Analyses and Validation of Conditional Dependencies with Built-in Predicates

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

- ¹ Peking University, China
- University of Edinburgh, UK
 Bell Laboratories, USA

Abstract. This paper proposes a natural extension of conditional functional dependencies (CFDs [11]) and conditional inclusion dependencies (CINDS [7]), denoted by CFD^ps and CIND^ps, respectively, by specifying patterns of data values with \neq , <, \leq , > and \geq predicates. As data quality rules, CFD^ps and CIND^ps 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^ps and CIND^ps. (a) One concerns the satisfiability and implication problems for CFD^ps and CIND^ps, 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^ps and CIND^ps. We show that given a set Σ of CFD^ps 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^ps and CIND^ps.

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 [11]) and conditional inclusion dependencies (CINDs [7]), 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 [11, 7].

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)

	id	name	type	price	shipping	sale	state	
t_1 :	b1	Harry Potter	book	25.99	0	Τ	WA	
t_2 :	c1	Snow White	CD	9.99	2	F	NY	
t_3 :	b2	Catch-22	book	34.99	20	F	DL	
t_4 :	a1	Sunflowers	art	$5\mathrm{m}$	500	F	DL	
(a) An item relation								

	state	rate		
t_5 :	PA	6		
t_6 :	NY	4		
t_7 :	DL	0		
t_8 :	NJ	3.5		
(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 [11, 7]) on the data include the following:

```
\mathsf{cfd}_1: item (id \to name, type, price, shipping, sale) \mathsf{cfd}_2: tax (state \to rate) \mathsf{cfd}_3: item (sale = 'T' \to 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 $t[\mathsf{sale}] = \mathsf{`T'}$ then $t[\mathsf{shipping}]$ 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, $\mathsf{sale} = \mathsf{`T'}$ and $\mathsf{shipping} = 0$. It is to hold only on item tuples that match the pattern $\mathsf{sale} = \mathsf{`T'}$. 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 [11] or CINDs of [7], but can be expressed as the following dependencies:

```
\begin{array}{lll} \mathsf{pfd}_1 \colon \mathsf{item} \; (\mathsf{sale} = `F' \; \& \; \mathsf{price} \leq 20 \to \mathsf{shipping} = 3) \\ \mathsf{pfd}_2 \colon \mathsf{item} \; (\mathsf{sale} = `F' \; \& \; \mathsf{price} > 20 \; \& \; \mathsf{price} \leq 40 \to \mathsf{shipping} = 6) \\ \mathsf{pfd}_3 \colon \mathsf{item} \; (\mathsf{sale} = `F' \; \& \; \mathsf{price} > 40 \to \mathsf{shipping} = 10) \\ \mathsf{pfd}_4 \colon \mathsf{item} \; (\mathsf{sale} = `T' \to \mathsf{price} \geq 2.99 \; \& \; \mathsf{price} < 9.99) \\ \mathsf{pind}_1 \colon \mathsf{item} \; (\mathsf{state}; \; \mathsf{type} \neq `\mathsf{art}') \subseteq \mathsf{tax} \; (\mathsf{state}; \; \mathsf{nil}) \end{array}
```

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 [11] by specifying patterns of semantically related data values in terms of predicates $<, \le, >$, and \ge . 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 [7] by specifying patterns with \ne : for any item

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

Using $\mathsf{pfd}_1\mathsf{-}\mathsf{pfd}_4$ 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 .

None of $\mathsf{pfd}_1\mathsf{-pfd}_4$ and pind_1 can be expressed as FDs or INDs [1], which do not allows constants, or as CFDs [11] or CINDs [7], which specify patterns with equality (=) only. While there have been extensions of CFDs [6, 15], 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^ps and CIND^ps, 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^ps or CIND^ps. These dependencies are capable of capturing errors in realworld 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., [19]), 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.

$$(1) \ \varphi_1 = \mathsf{tax} \ (\mathsf{state} \to \mathsf{rate}, \ T_1) \qquad \qquad (2) \ \varphi_2 = \mathsf{item} \ (\mathsf{sale} \to \mathsf{shipping}, \ T_2) \\ T_1: \ \ \frac{\mathsf{sale} \ | \mathsf{shipping} \ }{= \ T} = 0$$

Fig. 2. Example CFD^p s

Given a set Σ of CFD^ps and CIND^ps 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 CINDs, not to mention CFD^ps and CIND^ps taken together. These provide the capability of detecting errors in a single relation (CFD^ps) and across different relations (CIND^ps) 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. All proofs are in [13].

2 Incorporating Built-in Predicates into CFDs

We now define CFD^ps, 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 $\mathsf{attr}(R)$. For each attribute $A \in \mathsf{attr}(R)$, its domain is specified in R, denoted as $\mathsf{dom}(A)$, which is either finite $(e.g., \mathsf{bool})$ or infinite $(e.g., \mathsf{string})$. We assume w.l.o.g. that a domain is totally ordered if $<, \le, >$ or \ge is defined on it.

Syntax. A CFD^p φ on R is a pair $R(X \to Y, T_p)$, where (1) X, Y are sets of attributes in $\mathsf{attr}(R)$; (2) $X \to Y$ is a standard FD, referred to as the FD embedded $in \varphi$; 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 $\mathsf{dom}(A)$, or 'op a', where op is one of $=, \neq, <, \leq, >, \geq$, and 'a' is a constant in $\mathsf{dom}(A)$.

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 ' \parallel '. We write φ as $(X \to Y, T_p)$ when R is clear from the context, and denote X as $\mathsf{LHS}(\varphi)$ and Y as $\mathsf{RHS}(\varphi)$.

Example 2. The dependencies cfd1-cfd3 and pfd1-pfd4 that we have seen in Example 1 can all be expressed as CFD^ps. Figure 2 shows some of these CFD^ps: φ_1 (for FD cfd₂), φ_2 (for CFD cfd₃), φ_3 (for pfd₂), and φ_4 (for pfd₄).

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

A data tuple t of R is said to $match \ \mathsf{LHS}(\varphi)$, 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 $\mathsf{dom}(A)$), or (b) t[A] op a if $t_{pi}[A]$ is 'op a', where the operator $\mathsf{op} \ (=, \neq, <, \leq, > \ \mathsf{or} \ge)$ is interpreted by its standard semantics. Similarly, the notion that t $matches \ \mathsf{RHS}(\varphi)$ is defined, denoted by $t[Y] \asymp T_p[Y]$.

Intuitively, each pattern tuple t_{pi} specifies a condition via $t_{pi}[X]$, and $t[X] \approx T_p[X]$ if t[X] satisfies the *conjunction* of all these conditions. Similarly, $t[Y] \approx T_p[Y]$ 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 $t_1[X] = t_2[X] \times T_p[X]$, then $t_1[Y] = t_2[Y] \times T_p[Y]$. 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) $t_1[Y] \times T_p[Y]$, which assure the binding between the *constants* in $t_1[Y]$ and the *constants* in $t_{pi}[Y]$ for all t_{pi} in T_p .

An instance I of R satisfies a set Σ of CFD^p s, denoted by $I \models \Sigma$, if $I \models \varphi$ for $each CFD^p \varphi$ 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 $t_3[\mathsf{sale}] = \text{`F'}$ and $20 < t_3[\mathsf{price}] \le 40$, but $t_3[\mathsf{shipping}]$ is 20 instead of 6. Note that t_3 matches $\mathsf{LHS}(\varphi_3)$ since it satisfies the condition specified by the *conjunction* of the pattern tuples in T_3 . Similarly, t_1 violates φ_4 , since $t_1[\mathsf{sale}] = \text{`T'}$ but $t_1[\mathsf{price}] > 9.99$. Observe that while it takes two tuples to violate a standard FD, a single tuple may violate a CFD^p .

Special cases. (1) A standard FD $X \to Y$ [1] can be expressed as a CFD $(X \to Y, T_p)$ in which T_p contains a single tuple consisting of '-' only, without constants. (2) A CFD $(X \to Y, T_p)$ [11] with $T_p = \{t_{p1}, \ldots, t_{pk}\}$ can be expressed as a set $\{\varphi_1, \ldots, \varphi_k\}$ of CFD^ps such that for $i \in [1, k], \varphi_i = (X \to Y, T_{pi})$, 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 cfd2 and CFD cfd3 in Example 1, respectively. Note that all data quality rules in [8, 15] can be expressed as CFD^ps.

