ \parallel = a)\}$. Intuitively, φ_{n+1} requires that for any R tuple t , $t[Z] = a$.

Observe that Σ contains $8 * m + n + 2$ CFD^ps in total. Thus the reduction is in PTIME.

We now show that ϕ is satisfiable iff Σ is satisfiable.

Assume first that Σ is satisfiable, then we show that there exists a nonempty instance I of R such that $I \models \Sigma$. For any tuple $t \in I$, (a) Σ_{n+1} forces $t[Z] = a$, (b) Σ_0 forces $t[C_1, \dots, C_n] = (a, \dots, a)$, and (c) for each clause C_i ($i \in [1, n]$) with variables $X_{j_1}, X_{j_2}, X_{j_3}, X_{j_4}$ forces that $t[X_{j_1}, X_{j_2}, X_{j_3}]$ does not match the LHS of the CFD^ps that force $t[C_i] \neq a$. From the tuple t , we can construct a truth assignment ξ of ϕ such that $\xi(x_i) = \text{false}$ if $t[X_i] < a$ and $\xi(x_i) = \text{true}$ if $t[X_i] \geq a$ ($i \in [1, m]$). Since $\{t\} \models \Sigma$, it is easy to verify that the truth assignment ξ makes ϕ true.

Conversely, if ϕ is satisfiable, there exists a truth assignment ξ that makes ϕ true. We construct a tuple t of R as follows: (a) $t[C_1, \dots, C_n, Z] = (a, \dots, a)$ and (b) for each $i \in [1, m]$, $t[X_i] = a_i$ such that (a) $a_i \in \text{dom}$ and $a_i \geq a$ if $\xi(x_i) = \text{true}$ and (b) $a_i < a$ otherwise. Let $I = \{t\}$, then one can easily verify that $I \models \Sigma$.

Putting these together, we have the conclusion. \square

Satisfiability analysis of CIND^ps. Like FDs, one can specify arbitrary INDs or CINDs without worrying about their satisfiability. Below we show that CIND^ps preserve this nice property, by extending the proof of its counterpart in [31].

Proposition 3: Any set of CIND^ps is always satisfiable. \square

Proof: Given a set Σ of CIND^ps over a database schema $\mathcal{R}(R_1, \dots, R_n)$, we show that one can always construct a nonempty instance D of \mathcal{R} such that $D \models \Sigma$.

We build D as follows. First, for each attribute A , define the active domain of A to be a set $\text{adom}(A)$, which consists of certain data values in $\text{dom}(A)$. Second, using these active domains, we construct D .

(1) We start with the construction of active domains. (a) For each attribute A , initialize $\text{adom}(A)$ along the same lines as the one for CFD^ps in Proposition 1; (b) For each CIND^p $(R_a[A_1, A_2, \dots, A_m; X_p] \subseteq R_b[B_1, \dots, B_m; Y_p], T_p)$ in Σ , let $\text{adom}(B_i) = \text{adom}(B_i) \cup \text{adom}(A_i)$ for each $i \in [1, m]$, and this rule is repeatedly applied until a fixpoint of $\text{adom}(A)$ is reached for all attributes A in \mathcal{R} .

It is easy to verify that this process always terminates as we start with a finite set of data values.

(2) We next construct the database instance D . For each relation $R_i(A_1, \dots, A_k) \in \mathcal{R}$, we define $I_i = \text{adom}(A_1) \times \dots \times \text{adom}(A_k)$, where \times is the Cartesian Product operation [3].

Let $D = \{I_1, \dots, I_n\}$, then it is easy to verify that D is nonempty and $D \models \Sigma$. \square

Satisfiability analysis of CFD^ps and CIND^ps. The satisfiability problem for CFDs and CINDs taken together is undecidable [31]. Since CFD^ps and CIND^ps subsume CFDs and CINDs, respectively, we immediately have the following.

Corollary 4: The satisfiability problem for CFD^ps and CIND^ps is undecidable. \square

4.2 Implication Analyses

The implication problem is to determine, given a set Σ of dependencies and another dependency ϕ , whether or not Σ entails ϕ , denoted by $\Sigma \models \phi$. That is, whether or not for all databases D , if $D \models \Sigma$ then $D \models \phi$.

The implication analysis helps us remove redundant rules, and thus improve the performance of error detection and repairing based on the rules [23], [31].

Example 7: The CFD^ps in Fig. 2 imply another CFD^p $\varphi = \text{item}(\text{sale}, \text{price} \rightarrow \text{shipping}, T)$, where T consists of a single pattern tuple ($\text{sale} = 'F', \text{price} = 30 \parallel \text{shipping} = 6$). Thus in the presence of the CFD^ps in Fig. 2, φ is redundant. \square

Implication analysis of CFD^ps. We first show that the implication problem for CFD^ps retains the same complexity as their CFDs counterpart, verified by extending the proof of its counterpart in [23].

Proposition 5: The implication problem for CFD^ps is coNP-complete. \square

Proof: The lower bound follows from the coNP-hardness of their CFDs counterpart [23], since CFDs are a special case of CFD^ps. The coNP upper bound is verified by presenting an NP algorithm for its complement problem for determining whether $\Sigma \not\models \varphi$.

We next present the a NP algorithm for its complement problem. The algorithm is based on a small model property: if $\varphi = R(X \rightarrow Y, T_p)$ and $\Sigma \not\models \varphi$, then there exists an instance I of R with two tuples t_1 and t_2 such that $I \models \Sigma$ and $t_1[X] = t_2[X] \not\prec T_p[X]$, but either $t_1[Y] \neq t_2[Y]$ or $t_1[Y] \not\prec T_p[Y]$ (resp. $t_2[Y] \not\prec T_p[Y]$). Thus it suffices to consider instances I with two tuples only for deciding whether $\Sigma \not\models \varphi$.

Assume that the attributes $\text{attr}(R) = \{A_1, \dots, A_m\}$. For each $i \in [1, m]$, let $\text{adom}(A_i)$ be the active domain defined in Proposition 1. Then one can easily verify that $\Sigma \not\models \varphi$ iff there exist two mappings ρ_1 and ρ_2 from all attributes A_i to $\text{adom}(A_i)$ ($i \in [1, m]$) such that $I = \{(\rho_1(A_1), \dots, \rho_1(A_m)), (\rho_2(A_1), \dots, \rho_2(A_m))\}$, $I \models \Sigma$, but $I \not\models \varphi$.

Based on these, we give an NP algorithm as follows: (1) Guess two R tuples t_1 and t_2 such that $t_1[A_i], t_2[A_i] \in \text{adom}(A_i)$ for each $i \in [1, m]$. (2) Check whether $I = \{t_1, t_2\}$ satisfies Σ , but not φ . If so the algorithm returns 'yes', and otherwise it repeats steps (1) and (2). Obviously step (2) can be done in PTIME in the size of Σ . Hence the algorithm is in NP, and so is the problem. \square

Similar to the satisfiability analysis, it is known [23] that the implication analysis of CFDs is in PTIME when the

CFDs are defined only with attributes that have an infinite domain. Analogous to Theorem 2, the result below shows that this is no longer the case for CFD^ps, which does not find a counterpart in [23].

Proposition 6: *In the absence of finite domain attributes, the implication problem for CFD^ps is coNP-complete.* \square

Proof: The problem is in coNP by Proposition 5. The coNP-hardness is shown by reduction from the 3SAT problem to its complement problem, i.e., the problem for determining whether $\Sigma \not\models \varphi$.

We next show the reduction from the 3SAT problem to the complement problem of the implication problem for CFD^ps, where 3SAT is NP-complete (cf. Proposition 2). Given an instance ϕ of 3SAT, we construct a relational schema R and a set $\Sigma \cup \{\varphi\}$ of CFD^ps defined on R such that ϕ is satisfiable iff $\Sigma \not\models \varphi$.

The relational schema R and the set Σ of CFD^ps are the same as the corresponding ones in Proposition 2. Moreover, φ is defined as $(Z \rightarrow Z, T_p)$, where $T_p = \{(_ \parallel \neq a)\}$. Intuitively, φ requires that for any R tuple t , $t[Z] \neq a$. Along the same lines as Proposition 2, one can easily verify that ϕ is satisfiable iff $\Sigma \not\models \varphi$. Thus the problem is coNP-hard. \square

Implication analysis of CIND^ps. We next show that CIND^ps do not make their implication analysis harder, verified by extending the proof of their CINDs counterpart given in [31].

Proposition 7: *The implication problem for CIND^ps is EXPTIME-complete.* \square

Proof: The implication problem for CINDs is EXPTIME-hard [31]. Since CIND^ps subsume CINDs, the lower bound carries over to CIND^ps immediately. The EXPTIME upper bound is shown by presenting an EXPTIME algorithm that, given a set $\Sigma \cup \{\psi\}$ of CIND^ps over a database schema \mathcal{R} , determines whether $\Sigma \models \psi$ or not.

We next present the EXPTIME algorithm. Consider $\mathcal{R} = (R_1, \dots, R_n)$ and $\psi = (R_a[X; X_p] \subseteq R_b[Y; Y_p], T_p)$. And for each attribute A , define the active domain $\text{adom}(A)$ of A based on $\Sigma \cup \{\psi\}$ along the same line as the proof of Proposition 3. One can easily verify that if $\Sigma \not\models \psi$, there exists a non-empty instance D of \mathcal{R} such that (a) $D \models \Sigma$ and $D \not\models \psi$, and (b) D consists of data values from the active domains only.

The detailed EXPTIME algorithm is given as follows.

(1) We first build a labeled directed graph $G(V, E, l)$. Each node $u \in V$ is a possible tuple $\langle R_i : t_i \rangle$ such that $t_i[A] \in \text{adom}(A)$ for each attribute $A \in \text{attr}(R_i)$. There is an edge $e = (\langle R_i : t_i \rangle, \langle R_j : t_j \rangle)$ in E iff there exists a CIND^p $\phi = (R_i[U; U_p] \subseteq R_j[V; V_p], T_{p_\phi})$ in Σ such that $t_i[U_p] \succ T_{p_\phi}[U_p]$, $t_j[V] = t_i[U]$ and $t_j[V_p] \succ T_{p_\phi}[V_p]$, and e is labeled with the CIND^p ϕ , i.e., $\phi \in l(e)$. Note that an edge may have multiple labels.

(2) Let S_a be the set of nodes $\langle R_a : t_a \rangle$ such that $t_a[X_p] \succ T_p[X_p]$, and S_b be the set of nodes $\langle R_b : t_b \rangle$ such that $t_b[Y_p] \succ T_p[Y_p]$, respectively.

(3) For each node $u = \langle R_a : t_a \rangle$ in S_a , let G_u be the induced subgraph of G that contains all the nodes reachable from u ,

and exactly the edges that appear in G over the same set of nodes. We also refer to u as the root of G_u .

(4) For an induced subgraph G_u of G with root $u = \langle R_a : t_a \rangle$, we derive another graph G'_u by recursively removing edges as follows. For any v in G_u , if v has a child v' from which no nodes in $\langle R_b : t_b \rangle$ in S_b with $t_b[X] = t_b[Y]$ are reachable, then for all children v'' of v , we remove from labels $l(v, v'')$ all the labels in $l(v, v')$, and edge (v, v'') is removed when $l(v, v'')$ becomes empty.

(5) If there exists a subgraph G'_u derived from an induced subgraph G_u of G with root $u = \langle R_a : t_a \rangle$ such that no nodes $\langle R_b : t_b \rangle$ in S_b with $t_a[X] = t_b[Y]$ are reachable from u , we return 'no', and return 'yes', otherwise.

It can be verified that (a) if the algorithm returns 'no', we can construct an instance D such that $D \models \Sigma$, but not ψ , by collecting those tuples attached on the end nodes of edges whose labels become empty at step 4; and (b) if the algorithm returns 'yes', there exist no instances D such that $D \models \Sigma$, but not ψ .

We next show that the above algorithm indeed runs in exponential time: (a) The number of nodes in graph G is bounded by the maximum number of tuples in a database instance on \mathcal{R} . Let $|\Sigma \cup \{\psi\}|$ be the size of Σ and ψ , and $|\mathcal{R}|$ be the sum of arities of all relations in \mathcal{R} . Then the number of tuples in a database instance is bounded by $O(|\Sigma \cup \{\psi\}|^{|\mathcal{R}|})$; (b) The number of nodes in sets S_a or S_b is bounded by the maximum number of tuples in a database too; (c) The induced subgraph and the reachability testing can be done in linear-time in the size of the input [19].

Putting all these together, we have shown that the algorithm runs in exponential time. And, hence, the problem is in EXPTIME. \square

It is known [31] that the implication problem is PSPACE-complete for CINDs defined with infinite domain attributes. Similar to Theorem 6, below we show that this no longer holds for CIND^ps.

Theorem 8: *In the absence of finite domain attributes, the implication problem for CIND^ps remains EXPTIME-complete.* \square

Proof: The problem is in EXPTIME by Proposition 7. The EXPTIME-hardness is shown by reduction from the implication problem for CINDs in the general setting, in which finite-domain attributes may be present, that is known to be EXPTIME-complete [31].

We next present the reduction from the implication problem for CINDs in the general setting. Given a set $\Sigma \cup \{\psi\}$ of CINDs defined on a database schema $\mathcal{R} (R_1, \dots, R_n)$, we construct another database schema $\mathcal{R}' (R'_1, \dots, R'_n)$, in which each relation R'_i ($i \in [1, n]$) consists of infinite domain attributes only, and a set $\Sigma' \cup \{\psi'\}$ of CIND^ps on \mathcal{R}' such that $\Sigma \models \psi$ iff $\Sigma' \models \psi'$.

(1) We start with constructing \mathcal{R}' . For each $R_i (A_1, \dots, A_k)$ of \mathcal{R} , we define $R'_i (A'_1, \dots, A'_k)$ such that for each attribute A'_j ($j \in [1, k]$), let $\text{dom}(A'_j) = \text{dom}(A_j)$ if $\text{dom}(A_j)$ is infinite, and let $\text{dom}(A'_j)$ be integer, a totally ordered infinite domain, if $\text{dom}(A_j)$ is finite. Moreover, we define a mapping $\rho_{i,j}$ for each finite domain $\text{dom}(A_j) = \{a_1, \dots, a_h\}$ to integer: (a) Randomly choose h consecutive integers $\{b_1, \dots, b_h\}$ such that for each $i \in [1, h-1]$, $b_{i+1} = b_i + 1$.

Table 1
Summary of Complexity Results

Σ	General setting		Infinite domain only	
	Satisfiability	Implication	Satisfiability	Implication
CFDs [23]	NP-complete	coNP-complete	PTIME	PTIME
CFD ^p s	NP-complete	coNP-complete	NP-complete	coNP-complete
CINDs [31]	$O(1)$	EXPTIME-complete	$O(1)$	PSPACE-complete
CIND ^p s	$O(1)$	EXPTIME-complete	$O(1)$	EXPTIME-complete
CFDs + CINDs [31]	undecidable	undecidable	undecidable	undecidable
CFD ^p s + CIND ^p s	undecidable	undecidable	undecidable	undecidable

(b) We now define the mapping $\rho_{i,j}(a_i) = b_i$ for $i \in [1, h]$. Moreover, we require two extra integers $b_0 = b_1 - 1$ and $b_{h+1} = b_h + 1$, denoted as $\rho_{i,j}.b_0$ and $\rho_{i,j}.b_{h+1}$. Note that this is always doable. For clarity, we also denote $\rho_{i,j}$ as ρ when it is clear from the context.

(2) We next define Σ' and ψ' on \mathcal{R}' based on the mappings defined above. For each CIND $(R_a[X; A_1, \dots, A_{m1}] \subseteq R_b[Y; B_1, \dots, B_{m2}], T_p)$ in Σ and each $t_p \in T_p$, we define another CIND^p $(R'_a[X'; A'_1, \dots, A'_{m1}, X'_p] \subseteq R'_b[Y'; B'_1, \dots, B'_{m2}, Y'_p], T'_p)$, where (a) X' (resp. Y' , A'_1, \dots, A'_{m1} and B'_1, \dots, B'_{m2}) corresponds to X (resp. Y , A_1, \dots, A_{m1} and B_1, \dots, B_{m2}); (b) X'_p (resp. Y'_p) corresponds to those finite domain attributes in R_a (resp. R_b), but not in A_1, \dots, A_{m1} (resp. B_1, \dots, B_{m2}); and (c) $T'_p = \{t'_{p1}, t'_{p2}, t'_{p3}\}$ such that for each attribute A' in A'_1, \dots, A'_{m1} or B'_1, \dots, B'_{m2} , (i) $t'_{p1}[A'] = t_p[A]$ and $t'_{p2}[A'] = t'_{p3}[A'] = \text{'-'}'$ if $\text{dom}(A)$ is infinite, and (ii) $t'_{p1}[A'] = \rho(t_p[A])$ and $t'_{p2}[A'] = t'_{p3}[A'] = \text{'-'}'$ if $\text{dom}(A)$ is finite; and (iii) for the rest attributes A' in X'_p or Y'_p , $t'_{p1}[A'] = \text{'-'}'$, $t'_{p2}[A'] = \text{'> } \rho.b_0'$, and $t'_{p3}[A'] = \text{'< } \rho.b_{h+1}'$. Similarly, we can define a set Σ_ψ of CIND^ps from ψ based on the mappings.

Finally, one can easily verify that $\Sigma \models \psi$ iff $\Sigma' \models \Sigma_\psi$, i.e., $\Sigma' \models \psi'$ for each CIND^p $\psi' \in \Sigma_\psi$. Following from this, the problem is EXPTIME-hard. \square

Implication analysis of CFD^ps and CIND^ps. When CFD^ps and CIND^ps are taken together, their implication analysis is beyond reach in practice. This is not surprising since the implication problem for FDs and INDs is already undecidable [3]. Since CFD^ps and CIND^ps subsume FDs and INDs, respectively, from the undecidability result for FDs and INDs, the corollary below follows immediately.

Corollary 9: *The implication problem for CFD^ps and CIND^ps is undecidable.* \square

Summary. The complexity bounds for reasoning about CFD^ps and CIND^ps are summarized in Table 1. To give a complete picture we also include in Table 1 the complexity bounds for the static analyses of CFDs and CINDs, taken from [23], [31]. The results tell us the following.

(1) Despite the increased expressive power, CFD^ps and CIND^ps do not complicate the static analyses in the general case: the satisfiability and implication problems for CFD^ps and CIND^ps have the same complexity bounds as their counterparts for CFDs and CINDs, taken separately or together.

(2) In the special case when CFD^ps and CIND^ps are defined with infinite domain attributes only, however, their static analyses do not get simpler, as opposed to their counterparts

for CFDs and CINDs. That is, the increased expressive power of CFD^ps and CIND^ps comes at a price in this special case.

5 VALIDATION OF CFD^ps AND CIND^ps

If CFD^ps and CIND^ps are to be used as data quality rules, the first question we have to settle is how to effectively detect errors and inconsistencies as violations of these dependencies, by leveraging functionality supported by commercial DBMS. More specifically, consider a database schema $\mathcal{R} = (R_1, \dots, R_n)$, where R_i is a relational schema for $i \in [1, n]$. The error detection problem is stated as follows.

The *error detection problem* is to find, given a set Σ of CFD^ps and CIND^ps defined on \mathcal{R} , and a database instance $D = (I_1, \dots, I_n)$ of \mathcal{R} as input, the subset (I'_1, \dots, I'_n) of D such that for each $i \in [1, n]$, $I'_i \subseteq I_i$ and each tuple in I'_i violates at least one CFD^p or CIND^p in Σ . We denote the set as $\text{vio}(D, \Sigma)$, referred to it as *the violation set of D w.r.t. Σ* .

In this section we develop SQL-based techniques for error detection based on CFD^ps and CIND^ps. The main result of the section is as follows.

Theorem 10: *Given a set Σ of CFD^ps and CIND^ps defined on $\mathcal{R} = (R_1, \dots, R_n)$ and a database instance D of \mathcal{R} , a set of SQL queries can be automatically generated such that (a) the collection of the answers to the SQL queries in D is $\text{vio}(D, \Sigma)$, and (b) the number and size of the set of SQL queries depend only on the number n of relations and their arities in \mathcal{R} , regardless of Σ .* \square

We next present the main techniques for the query generation method. Let Σ_{cfdp}^i be the set of all CFD^ps in Σ defined on the same relational schema R_i , and $\Sigma_{\text{cindp}}^{(i,j)}$ the set of all CIND^ps in Σ from R_i to R_j , for $i, j \in [1, n]$. We show the following. (a) The violation set $\text{vio}(D, \Sigma_{\text{cfdp}}^i)$ can be computed by two SQL queries. (b) Similarly, $\text{vio}(D, \Sigma_{\text{cindp}}^{(i,j)})$ can be computed by a single SQL query. (c) These SQL queries encode pattern tableaux of CFD^ps (CIND^ps) with data tables, and hence their sizes are independent of Σ . From these Theorem 10 follows immediately.

5.1 Encoding CFD^ps and CIND^ps with Data Tables

We first show the following, by extending the encoding of [10], [23]. The pattern tableaux of all CFD^ps in Σ_{cfdp}^i can be encoded with *three data tables*, and the pattern tableaux of all CIND^ps in $\Sigma_{\text{cindp}}^{(i,j)}$ can be represented as *four data tables*, no matter how many dependencies are in the sets.

Encoding CFD^ps. We encode all pattern tableaux in Σ_{cfdp}^i with three tables enc_L , enc_R and enc_{\neq} , where enc_L

(1) enc_L					(2) enc_R					(3) enc_{\neq}			
cid	sale	price	price _{>}	price _{<}	cid	shipping	price	price _{>}	price _{<}	cid	pos	att	val
2	T	null	null	null	2	0	null	null	null				
3	F	—	20	40	3	6	null	null	null				
4	T	null	null	null	4	null	—	2.99	9.99				

Figure 4. Encoding example of CFD^p s

(1) enc			(2) enc_L			(3) enc_R		(4) enc_{\neq}			
cid	state _L	state _R	cid	type	state	cid	rate	cid	pos	att	val
1	1	1	1	—	null	1	null	1	LHS	type	art
2	1	1	2	—	DL	2	0	2	LHS	type	art

Figure 5. Encoding example of $CIND^p$ s

(resp. enc_R) encodes the non-negation ($=, <, \leq, >, \geq$) patterns in LHS (resp. RHS), and enc_{\neq} encodes those negation (\neq) patterns. More specifically, we associate a unique id cid with each CFD^p s in Σ_{cfdp}^i , and let enc_L consist of the following attributes: (a) cid , (b) each attribute A appearing in the LHS of some CFD^p s in Σ_{cfdp}^i , and (c) its four companion attributes $A_{>}, A_{\geq}, A_{<},$ and A_{\leq} . That is, for each attribute, there are five columns in enc_L , one for each non-negation operator. Similarly, enc_R is defined. We use an enc_{\neq} tuple to encode a pattern $A \neq c$ in a CFD^p , consisting of cid , att , pos , and val , encoding the CFD^p id, the attribute A , the position ('LHS' or 'RHS'), and the constant c , respectively. Note that the arity of enc_L (enc_R) is bounded by $5 * |R_i| + 1$, where $|R_i|$ is the arity of R_i , and the arity of enc_{\neq} is 4.

Before we populate these tables, let us first describe a preferred form of CFD^p s that would simplify the analysis to be given. Consider a CFD^p $\varphi = R(X \rightarrow Y, T_p)$. If φ is not satisfiable we can simply drop it from Σ . Otherwise it is equivalent to a CFD^p $\varphi' = R(X \rightarrow Y, T'_p)$ such that for any pattern tuples t_p, t'_p in T'_p and for any attribute A in $X \cup Y$, (a) if $t_p[A]$ is op a and $t'_p[A]$ is op b , where op is not \neq , then $a = b$, and (b) if $t_p[A]$ is ' $_$ ' then so is $t'_p[A]$. That is, for each non-negation op (resp. $_$), there is a *unique* constant a such that $t_p[A] = 'op$ a' (resp. $t_p[A] = _$) is the only op (resp. $_$) pattern appearing in the A column of T'_p . We refer to $t_p[A]$ as $T'_p(op, A)$ (resp. $T'_p(_, A)$), and consider *w.l.o.g.* CFD^p s of this form only. Note that there are possibly multiple $t_p[A] \neq c$ patterns in T'_p .

We populate enc_L , enc_R and enc_{\neq} as follows. For each CFD^p $\varphi = R(X \rightarrow Y, T_p)$ in Σ_{cfdp}^i , we generate a distinct cid_{φ} for it, and do the following.

- (1) Add a tuple t_1 to enc_L such that (a) $t[cid] = id_{\varphi}$; (b) for each $A \in X$, $t[A] = '_'$ if $T'_p(_, A)$ is ' $_$ ', and for each non-negation predicate op , $t[A_{op}] = 'a'$ if $T'_p(op, A)$ is ' op a' '; (c) we let $t[B] = \text{null}$ for all other attributes B in enc_L .
- (2) Similarly add a tuple t_2 to enc_R for attributes in Y .
- (3) For each attribute $A \in X \cup Y$ and each \neq a pattern in $T_p[A]$, add a tuple t to enc_{\neq} such that $t[cid] = id_{\varphi}$, $t[att] = 'A'$, $t[val] = 'a'$, and $t[pos] = 'LHS'$ (resp. $t[pos] = 'RHS'$) if attribute A appears in X (resp. Y).

Example 8: Recall from Fig. 2 CFD^p s φ_2 , φ_3 and φ_4 defined on relation item. The three CFD^p s are encoded with the tables shown in Fig. 4: (a) enc_L consists of attributes: cid , $sale$, $price$, $price_{>}$ and $price_{<}$; (b) enc_R consists of cid , $shipping$, $price$, $price_{>}$ and $price_{<}$; those attributes in a table with only

'null' pattern values do not contribute to error detection, and are thus omitted; And (c) enc_{\neq} is empty since all these CFD^p s have no negation patterns. One can easily reconstruct these CFD^p s from tables enc_L , enc_R and enc_{\neq} by collating the tuples based on cid . \square

Encoding $CIND^p$ s. All $CIND^p$ s in $\Sigma_{cindp}^{(i,j)}$ can be encoded with four tables enc , enc_L , enc_R and enc_{\neq} . Here enc_L (resp. enc_R) and enc_{\neq} encode non-negation patterns on relation R_i (resp. R_j) and negation patterns on relations R_i or R_j , respectively, along the same lines as their counterparts for CFD^p s. We use enc to encode the IND s embedded in $CIND^p$ s, which consists of the following attributes: (1) cid representing the id of a $CIND^p$, and (2) those X attributes of R_i and Y attributes of R_j appearing in some $CIND^p$ s in $\Sigma_{cindp}^{(i,j)}$. Note that the number of attributes in enc is bounded by $|R_i| + |R_j| + 1$, where $|R_i|$ is the arity of R_i .

For each $CIND^p$ $\psi = (R_i[A_1 \dots A_m; X_p] \subseteq R_j[B_1 \dots B_m; Y_p], T_p)$ in $\Sigma_{cindp}^{(i,j)}$, we generate a distinct cid_{ψ} for it, and do the following.

- (1) Add tuples t_1 and t_2 to enc_L and enc_R based on attributes X_p and Y_p , respectively, along the same lines as their CFD^p counterpart.
- (2) Add tuples to enc_{\neq} in the same way as their CFD^p counterparts.
- (3) Add tuple t to enc such that $t[cid] = id_{\psi}$. For each $k \in [1, m]$, let $t[A_k] = t[B_k] = k$, and $t[A] = \text{null}$ for the rest attributes A of enc .

Example 9: Figure 5 shows the coding of $CIND^p$ s ψ_1 and ψ_2 given in Fig. 3. We use $state_L$ and $state_R$ in enc to denote the occurrences of attribute state in item and tax, respectively. In enc_L and enc_R , the attributes with only 'null' patterns are omitted, for the same reason as CFD^p s mentioned above. \square

Putting these together, it is easy to verify that at most $O(n^2)$ data tables are needed to encode dependencies in Σ , regardless of the size of Σ . Recall that n is the number of relations in the database \mathcal{R} .

5.2 SQL-based Detection Methods

We next show how to generate SQL queries based on the encoding above. For each $i \in [1, n]$, we generate *two* SQL queries that, when evaluated on the I_i table of D , find $\text{vio}(D, \Sigma_{cfdp}^i)$. Similarly, for each $i, j \in [1, n]$, we generate a *single* SQL query $Q_{(i,j)}$ that, when evaluated on (I_i, I_j) of D ,

returns $\text{vio}(D, \Sigma_{\text{cindp}}^{(i,j)})$. Putting these query answers together, we get $\text{vio}(D, \Sigma)$, the violation set of D w.r.t. Σ .

SQL queries for CIND^ps. Below we show how the SQL query $Q_{(i,j)}$ is generated for validating CIND^ps in $\Sigma_{\text{cindp}}^{(i,j)}$, which has not been studied by previous work. For the lack of space, we put the generation of detection queries for CFD^ps in the supplementary material, which is an extension of the SQL techniques for CFDs and eCFDs discussed in [23] and [10], respectively.