3 Incorporating Built-in Predicates into CINDs

Along the same lines as CFD^p s, we next define $CIND^p$ s, 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 $\mathsf{attr}(R_1)$ and $\mathsf{attr}(R_2)$, 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

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

Fig. 3. Example $CIND^p$ s

unnamed variable '_' that draws values from dom(A), or 'op a', where op is one of $=, \neq, <, \leq, >, \geq$ and 'a' is a constant in dom(A).

We denote $X \cup X_p$ as $\mathsf{LHS}(\psi)$ and $Y \cup Y_p$ as $\mathsf{RHS}(\psi)$, and separate the X_p

and Y_p attributes in a pattern tuple with '||'. We use nil to denote an *empty* list. *Example 4.* Figure 3 shows two example CIND^ps: ψ expresses pind₁ of Example 1, and ψ_2 refines ψ 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.

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

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 LHS(ψ_1) since t_1 [type] \neq 'art', but there is no tuple t in tax such that t[state] = t_1 [state] = 'WA'. In contrast, D_0 satisfies ψ_2 .

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 $t'_p[\text{price}] \ge 19.99$.

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 CFD^ps 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 $T_p = \{t_{p1}, \ldots, t_{pk}\}$ can be expressed as a set $\{\psi_1, \ldots, \psi_k\}$ 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^p s and $CIND^p$ s

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^p s and $CIND^p$ s are used as data quality rules, one wants to know whether the rules make sense or are dirty themselves. The need for this is particularly evident when the rules are manually designed or discovered from various datasets [8, 15, 12].

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 [11], 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 $\varphi = (R : A \to B, T_p)$ 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 $t[B] \neq a$.

This problem is already NP-complete for CFDs [11]. Below we show that it has the same complexity for CFD p s despite their increased expressive power.

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

Proof sketch: The lower bound follows from the NP-hardness of their CFDs counterparts [11], 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. This requires the establishment of a small model property [13].

It is known [11] that the satisfiability problem for CFDs is in PTIME when the CFDs considered are defined over attributes that have an infinite domain. However, this is no longer the case for CFD^p s. This tells us that the increased expressive power of CFD^p s 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 [11], the result below and its proof are new [13].

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

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. [14]).

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 [7].

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

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 $D \models \Sigma$.

The satisfiability analysis of CFD^ps and $CIND^ps$. The satisfiability problem for CFDs and CINDs taken together is undecidable [7]. Since CFD^ps and CINDs subsume CFDs and CINDs, respectively, from these we immediately have:

Corollary 4. The satisfiability problem for CFD^ps and $CIND^ps$ is undecidable.

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

The implication analysis of CFD^p s. We first show that the implication problem for CFD^p s retains the same complexity as their CFDs counterpart. The result below is verified by extending the proof of its counterpart in [11].

Proposition 5. The implication problem for CFD^ps is conp-complete. \Box **Proof sketch:** The lower bound follows from the conp-hardness of their CFDs

counterpart [11], 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 $\Sigma \not\models \varphi$.

Similar to the satisfiability analysis, it is known [11] 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^p s, which does not find a counterpart in [11].

Theorem 6. In the absence of finite domain attributes, the implication problem for CFD^ps remains conp-complete.

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 $\Sigma \not\models \varphi$.

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 [7].

∇	Gene	eral setting	Infinite domain only			
	Satisfiability	Implication	Satisfiability	Implication		
CFDs [11]	NP-complete	conp-complete	PTIME	PTIME		
$\mathrm{CFD}^p\mathrm{s}$	NP-complete	conp-complete	NP-complete	conp-complete		
CINDs [7]	O(1)	EXPTIME-complete	O(1)	PSPACE-complete		
$CIND^p$ s	O(1)	EXPTIME-complete	O(1)	EXPTIME-complete		
CFDS + CINDS [7]	undecidable	undecidable	undecidable	undecidable		
$CFD^p s + CIND^p s$	undecidable	undecidable	undecidable	undecidable		

Table 1. Summary of complexity results

Proposition 7. The implication problem for CIND^ps is EXPTIME-complete. \square **Proof sketch:** The implication problem for CINDs is EXPTIME-hard [7]. 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 [7] 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 $CIND^p$ s.

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

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 [7].

The implication analysis of CFD^p s and $CIND^p$ s. When CFD^p s and $CIND^p$ s 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^p s and $CIND^p$ s 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^p s and $CIND^p$ s 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 [11, 7]. The results shown in Table 1 tell us the following.