The query $Q_{(i,j)}$ for the validation of $\Sigma_{\text{cindp}}^{(i,j)}$ is given as follows, which capitalizes on the data tables enc , enc_L , enc_R and enc_{\neq} that encode CIND^ps in $\Sigma_{\text{cindp}}^{(i,j)}$.

```
select Ri.* from Ri, encL L, enc≠ N
where Ri.X ≍ L and Ri.X ≍ N and not exists (
  select Rj.* from Rj, encH H, encR R, enc≠ N
  where Ri.X = Rj.Y and L.cid = R.cid and
    L.cid = H.cid and Rj.Y ≍ R and Rj.Y ≍ N)
```

Here (1) $X = \{A_1, \dots, A_{m_1}\}$ and $Y = \{B_1, \dots, B_{m_2}\}$ are the sets of attributes of R_i and R_j appearing in $\Sigma_{\text{cindp}}^{(i,j)}$, respectively; (2) $R_i.X \asymp L$ is the conjunction of

```
L.Ak is null or Ri.Ak = L.Ak or (L.Ak = ' ' and
(L.Ak > is null or Ri.Ak > L.Ak) and
L.Ak ≥ is null or Ri.Ak ≥ L.Ak) and
(L.Ak < is null or Ri.Ak < L.Ak) and
(L.Ak ≤ is null or Ri.Ak ≤ L.Ak))
```

for each $k \in [1, m_1]$; (3) $R_j.Y \asymp R$ is defined similarly for attributes in Y ; (4) $R_i.X \asymp N$ is a shorthand for the conjunction below, for each $k \in [1, m_1]$:

```
not exists (select * from N
  where L.cid = N.cid and N.pos = 'LHS' and
    N.att = 'Ak' and Ri.Ak = N.val);
```

(5) $R_j.Y \asymp N$ is defined similarly, but with $N.pos = \text{'RHS'}$; (6) $R_i.X = R_j.Y$ represents the following: for each A_k ($k \in [1, m_1]$) and each B_l ($l \in [1, m_2]$), $(H.A_k \text{ is null or } H.B_l \text{ is null or } H.B_l \neq H.A_k \text{ or } R_i.A_k = R_j.B_l)$.

Intuitively, (1) $R_i.X \asymp L$ and $R_i.X \asymp N$ ensure that the R_i tuples selected match the LHS patterns of some CIND^ps in $\Sigma_{\text{cindp}}^{(i,j)}$; (2) $R_j.Y \asymp R$ and $R_j.Y \asymp N$ check the corresponding RHS patterns of these CIND^ps on R_j tuples; (3) $R_i.X = R_j.Y$ enforces the *embedded* INDS; (4) $L.cid = R.cid$ and $L.cid = H.cid$ assure that the LHS and RHS patterns in the same CIND^p are correctly collated; and (5) **not exists** in $Q_{(i,j)}$ ensures that the R_i tuples selected violate CIND^ps in $\Sigma_{\text{cindp}}^{(i,j)}$.

Example 10: Using the coding of Fig. 5, an SQL query Q for checking CIND^ps ψ_1 and ψ_2 of Fig. 3 is given as follows:

```
select R1.* from item R1, encL L, enc≠ N
where (L.type is null or R1.type = L.type or L.type = ' ') and
  not exist (select * from N
    where N.cid = L.cid and N.pos = 'LHS' and
      N.att = 'type' and R1.type = N.val) and
  (L.state is null or R1.state = L.state or L.state = ' ') and
  not exist (select * from N
    where N.cid = L.cid and N.pos = 'LHS' and
      N.att = 'state' and R1.state = N.val) and
  not exist (select R2.* from tax R2, encH H, encR R
    where (H.stateL is null or H.stateR is null or
      H.stateL = H.stateR or R2.state = R1.state)
      and L.cid = H.cid and L.cid = R.cid and
      (R.rate is null or R2.rate = R.rate or
      R.rate = ' ') and not exist (select * from N
        where N.cid = R.cid and N.pos = 'RHS'
```

and $N.att = \text{'rate'}$ and $R_2.rate = N.val$))

The SQL queries generated can be simplified as follows. As shown in Example 10, when checking patterns imposed by enc , enc_L or enc_R , the queries need not consider attributes A if $t[A]$ is **null** for each tuple t in the table. Similarly, if an attribute A does not appear in any tuple in enc_{\neq} , the queries need not check A either. From this, it follows that we do not even need to generate those attributes with only **null** patterns for data tables enc , enc_L or enc_R when encoding CIND^ps or CFD^ps. \square

6 EXPERIMENTAL STUDY

We next present an extensive experimental study of CFD^ps and CIND^ps. Using real-life data, we conducted two sets of experiments to evaluate the efficiency and effectiveness of CFD^ps and CIND^ps vs. their counterparts CFDs and CINDs, separately and taken together.

6.1 Experimental Settings

We first present our experimental settings.

Datasets. We used two real-life datasets that were stored in an SQL Server 2012 database.

(1) HOSP (Hospital Compare) is a database publicly available from U.S. Department of Health & Human Services [1]. We used two tables hcahps and hcahps-state , which record the hospital level and state level ratings of the Hospital Consumer Assessment of Healthcare Providers and Systems (HCAHPS), respectively. For table hcahps , it records (a) the hospital information: hid (hospital ID), hname (hospital name), addr (address), city , state , zip , county , phn (phone number), and (b) the measure information: mid (measure ID), mq (question), mad (answer description), map (answer percentage), mnscs (number of completed surveys), msrrp (survey response rate percentage), mfn (footnote). And for table hcahps-state , it records state level measure information: state , mid , mq and map , among other things.

We designed 6 CFD^ps and 3 CIND^ps for HOSP, shown below in an informal way for easy of understanding.

```
φ1: hcahps (zip = ' ' and city = ' ' → state = ' ')
φ2: hcahps (hid = ' ' → hname = ' ' and county = ' ' and addr = ' ' and
  phn = ' ')
φ3: hcahps (hid = ' ' → msrrp = ' ')
φ4: hcahps (mid = ' ' → mq = ' ' and mad = ' ')
φ5: hcahps (hid = ' ' and mnscs = 'Not Available' → mfn ≥ 1 and mfn ≤ 14)
φ6: hcahps (hid = ' ' and mid = ' ' and mnscs ≠ 'Not Available' and
  mnscs ≠ 'Fewer than 100' → map ≥ 0 and map ≤ 100)
ψ1: hcahps (mid, mq; nil) ⊆ hcahps-state (mid, mq; nil)
ψ2: hcahps (mid, state; nil) ⊆ hcahps-state (mid, state; nil)
ψ3: hcahps (mid, state; mnscs ≠ 'Not Available') ⊆ hcahps-state (mid, state;
  map ≥ 0 and map ≤ 100)
```

For comparison, we also designed the CFDs and CINDs counterparts of the above CFD^ps and CIND^ps. Here φ_1 - φ_4 and ψ_1 - ψ_2 are indeed CFDs and CINDs, respectively, while φ_5 , φ_6 and ψ_3 are not. We hence further designed φ'_5 , φ'_6 and ψ'_3 to approximate φ_5 , φ_6 and ψ_3 , respectively.

```
φ'5: hcahps (hid = ' ' and mnscs = 'Not Available' → mfn = ' ')
φ'6: hcahps (hid = ' ' and mid = ' ' and mnscs = ' ' → map = ' ')
ψ'3: hcahps (mid, state; mnscs = ' ') ⊆ hcahps-state (mid, state; map = ' ')
```

(2) DBLP is a repository of computer science publications from 1946 to 2014 [2]. We further transformed its XML

format into two tables paper and proceeding that record the paper and proceeding information, respectively, such that paper (key, type, title, booktitle, year, crossref, isbn, publisher) records books, journal articles and conference papers, and proceeding (key, year, isbn, publisher) records the proceedings of conference papers, respectively.

We generated 3482 CFD^ps and 2568 CIND^ps for the DBLP data, with their representatives shown below.

ϕ_1 : paper (isbn = ' _ ' and booktitle = ' _ ' → publisher = ' _ ')
 ϕ_2 : paper (title = ' _ ' and year = ' _ ' and booktitle = ' _ ' → type = ' _ ')
 ϕ_3 : paper (booktitle = 'CleanDB' → year = 2006)
 ϕ_4 : paper (booktitle = 'VLDB' and year_L = ' _ ' → year_R ≥ 1975 and year_R ≤ 2007)
 ϕ_5 : paper (booktitle = 'PVLDB' and year_L = ' _ ' → year_R ≥ 2008)
 ρ_1 : paper (crossref, isbn, publisher; type = 'inproceedings' and booktitle = 'CIKM-CNIKM') ⊆ proceeding (key, isbn, publisher; year = 2009)
 ρ_2 : paper (crossref, isbn, publisher; type = 'inproceedings' and booktitle = 'VLDB') ⊆ proceeding (key, isbn, publisher; year ≥ 1975 and year ≤ 2007)
 ρ_3 : paper (crossref, isbn, publisher; type = 'inproceedings' and booktitle = 'ICDE') ⊆ proceeding (key, isbn, publisher; year ≥ 1984)

We collected all the booktitle and corresponding year from DBLP to generate the other CFD^ps and CIND^ps by instantiating the values of their booktitle and year attributes. Observe that ϕ_1 - ϕ_3 and ρ_1 are CFDs and CINDs, respectively. For comparison, we further designed the following CFDs and CINDs to approximate ϕ_4 - ϕ_5 and ρ_2 - ρ_3 .

ϕ'_1 : paper (booktitle = 'VLDB' and year_L = ' _ ' → year_R = ' _ ')
 ϕ'_2 : paper (booktitle = 'PVLDB' and year_L = ' _ ' → year_R = ' _ ')
 ρ'_2 : paper (crossref, isbn, publisher; type = 'inproceedings' and booktitle = 'VLDB') ⊆ proceeding (key, isbn, publisher; year = ' _ ')
 ρ'_3 : paper (crossref, isbn, publisher; type = 'inproceedings' and booktitle = 'ICDE') ⊆ proceeding (key, isbn, publisher; year = ' _ ')

Implementation. All the experiments were run within an SQL Server 2012 database installed on a machine with an Intel Core i5 (3.1GHz) CPU and 8GB of RAM. Each test was repeated 5 times, and the average is reported here.

6.2 Experimental Results

We next present our findings. Three parameters were used in our tests: (1) $|I_1|$, the number of tuples in table hcahps of HOSP or paper of DBLP, (2) $|I_2|$, the number of tuples in table hcahps-state of HOSP or proceeding of DBLP, and (3) *noise%*, the percentage of dirty tuples in table hcahps of HOSP or paper of DBLP, ranging from 0% to 9%. For easy of comparison, we deliberately dirty the tuples in hcahps of HOSP or paper of DBLP so that using the CFD^ps and CIND^ps together can detect all the dirty tuples. A clean copy of HOSP and DBLP is also kept to tell whether a tuple is dirty or clean.

6.2.1 Tests of Efficiency

In the first set of experiments, we evaluated the violation detection efficiency of CFD^ps and CIND^ps vs. their counterparts CFDs and CINDs, separately and taken together.

Exp-1.1: CFD^ps vs. CFDs. (1) To evaluate the impacts of $|I_1|$, we fixed *noise%* = 9%, and varied $|I_1|$ from 10K to 90K for HOSP (resp. from 100K to 900K for DBLP); And (2) to evaluate the impacts of *noise%*, we fixed $|I_1|$ = 90K for HOSP (resp. 900K for DBLP), and varied *noise%* from 0% to 9%. The results are reported in Figures 6(a) and 6(c) and Figures 6(b) and 6(d), respectively.

The results tell us that for CFDs and CFD^ps, both their running time (a) increases with the increment of the size of I_1 , and (b) is insensitive to the noise. Furthermore, (c) their running time is mainly affected by three factors: the size of I_1 , the LHS and RHS complexity of dependencies. For instance, (a) the LHS complexity of CFDs ϕ'_5 and ϕ'_6 is higher than CFD^ps ϕ_5 and ϕ_6 , as they match more I_1 tuples, but the RHS complexity of CFDs ϕ'_5 and ϕ'_6 is lower than CFD^ps ϕ_5 and ϕ_6 , as they are easier to check violations; And (b) the LHS complexity of CFDs ϕ'_4 and ϕ'_5 is the same as CFD^ps ϕ_4 and ϕ_5 , but the RHS complexity of CFDs ϕ'_4 and ϕ'_5 is similar to CFD^ps ϕ_4 and ϕ_5 , as they are easier to check violations. As a combined result, the running time of CFDs is lower than CFD^ps on HOSP, but close to CFD^ps on DBLP.

Exp-1.2: CIND^ps vs. CINDs. (1) To evaluate the impacts of $|I_1|$, we fixed *noise%* = 9% and $|I_2|$ = 1.6K for HOSP (resp. 16K for DBLP), and varied $|I_1|$ from 10K to 90K for HOSP (resp. from 100K to 900K for DBLP); (2) To evaluate the impacts of $|I_2|$, we fixed *noise%* = 9% and $|I_1|$ = 90K for HOSP (resp. 900K for DBLP), and varied $|I_2|$ from 1K to 1.6K for HOSP (resp. from 10K to 16K for DBLP); And (3) To evaluate the impacts of *noise%*, we fixed $|I_1|$ = 90K for HOSP (resp. 900K for DBLP) and $|I_2|$ = 1.6K for HOSP (resp. 16K for DBLP), and varied *noise%* from 0% to 9%. The results are reported in Figures 7(a) and 7(d), Figures 7(b) and 7(e), and Figures 7(c) and 7(f), respectively.

The results tell us that for CINDs and CIND^ps, both their running time (a) increases with the increment of the size of I_1 , (b) is not affected much by I_2 as $|I_2|$ is relatively small in the tests, and (c) is insensitive to the noise. Furthermore, (d) their running time is mainly affected by four factors: the size of I_1 , the size of I_2 , the LHS and RHS complexity of dependencies. For instance, (a) the LHS complexity of CIND ψ'_3 is higher than CFD^p ψ_3 , as they match more I_1 tuples, but the RHS complexity of CIND ψ'_3 is lower than CIND^p ψ_3 , as they are easier to check violations; And (b) the LHS complexity of CINDs ρ'_2 and ρ'_3 is the same as CFD^ps ϕ_4 and ϕ_5 , but the RHS complexity of CINDs ρ'_2 and ρ'_3 is lower than CIND^ps ρ_2 and ρ_3 , as they are easier to check violations. As a combined result, the running time of CINDs is close to CFD^ps on HOSP, but is lower on DBLP.

Exp-1.3: CFD^ps + CIND^ps vs. CFDs + CINDs. Using the same setting as Exp-1.2, we evaluated the impacts of $|I_1|$, $|I_2|$ and *noise%*. The results are reported in Figures 8(a) and 8(d), Figures 8(b) and 8(e) and Figures 8(c) and 8(f), respectively. The results show similar findings to Exp-1.1 and Exp-1.2, and are consistent with them.

6.2.2 Tests of Effectiveness

In the second set of experiments, we evaluated the violation detection effectiveness of CFD^ps and CIND^ps vs. their counterparts CFDs and CINDs, separately and taken together. Note that we did not report the results of varying $|I_2|$ as it has no impacts on the effectiveness tests in our setting.

Given one of CFDs, CFD^ps, CINDs, CIND^ps, CFDs + CINDs or CFD^ps + CIND^ps, denoted by x , its effectiveness of detecting violations is evaluated with the following measure:

$$\text{accuracy}(x) = \frac{\text{\#dirty tuples found by } x}{\text{\#dirty tuples found by cfdps + cindps}}.$$

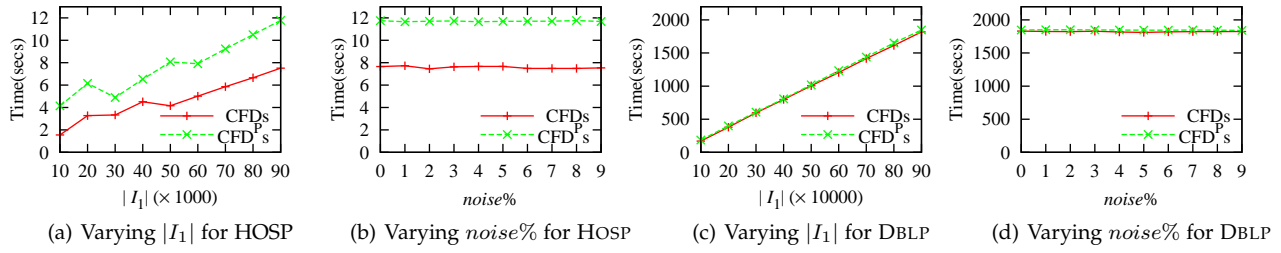


Figure 6. Efficiency of detecting violations: CFDPs vs. CFDs

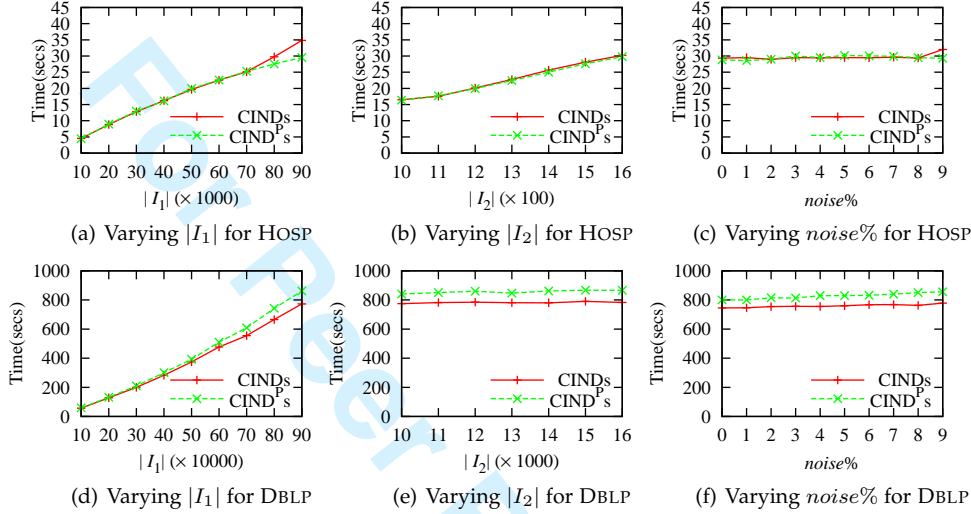


Figure 7. Efficiency of detecting violations: CINDPs vs. CINDs

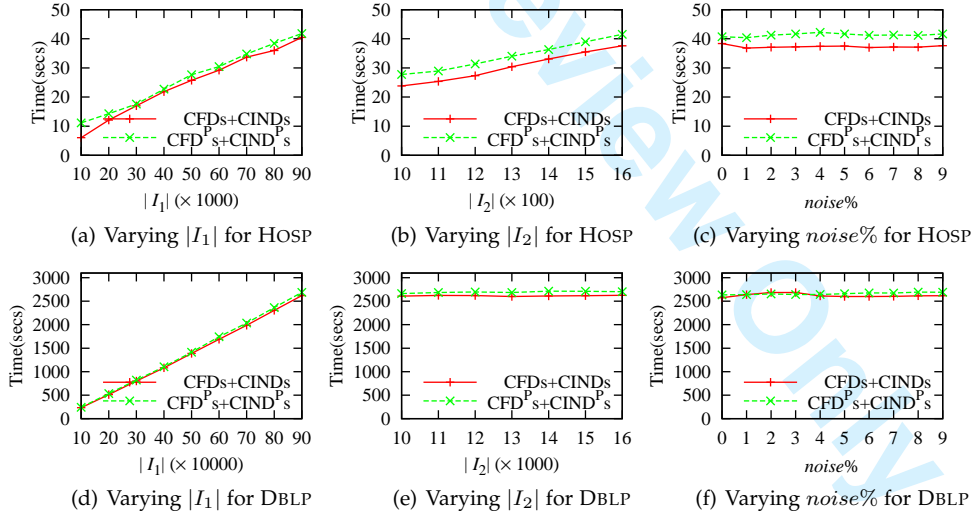


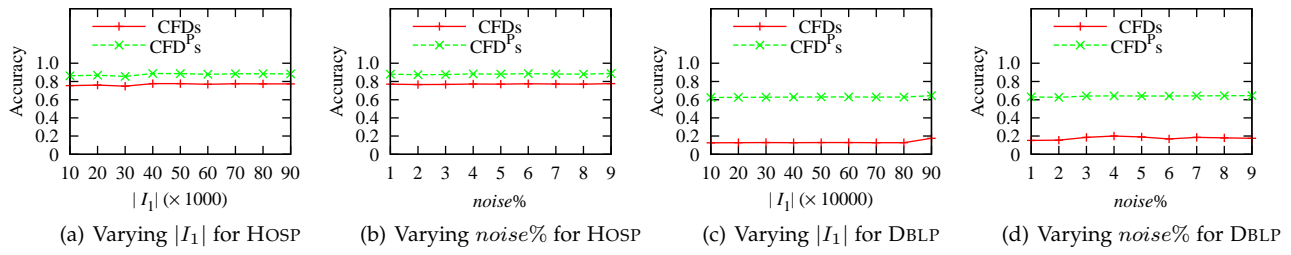
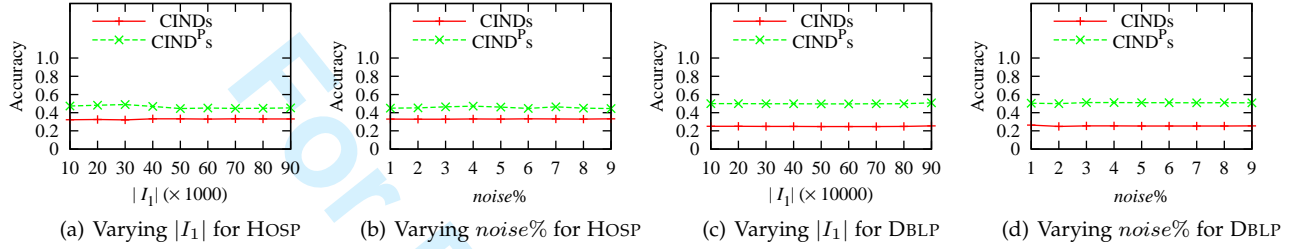
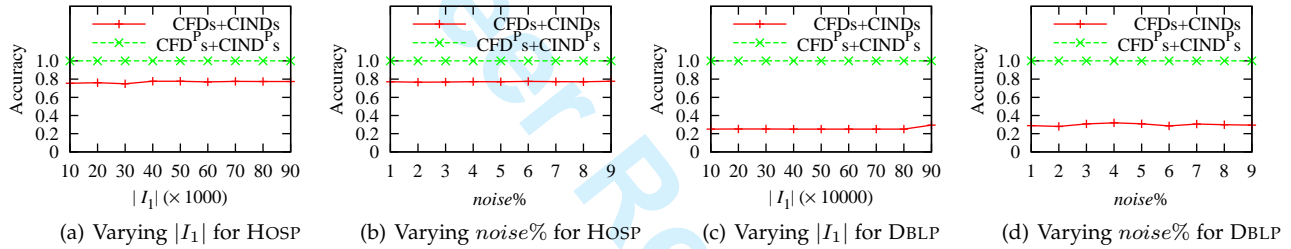
Figure 8. Efficiency of detecting violations: CFDPs + CINDPs vs. CFDs + CINDs

Exp-2. Using the same setting as Exp-1.1, Exp-1.2 and Exp-1.3, respectively, we evaluated the impacts of $|I_1|$ and $noise\%$ for (a) CFDPs vs. CFDs, (b) CINDPs vs. CINDs and (c) CFDPs + CINDPs vs. CFDs + CINDs, respectively. The results are reported in Figures 9, 10 and 11, respectively, and are summarized in Table 2.

The results tell us that (1) the effectiveness of detecting violations using all classes of dependencies are robust to $|I_1|$ and $noise\%$, (2) CFDPs, CINDPs and CFDPs + CINDPs obviously outperform their counterparts CFDs, CINDs and CFDs + CINDs, respectively, and (3) the increase of effectiveness

depends on the increase of the expressive power, and varies from 22% to 75% on HOSP and DBLP.

Summary. From these experimental results on real-life data HOSP and DBLP, we find the following. (1) The running time of CFDPs and CINDPs is comparable to their CFDs and CINDs counterparts, which is consistent with the static analyses: CFDPs and CINDPs retain the same complexity as their CFDs and CINDs counterparts. (2) CFDPs and CINDPs are able to capture more dirty tuples than CFDs and CINDs, due to the increased expressive power.

Figure 9. Effectiveness of detecting violations: CFD^p_s s vs. CFDsFigure 10. Effectiveness of detecting violations: $CIND^p_s$ s vs. CINDsFigure 11. Effectiveness of detecting violations: $CFD^p_s + CIND^p_s$ s vs. $CFDs + CINDs$ Table 2
Summary of violation detection accuracy

Datasets	varying	CFDs (%)	CFD^p_s (%)	CINDs (%)	$CIND^p_s$ (%)	CFDs + CINDs (%)	$CFD^p_s + CIND^p_s$ (%)
HOSP	$ I_1 $	[74.8, 77.7]	[85.5, 88.8]	[32.1, 33.3]	[44.6, 48.9]	[74.8, 77.7]	100
HOSP	$noise\%$	[76.6, 77.7]	[87.5, 88.8]	[32.8, 33.3]	[44.5, 47.2]	[76.6, 77.7]	100
DBLP	$ I_1 $	[12.6, 17.6]	[62.5, 64.5]	[24.8, 25.4]	[49.6, 50.9]	[25.0, 29.5]	100
DBLP	$noise\%$	[15.3, 20.1]	[62.5, 64.5]	[24.9, 26.3]	[50.0, 51.1]	[28.0, 31.9]	100

7 CONCLUSIONS

We have proposed CFD^p_s and $CIND^p_s$, which further extend CFDs and CINDs, respectively, by allowing patterns on data values to be expressed in terms of $\neq, <, \leq, >$ and \geq predicates. We have shown that CFD^p_s and $CIND^p_s$ are more powerful than CFDs and CINDs for detecting errors in real-life data. In addition, the satisfiability and implication problems for CFD^p_s and $CIND^p_s$ have the same complexity bounds as their counterparts for CFDs and CINDs, respectively. We have also provided automated methods to generate SQL queries for detecting errors based on CFD^p_s and $CIND^p_s$. These provide commercial DBMS with an immediate capability to capture errors commonly found in real-world data.

One topic for future work is to develop a dependency language that is capable of expressing various extensions of CFDs (e.g., CFD^p_s , $eCFDs$ [10] and CFD^c_s [13]), without increasing the complexity of static analyses. Second, we are to develop effective algorithms for discovering CFD^p_s

and $CIND^p_s$, along the same lines as [5], [30], [36]. Third, we plan to extend the methods of [8], [18] to repair data based on CFD^p_s and $CIND^p_s$, instead of using CFDs [18], traditional FDs and INDs [8], denial constraints [7], [16], and aggregate constraints [28]. Finally, distributed violation detection methods for CFD^p_s and $CIND^p_s$ are to be studied to deal with larger datasets, along the same lines as [24], [26].

ACKNOWLEDGMENTS

This work is supported in part by 973 program (No. 2014CB340300) and NSFC (No. 61322207). Fan is supported in part by 973 Program (No. 2012CB316200), NSFC (No. 61133002), Guangdong Innovative Research Team Program (2011D005), Shenzhen Peacock Program (1105100030834361) of China, and EPSRC (EP/J015377/1) of UK.

REFERENCES

- [1] <https://data.medicare.gov/data/hospital-compare>.