- (a) Despite the increased expressive power, CFD^p s and $CIND^p$ s do not complicate the static analyses: 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, taken separately or together.
- (b) In the special case when CFD^p s and $CIND^p$ s are defined with infinite-domain attributes only, however, the static analyses of CFD^p s and $CIND^p$ s 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^p s and $CIND^p$ s comes at a price.

5 Validation of CFD^p s and $CIND^p$ s

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 $\mathcal{R} = (R_1, \ldots, R_n)$, 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 $D = (I_1, \ldots, I_n)$ of \mathcal{R} as input, the set of all the tuples in D that violate at least one CFD^p or CIND^p in Σ . We denote the set as vio (D, Σ) , 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^p s and $CIND^p$ s. 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 $\mathcal{R} = (R_1, \ldots, R_n)$, a set of SQL queries can be automatically generated such that (a) the collection of the answers to the SQL queries in D is $vio(D, \Sigma)$, (b) the number and size of the set of SQL queries depend only on the number n of relations in \mathcal{R} , regardless of the cardinality and size of $\Sigma.\square$

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

5.1 Encoding CFD^p s and $CIND^p$ s with Data Tables

We first show the following, by extending the encoding of [11, 6]. (a) The pattern tableaux of all CFD^ps in $\Sigma^{i}_{\mathsf{cfd}^p}$ can be encoded with three data tables, and (b) the pattern tableaux of all CIND^ps in $\Sigma^{(i,j)}_{\mathsf{cind}^p}$ 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 $\Sigma^i_{\mathsf{cfd}^p}$ 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^p s in $\Sigma^i_{\mathsf{cfd}^p}$, 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 $\Sigma^i_{\mathsf{cfd}^p}$, and (b) its four companion attributes $A_>$, $A_>$, $A_>$, $A_<$, and $A_<$. 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.

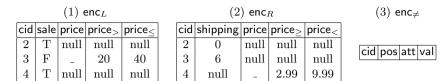


Fig. 4. Encoding example of CFD^p s

Before we populate these tables, let us first take a closer look at CFD^p $\varphi = R(X \to Y, T_p)$. If φ is not satisfiable we can simply drop it from Σ . Otherwise it is equivalent to a CFD^p $\varphi' = R(X \to Y, T_p')$ 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, (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 $T_p'(\mathsf{op}, A)$ (resp. $T_p'(\cdot, A)$), and consider w.l.o.g. CFD^ps 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 $\operatorname{enc}_{\neq}$ as follows. For each $\operatorname{CFD}^p \varphi = R(X \to Y, T_p)$ in $\Sigma^i_{\operatorname{cfd}^p}$, we generate a distinct $\operatorname{cid} \operatorname{id}_{\varphi}$ for it, and do the following.

- Add a tuple t_1 to enc_L such that (a) $t[\operatorname{cid}] = \operatorname{id}_{\varphi}$; (b) for each $A \in X$, $t[A] = -\operatorname{if} T'_p(-,A)$ is '-', and for each non-negation predicate op, $t[A_{\operatorname{op}}] = a$ ' if $T'_p(\operatorname{op},A)$ is 'op a'; (c) we let t[B] = a" in 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 $t[\mathsf{cid}] = \mathsf{id}_{\varphi}$, $t[\mathsf{att}] = `A'$, $t[\mathsf{val}] = `a'$, and $t[\mathsf{pos}] = `\mathsf{LHS}'$ (resp. $t[\mathsf{pos}] = `\mathsf{RHS}'$) if attribute A appears in X (resp. Y).

Example 8. Recall from Fig. 2 CFD^ps φ_2 , φ_3 and φ_4 defined on relation item. The three CFD^ps 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 with only 'null' pattern values do not contribute to error detection, and are thus omitted; (c) enc_{\neq} is empty since all these CFD^ps have no negation patterns. One can easily reconstruct these CFD^ps from tables enc_L, enc_R and enc_{\neq} by collating tuples based on cid.

Encoding CIND^ps. All CIND^ps in $\Sigma_{\mathsf{cindP}}^{(i,j)}$ 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^ps. We use enc to encode the INDs $\mathsf{embedded}$ in CINDP_s , which consists of the following attributes: (1) cid representing the id of a CINDP_s , and (2) those X attributes of R_i and Y attributes of R_j appearing in some CINDP_s in $\Sigma_{\mathsf{cindP}}^{(i,j)}$. 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 .

We populate these tables as follows. For each $\operatorname{CIND}^p \psi = (R_i[X; X_p] \subseteq R_j[Y; Y_p], T_p)$ in $\Sigma_{\mathsf{cind}^p}^{(i,j)}$, we generate a distinct $\mathsf{cid}\ \mathsf{id}_\psi$ 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.

(1) enc			(2) enc $_L$		(3)	(3) enc $_R$		(4) enc $_{\neq}$			
cid	stateL	state _R	cid	type	state	cid	rate	cio	pos	att	val
1	1	1	1	_	null	1	null	1	LHS	type	art
2	1	1	2	_	DL	2	0	2		type	1 1

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

- Add tuples to enc_{\neq} in the same way as their CFD^p counterparts.
- Add tuple t to enc such that $t[\operatorname{cid}] = \operatorname{id}_{\psi}$. For each $i \in [1, m]$, let $t[A_k] = t[B_k] = k$, and t[A] = '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.

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 cardinality of Σ .

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 $vio(D, \Sigma_{cfdP}^i)$. 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 $vio(D, \Sigma_{cindP}^{(i,j)})$. Putting these query answers together, we get $vio(D, \Sigma)$, 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 $\Sigma_{\mathsf{cind}^p}^{(i,j)}$, 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 [11, 6].

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