- [2] <http://dblp.uni-trier.de/xml/>.
- [3] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [4] M. Arenas, L. E. Bertossi, and J. Chomicki. Consistent query answers in inconsistent databases. In *PODS*, 1999.
- [5] J. Bauckmann, Z. Abedjan, U. Leser, H. Müller, and F. Naumann. Discovering conditional inclusion dependencies. In *CIKM*, 2012.
- [6] M. Baudinet, J. Chomicki, and P. Wolper. Constraint-Generating Dependencies. *J. Comput. Syst. Sci.*, 59(1):94–115, 1999.
- [7] L. E. Bertossi, L. Bravo, E. Franconi, and A. Lopatenko. The complexity and approximation of fixing numerical attributes in databases under integrity constraints. *Inf. Syst.*, 33(4-5):407–434, 2008.
- [8] P. Bohannon, W. Fan, M. Flaster, and R. Rastogi. A cost-based model and effective heuristic for repairing constraints by value modification. In *SIGMOD*, 2005.
- [9] P. D. Bra and J. Paredaens. Conditional dependencies for horizontal decompositions. In *ICALP*, 1983.
- [10] L. Bravo, W. Fan, F. Geerts, and S. Ma. Increasing the expressivity of conditional functional dependencies without extra complexity. In *ICDE*, 2008.
- [11] M. A. Casanova, R. Fagin, and C. H. Papadimitriou. Inclusion dependencies and their interaction with functional dependencies. *J. Comput. Syst. Sci.*, 28(1):29–59, 1984.
- [12] W. Chen, W. Fan, and S. Ma. Analyses and validation of conditional dependencies with built-in predicates. In *DEXA*, 2009.
- [13] W. Chen, W. Fan, and S. Ma. Incorporating cardinality constraints and synonym rules into conditional functional dependencies. *IPL*, 109(14):783–789, 2009.
- [14] F. Chiang and R. J. Miller. Discovering data quality rules. In *VLDB*, 2008.
- [15] J. Chomicki. Consistent query answering: Five easy pieces. In *ICDT*, 2007.
- [16] J. Chomicki and J. Marcinkowski. Minimal-change integrity maintenance using tuple deletions. *Inf. Comput.*, 197(1-2):90–121, 2005.
- [17] E. F. Codd. Relational completeness of data base sublanguages. In *Data Base Systems: Courant Computer Science Symposia Series 6*, pages 65–98. Prentice-Hall, 1972.
- [18] G. Cong, W. Fan, F. Geerts, X. Jia, and S. Ma. Improving data quality: Consistency and accuracy. In *VLDB*, 2007.
- [19] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. The MIT Press, 2001.
- [20] O. Curé. Improving the data quality of drug databases using conditional dependencies and ontologies. *J. Data and Information Quality*, 4(1):3, 2012.
- [21] W. Fan, H. Gao, X. Jia, J. Li, and S. Ma. Dynamic constraints for record matching. *VLDB J.*, 20(4):495–520, 2011.
- [22] W. Fan and F. Geerts. *Foundations of Data Quality Management*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2012.
- [23] W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies. *TODS*, 33(2), 2008.
- [24] W. Fan, F. Geerts, S. Ma, and H. Müller. Detecting inconsistencies in distributed data. In *ICDE*, 2010.
- [25] W. Fan, J. Li, S. Ma, N. Tang, and W. Yu. Towards certain fixes with editing rules and master data. *VLDB J.*, 21(2):213–238, 2012.
- [26] W. Fan, J. Li, N. Tang, and W. Yu. Incremental detection of inconsistencies in distributed data. In *ICDE*, 2012.
- [27] W. Fan, S. Ma, N. Tang, and W. Yu. Interaction between record matching and data repairing. *J. Data and Information Quality*, 4(4):16, 2014.
- [28] S. Flesca, F. Furfaro, and F. Parisi. Consistent query answers on numerical databases under aggregate constraints. In *DBPL*, 2005.
- [29] M. Garey and D. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Company, 1979.
- [30] L. Golab, H. J. Karloff, F. Korn, D. Srivastava, and B. Yu. On generating near-optimal tableaux for conditional functional dependencies. In *VLDB*, 2008.
- [31] S. Ma, W. Fan, and L. Bravo. Extending inclusion dependencies with conditions. *Theor. Comput. Sci.*, 515:64–95, 2014.
- [32] M. J. Maher. Constrained dependencies. *TCS*, 173(1):113–149, 1997.
- [33] M. J. Maher and D. Srivastava. Chasing Constrained Tuple-Generating Dependencies. In *PODS*, 1996.
- [34] B. Saha and D. Srivastava. Data quality: The other face of big data. In *ICDE*, 2014.
- [35] R. van der Meyden. The complexity of querying indefinite data about linearly ordered domains. *JCSS*, 54(1), 1997.
- [36] A. Zanzi and A. Trombetta. Discovering non-constant conditional functional dependencies with built-in predicates. In *DEXA*, 2014.

PLACE
PHOTO
HERE

Shuai Ma is a professor at the School of Computer Science and Engineering, Beihang University. He obtained his PhD degrees from University of Edinburgh in 2010, and from Peking University in 2004, respectively. He was a post-doctoral research fellow in the database group, University of Edinburgh, a summer intern at Bell labs, Murray Hill, USA, in the summer of 2008, and a visiting researcher of MRSA in 2012. He is a recipient of the Best Paper Award for VLDB 2010 and the Best Challenge Paper Award for WISE 2013. His current research interests include database theory and systems, social data analysis and data intensive computing.

PLACE
PHOTO
HERE

Liang Duan is a PhD student in the School of Computer Science and Engineering, Beihang University, supervised by Prof. Shuai Ma. He received his MS degree in computer software and theory from Yunnan University in 2014, and BS degree in computer science and technology from Beihang University in 2009. His current research interests include databases and social data analysis.

PLACE
PHOTO
HERE

Wenfei Fan is Professor (Chair) of Web Data Management in the School of Informatics, University of Edinburgh, UK. He is a Fellow of ACM, a Fellow of the Royal Society of Edinburgh, UK, a National Professor of the Thousand-Talent Program and a Yangtze River Scholar, China. He received his PhD degree from the University of Pennsylvania. He is a recipient of the Runner-up for Best Paper Award of ICDE 2014, the Alberto O. Mendelzon Test-of-Time Award of ACM PODS 2010, the Best Paper Award for VLDB 2010, the Roger Needham Award in 2008 (UK), the Best Paper Award for ICDE 2007, the Best Paper of the Year Award for Computer Networks in 2002, and the Career Award in 2001 (USA). His current research interests include database theory and systems, in particular data quality, data integration, distributed query processing, query languages, recommender systems, social networks and Web services.

PLACE
PHOTO
HERE

Chunming Hu is an associate professor at the School of Computer Science and Engineering, Beihang University. He received his PhD degree from Beihang University in 2006. His current research interests include distributed systems, system virtualization, large scale data management and processing systems.

PLACE
PHOTO
HERE

Wenguang Chen is an associate professor at the Department of Computer Science, Peking University. He obtained his PhD degree from Peking University in 2009. He was a visiting scholar at University of Alberta from 2011 to 2012 and at University of Hawaii at Manoa in from 2012 to 2013. His current research interests include data management, data quality and intelligent HCI.

Analyses and Validation of Conditional Dependencies with Built-in Predicates

Wenguang Chen¹, Wenfei Fan^{2,3}, and Shuai Ma²

¹ Peking University, China

² University of Edinburgh, UK

³ Bell Laboratories, USA

Abstract. This paper proposes a natural extension of conditional functional dependencies (CFDs [14]) and conditional inclusion dependencies (CINDs [8]), denoted by CFD^P s and $CIND^P$ s, respectively, by specifying patterns of data values with $\neq, <, \leq, >$ and \geq predicates. As data quality rules, CFD^P s and $CIND^P$ s are able to capture errors that commonly arise in practice but cannot be detected by CFDs and CINDs. We establish two sets of results for central technical problems associated with CFD^P s and $CIND^P$ s. (a) One concerns the satisfiability and implication problems for CFD^P s and $CIND^P$ s, taken separately or together. These are important for, *e.g.*, deciding whether data quality rules are dirty themselves, and for removing redundant rules. We show that despite the increased expressive power, the static analyses of CFD^P s and $CIND^P$ s retain the same complexity as their CFDs and CINDs counterparts. (b) The other concerns validation of CFD^P s and $CIND^P$ s. We show that given a set Σ of CFD^P s and $CIND^P$ s on a database D , a set of SQL queries can be automatically generated that, when evaluated against D , return all tuples in D that violate some dependencies in Σ . This provides commercial DBMS with an immediate capability to detect errors based on CFD^P s and $CIND^P$ s.

Key words: functional dependency, inclusion dependency, data quality

1 Introduction

Extensions of functional dependencies (FDs) and inclusion dependencies (INDs), known as *conditional functional dependencies* (CFDs [14]) and *conditional inclusion dependencies* (CINDs [8]), respectively, have recently been proposed for improving data quality. These extensions enforce patterns of semantically related data values, and detect errors as violations of the dependencies. Conditional dependencies are able to capture more inconsistencies than FDs and INDs [14, 8].

Conditional dependencies specify constant patterns in terms of equality ($=$). In practice, however, the semantics of data often needs to be specified in terms of other predicates such as $\neq, <, \leq, >$ and \geq , as illustrated by the example below.

Example 1. An online store maintains a database of two relations: (a) *item* for items sold by the store, and (b) *tax* for the sale tax rates for the items, except artwork, in various states. The relations are specified by the following schemas:

item (id: string, name: string, type: string, price: float, shipping: float,
sale: bool, state: string)
tax (state: string, rate: float)

Wenguang Chen, Wenfei Fan, and Shuai Ma

	id	name	type	price	shipping	sale	state		state	rate
t_1 :	b1	Harry Potter	book	25.99	0	T	WA	t_5 :	PA	6
t_2 :	c1	Snow White	CD	9.99	2	F	NY	t_6 :	NY	4
t_3 :	b2	Catch-22	book	34.99	20	F	DL	t_7 :	DL	0
t_4 :	a1	Sunflowers	art	5m	500	F	DL	t_8 :	NJ	3.5

(a) An item relation

(b) tax rates

Fig. 1. Example instance D_0 of item and tax

where each item is specified by its id, name, type (*e.g.*, book, CD), price, shipping fee, the state to which it is shipped, and whether it is on sale. A tax tuple specifies the sale tax rate in a state. An instance D_0 of item and tax is shown in Fig. 1.

One wants to specify dependencies on the relations as data quality rules to detect errors in the data, such that inconsistencies emerge as violations of the dependencies. Traditional dependencies (FDs, INDs; see, *e.g.*, [1]) and conditional dependencies (CFDs, CINDs [14, 8]) on the data include the following:

- cfd_1 : item ($\text{id} \rightarrow \text{name, type, price, shipping, sale}$)
- cfd_2 : tax ($\text{state} \rightarrow \text{rate}$)
- cfd_3 : item ($\text{sale} = \text{'T'} \rightarrow \text{shipping} = 0$)

These are CFDs: (a) cfd_1 assures that the id of an item uniquely determines the name, type, price, shipping, sale of the item; (b) cfd_2 states that state is a key for tax, *i.e.*, for each state there is a unique sale tax rate; and (c) cfd_3 is to ensure that for any item tuple t , if $t[\text{sale}] = \text{'T'}$ then $t[\text{shipping}]$ must be 0; *i.e.*, the store provides free shipping for items on sale. Here cfd_3 is specified in terms of patterns of semantically related data values, namely, $\text{sale} = \text{'T'}$ and $\text{shipping} = 0$. It is to hold only on item tuples that match the pattern $\text{sale} = \text{'T'}$. In contrast, cfd_1 and cfd_2 are traditional FDs without constant patterns, a special case of CFDs. One can verify that no sensible INDs or CINDs can be defined across item and tax.

Note that D_0 of Fig. 1 satisfies cfd_1 , cfd_2 and cfd_3 . That is, when these dependencies are used as data quality rules, no errors are found in D_0 .

In practice, the shipment fee of an item is typically determined by the price of the item. Moreover, when an item is on sale, the price of the item is often in a certain range. Furthermore, for any item sold by the store to a customer in a state, if the item is *not* artwork, then one expects to find the sale tax rate in the state from the tax table. These semantic relations cannot be expressed as CFDs of [14] or CINDs of [8], but can be expressed as the following dependencies:

- pfd_1 : item ($\text{sale} = \text{'F'} \ \& \ \text{price} \leq 20 \rightarrow \text{shipping} = 3$)
- pfd_2 : item ($\text{sale} = \text{'F'} \ \& \ \text{price} > 20 \ \& \ \text{price} \leq 40 \rightarrow \text{shipping} = 6$)
- pfd_3 : item ($\text{sale} = \text{'F'} \ \& \ \text{price} > 40 \rightarrow \text{shipping} = 10$)
- pfd_4 : item ($\text{sale} = \text{'T'} \rightarrow \text{price} \geq 2.99 \ \& \ \text{price} < 9.99$)
- pind_1 : item ($\text{state; type} \neq \text{'art'} \subseteq \text{tax (state; nil)}$)

Here pfd_2 states that for any item tuple, if it is not on sale and its price is in the range (20, 40], then its shipment fee must be 6; similarly for pfd_1 and pfd_3 . These dependencies extend CFDs [14] by specifying patterns of semantically related data values in terms of predicates $<$, \leq , $>$, and \geq . Similarly, pfd_4 assures that for any item tuple, if it is on sale, then its price must be in the range [2.99, 9.99). Dependency pind_1 extends CINDs [8] by specifying patterns with \neq : for any item

tuple t , if $t[\text{type}]$ is *not* artwork, then there must exist a **tax** tuple t' such that $t[\text{state}] = t'[\text{state}]$, *i.e.*, the sale tax of the item can be found from the **tax** relation.

Using pfd_1 – pfd_4 and pind_1 as data quality rules, we find that D_0 of Fig. 1 is *not* clean. Indeed, (a) t_2 violates pfd_1 : its price is less than 20, but its shipping fee is 2 rather than 3; similarly, t_3 violates pfd_2 , and t_4 violates pfd_3 . (b) Tuple t_1 violates pfd_4 : it is on sale but its price is not in the range $[2.99, 9.99]$. (c) The database D_0 also violates pind_1 : t_1 is not artwork, but its state cannot find a match in the **tax** relation, *i.e.*, no tax rate for WA is found in D_0 . \square

None of pfd_1 – pfd_4 and pind_1 can be expressed as FDs or INDs [1], which do not allow constants, or as CFDs [14] or CINDs [8], which specify patterns with equality ($=$) only. While there have been extensions of CFDs [7, 18], none of these allows dependencies to be specified with patterns on data values in terms of built-in predicates $\neq, <, \leq, >$ or \geq . To the best of our knowledge, no previous work has studied extensions of CINDs (see Section 6 for detailed discussions).

These highlight the need for extending CFDs and CINDs to capture errors commonly found in real-life data. While one can consider arbitrary extensions, it is necessary to strike a balance between the expressive power of the extensions and their complexity. In particular, we want to be able to reason about data quality rules expressed as extended CFDs and CINDs. Furthermore, we want to have effective algorithms to detect inconsistencies based on these extensions.

Contributions. This paper proposes a natural extension of CFDs and CINDs, provides complexity bounds for reasoning about the extension, and develops effective SQL-based techniques for detecting errors based on the extension.

(1) We propose two classes of dependencies, denoted by CFD^p s and CIND^p s, which respectively extend CFDs and CINDs by supporting $\neq, <, \leq, >$, \geq predicates. For example, all the dependencies we have encountered so far can be expressed as CFD^p s or CIND^p s. These dependencies are capable of capturing errors in real-world data that cannot be detected by CFDs or CINDs.

(2) We establish complexity bounds for the satisfiability problem and the implication problem for CFD^p s and CIND^p s, taken separately or together. The satisfiability problem is to determine whether a set Σ of dependencies has a nonempty model, *i.e.*, whether the rules in Σ are consistent themselves. The implication problem is to decide whether a set Σ of dependencies entails another dependency φ , *i.e.*, whether the rule φ is redundant in the presence of the rules in Σ . These are the central technical problems associated with any dependency language.

We show that despite the increased expressive power, CFD^p s and CIND^p s do not increase the complexity for reasoning about them. In particular, we show that the satisfiability and implication problems remain (a) NP-complete and coNP-complete for CFD^p s, respectively, (b) in $O(1)$ -time (constant-time) and EXPTIME-complete for CIND^p s, respectively, and (c) are undecidable when CFD^p s and CIND^p s are taken together. These are *the same as* their CFDs and CINDs counterparts. While data with linearly ordered domains often makes our lives harder (see, *e.g.*, [21]), CFD^p s and CIND^p s do not complicate their static analyses.

(3) We provide SQL-based techniques to detect errors based on CFD^p s and CIND^p s.

Wenguang Chen, Wenfei Fan, and Shuai Ma

$$\begin{array}{ll}
 (1) \varphi_1 = \text{tax}(\text{state} \rightarrow \text{rate}, T_1) & (2) \varphi_2 = \text{item}(\text{sale} \rightarrow \text{shipping}, T_2) \\
 T_1: \begin{array}{c|c} \text{state} & \text{rate} \\ \hline - & - \end{array} & T_2: \begin{array}{c|c} \text{sale} & \text{shipping} \\ \hline = T & = 0 \end{array} \\
 (3) \varphi_3 = \text{item}(\text{sale}, \text{price} \rightarrow \text{shipping}, T_3) & (4) \text{CFD}^P \varphi_4 = \text{item}(\text{sale} \rightarrow \text{price}, T_4) \\
 T_3: \begin{array}{c|c|c} \text{sale} & \text{price} & \text{shipping} \\ \hline = F & > 20 & = 6 \\ = F & \leq 40 & = 6 \end{array} & T_4: \begin{array}{c|c} \text{sale} & \text{price} \\ \hline = T & \geq 2.99 \\ = T & < 9.99 \end{array}
 \end{array}$$

Fig. 2. Example CFD^P s

Given a set Σ of CFD^P s and CIND^P s on a database D , we automatically generate a set of SQL queries that, when evaluated on D , find all tuples in D that violate some dependencies in Σ . Further, the SQL queries are independent of the size and cardinality of Σ . No previous work has been studied error detection based on CIND s, not to mention CFD^P s and CIND^P s taken together. These provide the capability of detecting errors in a single relation (CFD^P s) and across different relations (CIND^P s) within the immediate reach of commercial DBMS.

Organizations. Sections 2 and 3 introduce CFD^P s and CIND^P s, respectively. Section 4 establishes complexity bounds for reasoning about CFD^P s and CIND^P s. Section 5 provides SQL techniques for error detection. Related work is discussed in Section 6, followed by topics for future work in Section 7.

2 Incorporating Built-in Predicates into CFDs

We now define CFD^P s, also referred to as *conditional functional dependencies*, by extending CFDs with predicates ($\neq, <, \leq, >, \geq$) in addition to equality ($=$).

Consider a relation schema R defined over a finite set of attributes, denoted by $\text{attr}(R)$. For each attribute $A \in \text{attr}(R)$, its domain is specified in R , denoted as $\text{dom}(A)$, which is either finite (e.g., `bool`) or infinite (e.g., `string`). We assume w.l.o.g. that a domain is totally ordered if $<, \leq, >$ or \geq is defined on it.

Syntax. A $\text{CFD}^P \varphi$ on R is a pair $R(X \rightarrow Y, T_p)$, where (1) X, Y are sets of attributes in $\text{attr}(R)$; (2) $X \rightarrow Y$ is a standard FD, referred to as the FD *embedded in* φ ; and (3) T_p is a tableau with attributes in X and Y , referred to as the *pattern tableau* of φ , where for each A in $X \cup Y$ and each tuple $t_p \in T_p$, $t_p[A]$ is either an unnamed variable ‘ $_$ ’ that draws values from $\text{dom}(A)$, or ‘ $\text{op } a$ ’, where op is one of $=, \neq, <, \leq, >, \geq$, and ‘ a ’ is a constant in $\text{dom}(A)$.

If attribute A occurs in both X and Y , we use A_L and A_R to indicate the occurrence of A in X and Y , respectively, and separate the X and Y attributes in a pattern tuple with ‘ $\|$ ’. We write φ as $(X \rightarrow Y, T_p)$ when R is clear from the context, and denote X as $\text{LHS}(\varphi)$ and Y as $\text{RHS}(\varphi)$.

Example 2. The dependencies `cfid1`–`cfid3` and `pfid1`–`pfid4` that we have seen in Example 1 can all be expressed as CFD^P s. Figure 2 shows some of these CFD^P s: φ_1 (for FD `cfid2`), φ_2 (for CFD `cfid3`), φ_3 (for `pfid2`), and φ_4 (for `pfid4`). \square

Semantics. Consider $\text{CFD}^P \varphi = (R : X \rightarrow Y, T_p)$, where $T_p = \{t_{p1}, \dots, t_{pk}\}$.

A data tuple t of R is said to *match* $\text{LHS}(\varphi)$, denoted by $t[X] \asymp T_p[X]$, if for each tuple t_{pi} in T_p and each attribute A in X , either (a) $t_{pi}[A]$ is the wildcard ‘ $_$ ’

(which matches any value in $\text{dom}(A)$), or (b) $t[A] \text{ op } a$ if $t_{pi}[A]$ is ‘ $\text{op } a$ ’, where the operator op ($=, \neq, <, \leq, >$ or \geq) is interpreted by its standard semantics. Similarly, the notion that t matches $\text{RHS}(\varphi)$ is defined, denoted by $t[Y] \asymp T_p[Y]$.

Intuitively, each pattern tuple t_{pi} specifies a condition via $t_{pi}[X]$, and $t[X] \asymp T_p[X]$ if $t[X]$ satisfies the *conjunction* of all these conditions. Similarly, $t[Y] \asymp T_p[Y]$ if $t[Y]$ matches all the patterns specified by $t_{pi}[Y]$ for all t_{pi} in T_p .

An instance I of R satisfies the CFD^p φ , denoted by $I \models \varphi$, if for *each pair* of tuples t_1, t_2 in the instance I , if $t_1[X] = t_2[X] \asymp T_p[X]$, then $t_1[Y] = t_2[Y] \asymp T_p[Y]$. That is, if $t_1[X]$ and $t_2[X]$ are equal and in addition, they both match the pattern tableau $T_p[X]$, then $t_1[Y]$ and $t_2[Y]$ must also be equal to each other and they both match the pattern tableau $T_p[Y]$.

Observe that φ is imposed only on the subset of tuples in I that match $\text{LHS}(\varphi)$, rather than on the entire I . For all tuples t_1, t_2 in this subset, if $t_1[X] = t_2[X]$, then (a) $t_1[Y] = t_2[Y]$, *i.e.*, the semantics of the embedded FDs is enforced; and (b) $t_1[Y] \asymp T_p[Y]$, which assures that the *constants* in $t_1[Y]$ match the *constants* in $t_{pi}[Y]$ for all t_{pi} in T_p . Note that here tuples t_1 and t_2 can be the same.

An instance I of R satisfies a set Σ of CFD^ps, denoted by $I \models \Sigma$, if $I \models \varphi$ for *each* CFD^p φ in Σ .

Example 3. The instance D_0 of Fig. 1 satisfies φ_1 and φ_2 of Fig. 2, but neither φ_3 nor φ_4 . Indeed, tuple t_3 violates (*i.e.*, does not satisfy) φ_3 , since $t_3[\text{sale}] = \text{‘F’}$ and $20 < t_3[\text{price}] \leq 40$, but $t_3[\text{shipping}]$ is 20 instead of 6. Note that t_3 matches $\text{LHS}(\varphi_3)$ since it satisfies the condition specified by the *conjunction* of the pattern tuples in T_3 . Similarly, t_1 violates φ_4 , since $t_1[\text{sale}] = \text{‘T’}$ but $t_1[\text{price}] > 9.99$. Observe that while it takes two tuples to violate a standard FD, a single tuple may violate a CFD^p. \square

Special cases. (1) A standard FD $X \rightarrow Y$ [1] can be expressed as a CFD $(X \rightarrow Y, T_p)$ in which T_p contains a single tuple consisting of ‘.’ only, without constants. (2) A CFD $(X \rightarrow Y, T_p)$ [14] with $T_p = \{t_{p1}, \dots, t_{pk}\}$ can be expressed as a set $\{\varphi_1, \dots, \varphi_k\}$ of CFD^ps such that for $i \in [1, k]$, $\varphi_i = (X \rightarrow Y, T_{pi})$, where T_{pi} contains a single pattern tuple t_{pi} of T_p , with equality ($=$) only. For example, φ_1 and φ_2 in Fig. 2 are CFD^ps representing FD *cfid2* and CFD *cfid3* in Example 1, respectively. Note that all data quality rules in [10, 18] can be expressed as CFD^ps.

3 Incorporating Built-in Predicates into CINDs

Along the same lines as CFD^ps, we next define CIND^ps, also referred to as *conditional inclusion dependencies*. Consider two relation schemas R_1 and R_2 .

Syntax. A CIND^p ψ is a pair $(R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$, where (1) X, X_p and Y, Y_p are lists of attributes in $\text{attr}(R_1)$ and $\text{attr}(R_2)$, respectively; (2) $R_1[X] \subseteq R_2[Y]$ is a standard IND, referred to as the IND *embedded* in ψ ; and (3) T_p is a tableau, called the *pattern tableau* of ψ defined over attributes $X_p \cup Y_p$, and for each A in X_p or Y_p and each pattern tuple $t_p \in T_p$, $t_p[A]$ is either an

6 Wenguang Chen, Wenfei Fan, and Shuai Ma

$$\begin{aligned}
 (1) \quad \psi_1 &= (\text{item} [\text{state}; \text{type}] \subseteq \text{tax} [\text{state}; \text{nil}], T_1), \\
 (2) \quad \psi_2 &= (\text{item} [\text{state}; \text{type}, \text{state}] \subseteq \text{tax} [\text{state}; \text{rate}], T_2) \\
 T_1: \quad &\frac{\text{type} \parallel \text{nil}}{\neq \text{art}} \quad T_2: \quad \frac{\text{type} \parallel \text{state} \parallel \text{rate}}{\neq \text{art} \parallel = \text{DL} \parallel = 0}
 \end{aligned}$$

Fig. 3. Example CIND^ps

unnamed variable ‘ \cdot ’ that draws values from $\text{dom}(A)$, or ‘ $\text{op } a$ ’, where op is one of $=, \neq, <, \leq, >, \geq$ and ‘ a ’ is a constant in $\text{dom}(A)$.

We denote $X \cup X_p$ as $\text{LHS}(\psi)$ and $Y \cup Y_p$ as $\text{RHS}(\psi)$, and separate the X_p and Y_p attributes in a pattern tuple with ‘ \parallel ’. We use nil to denote an *empty* list.

Example 4. Figure 3 shows two example CIND^ps: ψ_1 expresses pind_1 of Example 1, and ψ_2 refines ψ_1 by stating that for any *item* tuple t_1 , if its *type* is not *art* and its *state* is *DL*, then there must be a *tax* tuple t_2 such that its *state* is *DL* and *rate* is 0, i.e., ψ_2 assures that the sale tax rate in Delaware is 0. \square

Semantics. Consider CIND^p $\psi = (R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$. An instance (I_1, I_2) of (R_1, R_2) satisfies the CIND^p ψ , denoted by $(I_1, I_2) \models \psi$, iff for each tuple $t_1 \in I_1$, if $t_1[X_p] \simeq T_p[X_p]$, then there exists a tuple $t_2 \in I_2$ such that $t_1[X] = t_2[Y]$ and moreover, $t_2[Y_p] \simeq T_p[Y_p]$.

That is, if $t_1[X_p]$ matches the pattern tableau $T_p[X_p]$, then ψ requires the existence of t_2 such that (1) $t_1[X] = t_2[Y]$ as required by the standard IND embedded in ψ ; and (2) $t_2[Y_p]$ must match the pattern tableau $T_p[Y_p]$. In other words, ψ is “conditional” since its embedded IND is applied only to the subset of tuples in I_1 that match $T_p[X_p]$, and moreover, the pattern $T_p[Y_p]$ is enforced on the tuples in I_2 that match those tuples in I_1 . As remarked in Section 2, the pattern tableau T_p specifies the *conjunction* of patterns of all tuples in T_p .

Example 5. The instance D_0 of *item* and *tax* in Fig. 1 violates CIND^p ψ_1 . Indeed, tuple t_1 in *item* matches $\text{LHS}(\psi_1)$ since $t_1[\text{type}] \neq \text{‘art’}$, but there is no tuple t in *tax* such that $t[\text{state}] = t_1[\text{state}] = \text{‘WA’}$. In contrast, D_0 satisfies ψ_2 . \square

We say that a database D satisfies a set Σ of CINDs, denoted by $D \models \Sigma$, if $D \models \varphi$ for each $\varphi \in \Sigma$.

Safe CIND^ps. We say a CIND^p $(R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$ is *unsafe* if there exist pattern tuples t_p, t'_p in T_p such that either (a) there exists $B \in Y_p$, such that $t_p[B]$ and $t'_p[B]$ are not satisfiable when taken together, or (b) there exist $C \in Y, A \in X$ such that A corresponds to B in the IND and $t_p[C]$ and $t'_p[A]$ are not satisfiable when taken together; e.g., $t_p[\text{price}] = 9.99$ and $t'_p[\text{price}] \geq 19.99$.

Obviously unsafe CIND^ps do not make sense: there exist no nonempty database that satisfies unsafe CIND^ps. It takes $O(|T_p|^2)$ -time in the size $|T_p|$ of T_p to decide whether a CIND^p is unsafe. Thus in the sequel we consider safe CIND^p only.

Special cases. Observe that (1) a standard CIND $(R_1[X] \subseteq R_2[Y])$ can be expressed as a CIND^p $(R_1[X; \text{nil}] \subseteq R_2[Y; \text{nil}], T_p)$ such that T_p is simply a *empty* set; and (2) a CIND $(R_1[X; X_p] \subseteq R_2[Y; Y_p], T_p)$ with $T_p = \{t_{p1}, \dots, t_{pk}\}$ can be expressed as a set $\{\psi_1, \dots, \psi_k\}$ of CIND^ps, where for $i \in [1, k]$, $\psi_i = (R_1[X; X_p] \subseteq R_2[Y; Y_p], T_{pi})$ such that T_{pi} consists of a single pattern tuple t_{pi} of T_p defined in terms of equality ($=$) only.

4 Reasoning about CFD^ps and CIND^ps

The satisfiability problem and the implication problem are the two central technical questions associated with any dependency languages. In this section we investigate these problems for CFD^ps and CIND^ps, separately and taken together.

4.1 The Satisfiability Analysis

The satisfiability problem is to determine, given a set Σ of constraints, whether there exists a *nonempty* database that satisfies Σ .

The satisfiability analysis of conditional dependencies is not only of theoretical interest, but is also important in practice. Indeed, when CFD^ps and CIND^ps are used as data quality rules, this analysis helps one check whether the rules make sense themselves. The need for this is particularly evident when the rules are manually designed or discovered from various datasets [10, 18, 15].

The satisfiability analysis of CFD^ps. Given any FDs, one does not need to worry about their satisfiability since any set of FDs is always satisfiable. However, as observed in [14], for a set Σ of CFDs on a relational schema R , there may not exist a *nonempty* instance I of R such that $I \models \Sigma$. As CFDs are a special case of CFD^ps, the same problem exists when it comes to CFD^ps.

Example 6. Consider CFD^p $\varphi = (R : A \rightarrow B, T_p)$ such that $T_p = \{(- \parallel a), (- \parallel \neq a)\}$. Then there exists no *nonempty* instance I of R that satisfies φ . Indeed, for any tuple t of R , φ requires that both $t[B] = a$ and $t[B] \neq a$. \square

This problem is already NP-complete for CFDs [14]. Below we show that it has the same complexity for CFD^ps despite their increased expressive power.

Proposition 1. *The satisfiability problem for CFD^ps is NP-complete.* \square

Proof sketch: The lower bound follows from the NP-hardness of their CFDs counterparts [14], since CFDs are a special case of CFD^ps. The upper bound is verified by presenting an NP algorithm that, given a set Σ of CFD^ps defined on a relation schema R , determines whether Σ is satisfiable. \square

It is known [14] that the satisfiability problem for CFDs is in PTIME when the CFDs considered are defined over attributes that have an infinite domain, *i.e.*, in the absence of finite domain attributes. However, this is no longer the case for CFD^ps. This tells us that the increased expressive power of CFD^ps does take a toll in this special case. It should be remarked that while the proof of Proposition 1 is an extension of its counterpart in [14], the result below is new.

Theorem 2. *In the absence of finite domain attributes, the satisfiability problem for CFD^ps remains NP-complete.* \square

Proof sketch: The problem is in NP by Proposition 1. Its NP-hardness is shown by reduction from the 3SAT problem, which is NP-complete (cf. [17]). \square

The satisfiability analysis of CIND^ps. Like FDs, one can specify arbitrary INDs or CINDs without worrying about their satisfiability. Below we show that CIND^ps also have this property, by extending the proof of its counterpart in [8].

8 Wenguang Chen, Wenfei Fan, and Shuai Ma

Proposition 3. *Any set Σ of CIND^Ps is always satisfiable.* \square

Proof sketch: Given a set Σ of CIND^Ps over a database schema \mathcal{R} , one can always construct a *nonempty* instance D of \mathcal{R} such that $D \models \Sigma$. \square

The satisfiability analysis of CFD^Ps and CIND^Ps. The satisfiability problem for CFDs and CINDs taken together is undecidable [8]. Since CFD^Ps and CIND^Ps subsume CFDs and CINDs, respectively, from these we immediately have:

Corollary 4. *The satisfiability problem for CFD^Ps and CIND^Ps is undecidable.* \square

4.2 The Implication Analysis

The implication problem is to determine, given a set Σ of dependencies and another dependency ϕ , whether or not Σ entails ϕ , denoted by $\Sigma \models \phi$. That is, whether or not for all databases D , if $D \models \Sigma$ then $D \models \phi$.

The implication analysis helps us remove redundant data quality rules, and thus improve the performance of error detection and repairing based on the rules.

Example 7. The CFD^Ps of Fig. 2 imply CFD^Ps $\varphi = \text{item}(\text{sale}, \text{price} \rightarrow \text{shipping}, T)$, where T consists of a single pattern tuple ($\text{sale} = \text{'F'}$, $\text{price} = 30 \parallel \text{shipping} = 6$). Thus in the presence of the CFD^Ps of Fig. 2, φ is redundant. \square

The implication analysis of CFD^Ps. We first show that the implication problem for CFD^Ps retains the same complexity as their CFDs counterpart. The result below is verified by extending the proof of its counterpart in [14].

Proposition 5. *The implication problem for CFD^Ps is coNP-complete.* \square

Proof sketch: The lower bound follows from the coNP-hardness of their CFDs counterpart [14], since CFDs are a special case of CFD^Ps. The coNP upper bound is verified by presenting an NP algorithm for its complement problem, *i.e.*, the problem for determining whether $\Sigma \not\models \varphi$. \square

Similar to the satisfiability analysis, it is known [14] that the implication analysis of CFDs is in PTIME when the CFDs are defined only with attributes that have an infinite domain. Analogous to Theorem 2, the result below shows that this is no longer the case for CFD^Ps, which does not find a counterpart in [14].

Theorem 6. *In the absence of finite domain attributes, the implication problem for CFD^Ps remains coNP-complete.* \square

Proof sketch: It is in coNP by Proposition 5. The coNP-hardness is shown by reduction from the 3SAT problem to its complement problem, *i.e.*, the problem for determining whether $\Sigma \not\models \varphi$. \square

The implication analysis of CIND^Ps. We next show that CIND^Ps do not make their implication analysis harder. This is verified by extending the proof of their CINDs counterpart given in [8].

Proposition 7. *The implication problem for CIND^Ps is EXPTIME-complete.* \square

Σ	General setting		Infinite domain only	
	Satisfiability	Implication	Satisfiability	Implication
CFDs [14]	NP-complete	coNP-complete	PTIME	PTIME
CFD ^p s	NP-complete	coNP-complete	NP-complete	coNP-complete
CINDs [8]	$O(1)$	EXPTIME-complete	$O(1)$	PSPACE-complete
CIND ^p s	$O(1)$	EXPTIME-complete	$O(1)$	EXPTIME-complete
CFDs + CINDs [8]	undecidable	undecidable	undecidable	undecidable
CFD ^p s + CIND ^p s	undecidable	undecidable	undecidable	undecidable

Table 1. Summary of complexity results

Proof sketch: The implication problem for CINDs is EXPTIME-hard [8]. The lower bound carries over to CIND^ps since CIND^ps subsume CINDs. The EXPTIME upper bound is shown by presenting an EXPTIME algorithm that, given a set $\Sigma \cup \{\psi\}$ of CIND^ps over a database schema \mathcal{R} , determines whether $\Sigma \models \psi$. \square

It is known [8] that the implication problem is PSPACE-complete for CINDs defined with infinite-domain attributes. Similar to Theorem 6, below we present a new result showing that this no longer holds for CIND^ps.

Theorem 8. *In the absence of finite domain attributes, the implication problem for CIND^ps remains EXPTIME-complete.* \square

Proof sketch: The EXPTIME upper bound follows from Proposition 7. The EXPTIME-hardness is shown by reduction from the implication problem for CINDs in the general setting, in which finite-domain attributes may be present; the latter is known to be EXPTIME-complete [8]. \square

The implication analysis of CFD^ps and CIND^ps. When CFD^ps and CIND^ps are taken together, their implication analysis is beyond reach in practice. This is not surprising since the implication problem for FDs and INDs is already undecidable [1]. Since CFD^ps and CIND^ps subsume FDs and INDs, respectively, from the undecidability result for FDs and INDs, the corollary below follows immediately.

Corollary 9. *The implication problem for CFD^ps and CIND^ps is undecidable.* \square

Summary. The complexity bounds for reasoning about CFD^ps and CIND^ps are summarized in Table 1. To give a complete picture we also include in Table 1 the complexity bounds for the static analyses of CFDs and CINDs, taken from [14, 8]. The results shown in Table 1 tell us the following.

(a) Despite the increased expressive power, CFD^ps and CIND^ps do not complicate the static analyses: the satisfiability and implication problems for CFD^ps and CIND^ps have the same complexity bounds as their counterparts for CFDs and CINDs, taken separately or together.

(b) In the special case when CFD^ps and CIND^ps are defined with infinite-domain attributes only, however, the static analyses of CFD^ps and CIND^ps do not get simpler, as opposed to their counterparts for CFDs and CINDs. That is, in this special case the increased expressive power of CFD^ps and CIND^ps comes at a price.

10 Wenguang Chen, Wenfei Fan, and Shuai Ma

5 Validation of CFD^ps and CIND^ps

If CFD^ps and CIND^ps are to be used as data quality rules, the first question we have to settle is how to effectively detect errors and inconsistencies as violations of these dependencies, by leveraging functionality supported by commercial DBMS. More specifically, consider a database schema $\mathcal{R} = (R_1, \dots, R_n)$, where R_i is a relation schema for $i \in [1, n]$. The error detection problem is stated as follows.

The *error detection problem* is to find, given a set Σ of CFD^ps and CIND^ps defined on \mathcal{R} , and a database instance $D = (I_1, \dots, I_n)$ of \mathcal{R} as input, the subset (I'_1, \dots, I'_n) of D such that for each $i \in [1, n]$, $I'_i \subseteq I_i$ and each tuple in I'_i violates at least one CFD^p or CIND^p in Σ . We denote the set as $\text{vio}(D, \Sigma)$, referred to it as the *violation set* of D w.r.t. Σ .

In this section we develop SQL-based techniques for error detection based on CFD^ps and CIND^ps. The main result of the section is as follows.

Theorem 10. *Given a set Σ of CFD^ps and CIND^ps defined on \mathcal{R} and a database instance D of \mathcal{R} , where $\mathcal{R} = (R_1, \dots, R_n)$, a set of SQL queries can be automatically generated such that (a) the collection of the answers to the SQL queries in D is $\text{vio}(D, \Sigma)$, (b) the number and size of the set of SQL queries depend only on the number n of relations and their arities in \mathcal{R} , regardless of Σ . \square*

We next present the main techniques for the query generation method. Let $\Sigma_{\text{cfd}^p}^i$ be the set of all CFD^ps in Σ defined on the same relation schema R_i , and $\Sigma_{\text{cind}^p}^{(i,j)}$ the set of all CIND^ps in Σ from R_i to R_j , for $i, j \in [1, n]$. We show the following. (a) The violation set $\text{vio}(D, \Sigma_{\text{cfd}^p}^i)$ can be computed by *two* SQL queries. (b) Similarly, $\text{vio}(D, \Sigma_{\text{cind}^p}^{(i,j)})$ can be computed by a *single* SQL query. (c) These SQL queries encode pattern tableaux of CFD^ps (CIND^ps) with data tables, and hence their sizes are independent of Σ . From these Theorem 10 follows immediately.

5.1 Encoding CFD^ps and CIND^ps with Data Tables

We first show the following, by extending the encoding of [14, 7]. (a) The pattern tableaux of all CFD^ps in $\Sigma_{\text{cfd}^p}^i$ can be encoded with *three data tables*, and (b) the pattern tableaux of all CIND^ps in $\Sigma_{\text{cind}^p}^{(i,j)}$ can be represented as *four data tables*, no matter how many dependencies are in the sets and how large they are.

Encoding CFD^ps. We encode all pattern tableaux in $\Sigma_{\text{cfd}^p}^i$ with three tables enc_L , enc_R and enc_{\neq} , where enc_L (resp. enc_R) encodes the non-negation ($=, <, \leq, >, \geq$) patterns in LHS (resp. RHS), and enc_{\neq} encodes those negation (\neq) patterns. More specifically, we associate a unique id cid with each CFD^ps in $\Sigma_{\text{cfd}^p}^i$, and let enc_L consist of the following attributes: (a) cid , (b) each attribute A appearing in the LHS of some CFD^ps in $\Sigma_{\text{cfd}^p}^i$, and (b) its four companion attributes $A_>, A_{\geq}, A_<, \text{ and } A_{\leq}$. That is, for each attribute, there are five columns in enc_L , one for each non-negation operator. Similarly, enc_R is defined. We use an enc_{\neq} tuple to encode a pattern $A \neq c$ in a CFD^p, consisting of cid , att , pos , and val , encoding the CFD^p id, the attribute A , the position ('LHS' or 'RHS'),

(1) enc_L					(2) enc_R					(3) enc_{\neq}			
cid	sale	price	price _{>}	price _≤	cid	shipping	price	price _≥	price _{<}	cid	pos	att	val
2	T	null	null	null	2	0	null	null	null				
3	F	-	20	40	3	6	null	null	null				
4	T	null	null	null	4	null	-	2.99	9.99				

Fig. 4. Encoding example of CFD^p s

and the constant c , respectively. Note that the arity of enc_L (enc_R) is bounded by $5 * |R_i| + 1$, where $|R_i|$ is the arity of R_i , and the arity of enc_{\neq} is 4.

Before we populate these tables, let us first describe a preferred form of CFD^p s that would simplify the analysis to be given. Consider a $\text{CFD}^p \varphi = R(X \rightarrow Y, T_p)$. If φ is not satisfiable we can simply drop it from Σ . Otherwise it is equivalent to a $\text{CFD}^p \varphi' = R(X \rightarrow Y, T'_p)$ such that for any pattern tuples t_p, t'_p in T'_p and for any attribute A in $X \cup Y$, (a) if $t_p[A]$ is $\text{op } a$ and $t'_p[A]$ is $\text{op } b$, where op is not \neq , then $a = b$, (b) if $t_p[A]$ is '-' then so is $t'_p[A]$. That is, for each non-negation op (resp. '-'), there is a *unique* constant a such that $t_p[A] = \text{'op } a'$ (resp. $t_p[A] = \text{'-'}$) is the only op (resp. '-') pattern appearing in the A column of T'_p . We refer to $t_p[A]$ as $T'_p(\text{op}, A)$ (resp. $T'_p(\text{'-'}, A)$), and consider *w.l.o.g.* CFD^p s of this form only. Note that there are possibly multiple $t_p[A] \neq c$ patterns in T'_p .

We populate enc_L , enc_R and enc_{\neq} as follows. For each $\text{CFD}^p \varphi = R(X \rightarrow Y, T_p)$ in $\Sigma_{\text{cfid}^p}^i$, we generate a distinct cid id_{φ} for it, and do the following.

- Add a tuple t_1 to enc_L such that (a) $t[\text{cid}] = \text{id}_{\varphi}$; (b) for each $A \in X$, $t[A] = \text{'-'}$ if $T'_p(\text{'-'}, A)$ is '-' , and for each non-negation predicate op , $t[A_{\text{op}}] = \text{'a'}$ if $T'_p(\text{op}, A)$ is $\text{'op } a'$; (c) we let $t[B] = \text{'null'}$ for all other attributes B in enc_L .
- Similarly add a tuple t_2 to enc_R for attributes in Y .
- For each attribute $A \in X \cup Y$ and each $\neq a$ pattern in $T_p[A]$, add a tuple t to enc_{\neq} such that $t[\text{cid}] = \text{id}_{\varphi}$, $t[\text{att}] = \text{'A'}$, $t[\text{val}] = \text{'a'}$, and $t[\text{pos}] = \text{'LHS'}$ (resp. $t[\text{pos}] = \text{'RHS'}$) if attribute A appears in X (resp. Y).

Example 8. Recall from Fig. 2 CFD^p s φ_2 , φ_3 and φ_4 defined on relation *item*. The three CFD^p s are encoded with tables shown in Fig. 4: (a) enc_L consists of attributes: *cid*, *sale*, *price*, *price_>* and *price_≤*; (b) enc_R consists of *cid*, *shipping*, *price*, *price_≥* and *price_<*; those attributes in a table with only 'null' pattern values do not contribute to error detection, and are thus omitted; (c) enc_{\neq} is empty since all these CFD^p s have no negation patterns. One can easily reconstruct these CFD^p s from tables enc_L , enc_R and enc_{\neq} by collating tuples based on *cid*. \square

Encoding CIND^p s. All CIND^p s in $\Sigma_{\text{cind}^p}^{(i,j)}$ can be encoded with four tables enc , enc_L , enc_R and enc_{\neq} . Here enc_L (resp. enc_R) and enc_{\neq} encode non-negation patterns on relation R_i (resp. R_j) and negation patterns on relations R_i or R_j , respectively, along the same lines as their counterparts for CFD^p s. We use enc to encode the *INDs embedded* in CIND^p s, which consists of the following attributes: (1) *cid* representing the id of a CIND^p , and (2) those X attributes of R_i and Y attributes of R_j appearing in some CIND^p s in $\Sigma_{\text{cind}^p}^{(i,j)}$. Note that the number of attributes in enc is bounded by $|R_i| + |R_j| + 1$, where $|R_i|$ is the arity of R_i .

12 Wenguang Chen, Wenfei Fan, and Shuai Ma

(1) enc			(2) enc _L			(3) enc _R		(4) enc _≠			
cid	state _L	state _R	cid	type	state	cid	rate	cid	pos	att	val
1	1	1	1	-	null	1	null	1	LHS	type	art
2	1	1	2	-	DL	2	0	2	LHS	type	art

Fig. 5. Encoding example of CIND^ps

For each CIND^p $\psi = (R_i[A_1 \dots A_m; X_p] \subseteq R_j[B_1 \dots B_m; Y_p], T_p)$ in $\Sigma_{\text{cind}^p}^{(i,j)}$, we generate a distinct cid id_ψ for it, and do the following.

- Add tuples t_1 and t_2 to enc_L and enc_R based on attributes X_p and Y_p , respectively, along the same lines as their CFD^p counterpart.
- Add tuples to enc_\neq in the same way as their CFD^p counterparts.
- Add tuple t to enc such that $t[\text{cid}] = \text{id}_\psi$. For each $k \in [1, m]$, let $t[A_k] = t[B_k] = k$, and $t[A] = \text{'null'}$ for the rest attributes A of enc .

Example 9. Figure 4 shows the coding of CIND^ps ψ_1 and ψ_2 given in Fig. 3. We use state_L and state_R in enc to denote the occurrences of attribute state in item and tax , respectively. In tables enc_L and enc_R , attributes with only ‘null’ patterns are omitted, for the same reason as for CFD^ps mentioned above. \square

Putting these together, it is easy to verify that at most $O(n^2)$ data tables are needed to encode dependencies in Σ , regardless of the size of Σ . Recall that n is the number of relations in database \mathcal{R} .

5.2 SQL-based Detection Methods

We next show how to generate SQL queries based on the encoding above. For each $i \in [1, n]$, we generate *two* SQL queries that, when evaluated on the I_i table of D , find $\text{vio}(D, \Sigma_{\text{cfd}^p}^i)$. Similarly, for each $i, j \in [1, n]$, we generate a *single* SQL query $Q_{(i,j)}$ that, when evaluated on (I_i, I_j) of D , returns $\text{vio}(D, \Sigma_{\text{cind}^p}^{(i,j)})$. Putting these query answers together, we get $\text{vio}(D, \Sigma)$, the violation set of D w.r.t. Σ .

Below we show how the SQL query $Q_{(i,j)}$ is generated for validating CIND^ps in $\Sigma_{\text{cind}^p}^{(i,j)}$, which has not been studied by previous work. For the lack of space we omit the generation of detection queries for CFD^ps, which is an extension of the SQL techniques for CFDs discussed in [14, 7].

The query $Q_{(i,j)}$ for the validation of $\Sigma_{\text{cind}^p}^{(i,j)}$ is given as follows, which capitalizes on the data tables enc , enc_L , enc_R and enc_\neq that encode CIND^ps in $\Sigma_{\text{cind}^p}^{(i,j)}$.

```

select Ri.*
from Ri, encL L, enc≠ N
where Ri.X ≍ L and Ri.X ≍ N and not exists (
  select Rj.*
  from Rj, enc H, encR R, enc≠ N
  where Ri.X = Rj.Y and L.cid = R.cid and L.cid = H.cid and
    Rj.Y ≍ R and Rj.Y ≍ N)

```

Here (1) $X = \{A_1, \dots, A_{m1}\}$ and $Y = \{B_1, \dots, B_{m2}\}$ are the sets of attributes of R_i and R_j appearing in $\Sigma_{\text{cind}^p}^{(i,j)}$, respectively; (2) $R_i.X \asymp L$ is the conjunction of

$L.A_k$ is null or $R_i.A_k = L.A_k$ or ($L.A_k = \text{'_'} \wedge$
 and ($L.A_{i>} \text{ is null or } R_i.A_k > L.A_{i>}) \wedge (L.A_{i\geq} \text{ is null or } R_i.A_k \geq L.A_{i\geq})$
 and ($L.A_{k<} \text{ is null or } R_i.A_k < L.A_{k<}) \wedge (L.A_{i\leq} \text{ is null or } R_i.A_k \leq L.A_{i\leq})$)

for $k \in [1, m_1]$; (3) $R_j.Y \asymp R$ is defined similarly for attributes in Y ; (4) $R_i.X \asymp N$ is a shorthand for the conjunction below, for $k \in [1, m_1]$:

not exists (select * from N where $L.cid = N.cid$ and $N.pos = \text{'LHS'}$ and
 $N.att = \text{'A}_k$ and $R_i.A_k = N.val$);

(5) $R_j.Y \asymp N$ is defined similarly, but with $N.pos = \text{'RHS'}$; (6) $R_i.X = R_j.Y$ represents the following: for each A_k ($k \in [1, m_1]$) and each B_l ($l \in [1, m_2]$),
 ($H.A_k \text{ is null or } H.B_l \text{ is null or } H.B_l \neq H.A_k \text{ or } R_i.A_k = R_j.B_l$).

Intuitively, (1) $R_i.X \asymp L$ and $R_i.X \asymp N$ ensure that the R_i tuples selected match the LHS patterns of some $CIND^p$ s in $\Sigma_{cind^p}^{(i,j)}$; (2) $R_j.Y \asymp R$ and $R_j.Y \asymp N$ check the corresponding RHS patterns of these $CIND^p$ s on R_j tuples; (3) $R_i.X = R_j.Y$ enforces the *embedded* INDS; (4) $L.cid = R.cid$ and $L.cid = H.cid$ assure that the LHS and RHS patterns in the same $CIND^p$ are correctly collated; and (5) not exists in Q ensures that the R_i tuples selected violate $CIND^p$ s in $\Sigma_{cind^p}^{(i,j)}$.

Example 10. Using the coding of Fig. 5, an SQL query Q for checking $CIND^p$ s ψ_1 and ψ_2 of Fig. 3 is given as follows:

```
select R1.* from item R1, encL L, enc≠ N
where (L.type is null or R1.type = L.type or L.type = \_') and not exist (
  select * from N
  where N.cid = L.cid and N.pos = \text{'LHS'} and N.att = \text{'type'}
  and (L.state is null or R1.state = L.state or L.state = \_') and not exist (
    select * from N
    where N.cid = L.cid and N.pos = \text{'LHS'} and N.att = \text{'state'} and R1.state = N.val)
  and not exists (
    select R2.* from tax R2, encH H, encR R
    where (H.stateL is null or H.stateR is null or H.stateL! = H.stateR or
      R2.state = R1.state) and L.cid = H.cid and L.cid = R.cid and
      (R.rate is null or R2.rate = R.rate or R.rate = \_') and not exist (
        select * from N
        where N.cid = R.cid and N.pos = \text{'RHS'} and N.att = \text{'rate'} and R2.rate = N.val))
```

The SQL queries generated for error detection can be simplified as follows. As shown in Example 10, when checking patterns imposed by enc , enc_L or enc_R , the queries need not consider attributes A if $t[A]$ is 'null' for each tuple t in the table. Similarly, if an attribute A does not appear in any tuple in enc_{\neq} , the queries need not check A either. From this, it follows that we do not even need to generate those attributes with only 'null' patterns for data tables enc , enc_L or enc_R when encoding $CIND^p$ s or CFD^p s. \square

6 Related Work

Constraint-based data cleaning was introduced in [2], which proposed to use dependencies, *e.g.*, FDs, INDS and denial constraints, to detect and repair errors

14 Wenguang Chen, Wenfei Fan, and Shuai Ma

in real-life data (see, *e.g.*, [11] for a comprehensive survey). As an extension of traditional FDs, CFDs were developed in [14], for improving the quality of data. It was shown in [14] that the satisfiability and implication problems for CFDs are NP-complete and coNP-complete, respectively. Along the same lines, CINDs were proposed in [8] to extend INDs. It was shown [8] that the satisfiability and implication problems for CINDs are in constant time and EXPTIME-complete, respectively. SQL techniques were developed in [14] to detect errors by using CFDs, but have not been studied for CINDs. This work extends the static analyses of conditional dependencies of [14, 8], and has established several new complexity results, notably in the absence of finite-domain attributes (*e.g.*, Theorems 2, 6, 8). In addition, it is the first work to develop SQL-based techniques for checking violations of CINDs and violations of CFD^ps and CIND^ps taken together.

Extensions of CFDs have been proposed to support disjunction and negation [7], cardinality constraints and synonym rules [9], and to specify patterns in terms of value ranges [18]. While CFD^ps are more powerful than the extension of [18], they cannot express disjunctions [7], cardinality constraints and synonym rules [9]. To our knowledge no extensions of CINDs have been studied. This work is the first full treatment of extensions of CFDs and CINDs by incorporating built-in predicates ($\neq, <, \leq, >, \geq$), from static analyses to error detection.

Methods have been developed for discovering CFDs [10, 18, 15] and for repairing data based on either CFDs [13], traditional FDs and INDs taken together [5], denial constraints [4, 12], or aggregate constraints [16]. We defer the treatment of these topics for CFD^ps and CIND^ps to future work.

A variety of extensions of FDs and INDs have been studied for specifying constraint databases and constraint logic programs [3, 6, 19, 20]. While the languages of [3, 19] cannot express CFDs, constraint-generating dependencies (CGDs) of [3] and constrained tuple-generating dependencies (CTGDs) of [20] can express CFD^ps, and CTGDs can also express CIND^ps. The increased expressive power of CTGDs comes at the price of a higher complexity: both their satisfiability and implication problems are undecidable. Built-in predicates and arbitrary constraints are supported by CGDs, for which it is not clear whether effective SQL queries can be developed to detect errors. It is worth mentioning that Theorems 2 and 6 of this work provide lower bounds for the consistency and implication analyses of CGDs, by using patterns with built-in predicates only.

7 Conclusions

We have proposed CFD^ps and CIND^ps, which further extend CFDs and CINDs, respectively, by allowing patterns on data values to be expressed in terms of $\neq, <, \leq, >$ and \geq predicates. We have shown that CFD^ps and CIND^ps are more powerful than CFDs and CINDs for detecting errors in real-life data. In addition, the satisfiability and implication problems for CFD^ps and CIND^ps have the same complexity bounds as their counterparts for CFDs and CINDs, respectively. We have also provided automated methods to generate SQL queries for detecting errors based on CFD^ps and CIND^ps. These provide commercial DBMS with an immediate capability to capture errors commonly found in real-world data.

One topic for future work is to develop a dependency language that is capable of expressing various extensions of CFDs (*e.g.*, CFD^p s, eCFD s [7] and CFD^c s [9]), without increasing the complexity of static analyses. Second, we are developing effective algorithms for discovering CFD^p s and CIND^p s, along the same lines as [10, 18, 15]. Third, we plan to extend the methods of [5, 13] to repair data based on CFD^p s and CIND^p s, instead of using CFDs [13], traditional FDs and INs [5], denial constraints [4, 12], and aggregate constraints [16].

Acknowledgments. Fan and Ma are supported in part by EPSRC E029213/1. Fan is a Yangtze River Scholar at Harbin Institute of Technology.

References

1. S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
2. M. Arenas, L. E. Bertossi, and J. Chomicki. Consistent query answers in inconsistent databases. In *PODS*, 1999.
3. M. Baudinet, J. Chomicki, and P. Wolper. Constraint-Generating Dependencies. *J. Comput. Syst. Sci.*, 59(1):94–115, 1999.
4. L. E. Bertossi, L. Bravo, E. Franconi, and A. Lopatenko. The complexity and approximation of fixing numerical attributes in databases under integrity constraints. *Inf. Syst.*, 33(4-5):407–434, 2008.
5. P. Bohannon, W. Fan, M. Flaster, and R. Rastogi. A cost-based model and effective heuristic for repairing constraints by value modification. In *SIGMOD*, 2005.
6. P. D. Bra and J. Paredaens. Conditional dependencies for horizontal decompositions. In *ICALP*, 1983.
7. L. Bravo, W. Fan, F. Geerts, and S. Ma. Increasing the expressivity of conditional functional dependencies without extra complexity. In *ICDE*, 2008.
8. L. Bravo, W. Fan, and S. Ma. Extending dependencies with conditions. In *VLDB*, 2007.
9. W. Chen, W. Fan, and S. Ma. Incorporating cardinality constraints and synonym rules into conditional functional dependencies. *IPL*, 109(14):783–789, 2009.
10. F. Chiang and R. J. Miller. Discovering data quality rules. In *VLDB*, 2008.
11. J. Chomicki. Consistent query answering: Five easy pieces. In *ICDT*, 2007.
12. J. Chomicki and J. Marcinkowski. Minimal-change integrity maintenance using tuple deletions. *Inf. Comput.*, 197(1-2):90–121, 2005.
13. G. Cong, W. Fan, F. Geerts, X. Jia, and S. Ma. Improving data quality: Consistency and accuracy. In *VLDB*, 2007.
14. W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies. *TODS*, 33(2), 2008.
15. W. Fan, F. Geerts, L. V. Lakshmanan, and M. Xiong. Discovering conditional functional dependencies. In *ICDE*, 2009.
16. S. Flesca, F. Furfaro, and F. Parisi. Consistent query answers on numerical databases under aggregate constraints. In *DBPL*, 2005.
17. M. Garey and D. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Company, 1979.
18. L. Golab, H. J. Karloff, F. Korn, D. Srivastava, and B. Yu. On generating near-optimal tableaux for conditional functional dependencies. In *VLDB*, 2008.
19. M. J. Maher. Constrained dependencies. *TCS*, 173(1):113–149, 1997.
20. M. J. Maher and D. Srivastava. Chasing Constrained Tuple-Generating Dependencies. In *PODS*, 1996.
21. R. van der Meyden. The complexity of querying indefinite data about linearly ordered domains. *JCSS*, 54(1), 1997.