```
select R_i.*

from R_i, enc<sub>L</sub> L, enc<sub>\neq</sub> N

where R_i.X \asymp L and R_i.X \asymp N and not exists (
    select R_j.*
    from R_j, enc H, enc<sub>R</sub> R, enc<sub>\neq</sub> N
    where R_i.X = R_j.Y and L.\operatorname{cid} = R.\operatorname{cid} and L.\operatorname{cid} = H.\operatorname{cid} and R_j.Y \asymp R and R_j.Y \asymp N)
```

Here (1) $X = \{A_1, \ldots, A_{m1}\}$ and $Y = \{B_1, \ldots, B_{m2}\}$ are the sets of attributes of R_i and R_j appearing in $\Sigma_{\mathsf{cind}^p}^{(i,j)}$, respectively; (2) $R_i.X \times L$ is the conjunction of $L.A_k$ is null or $R_i.A_k = L.A_k$ or $(L.A_k = `-'$

```
and (L.A_{i>} is null or R_i.A_k > L.A_{i>}) and (L.A_{i\geq} is null or R_i.A_k \geq L.A_{k\geq}) and (L.A_{k<} is null or R_i.A_k \leq L.A_{k<}) and (L.A_{i<} is null or R_i.A_k \leq L.A_{i<})
```

for $k \in [1, m_1]$; (3) $R_j \cdot Y \times R$ is defined similarly for attributes in Y; (4) $R_i \cdot X \times N$ 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 = 'LHS' and N.att = 'A_k' and R_i.A_k = N.val);
```

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

Intuitively, (1) $R_i.X \approx L$ and $R_i.X \approx N$ ensure that the R_i tuples selected match the LHS patterns of some CIND^ps in $\Sigma_{\text{cind}^p}^{(i,j)}$; (2) $R_j.Y \approx R$ and $R_j.Y \approx N$ 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 ensures that the R_i tuples selected violate CIND^ps in $\Sigma_{\text{cind}^p}^{(i,j)}$.

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 R_1.* from item R_1, enc_L L, enc_{\neq} N where (L.type is null or R_1.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 (L.state is null or R_1.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 not exists ( select R_2.* from tax R_2, enc H, enc_R R where (H.state_L is null or H.state_R is null or H.state_L! = H.state_R or R_2.state = R_1.state) and L.cid = H.cid and L.cid = R.cid and (R.rate is null or R_2.rate = R.rate or R.rate = '_') and not exist ( select * from N where N.cid = R.cid and N.pos = 'RHS' and N.att = 'rate'))
```

Optimization. The SQL queries generated for error detection can be simplified as follows. (1) As shown in Example 10, when checking patterns imposed by 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 $\operatorname{enc}_{\neq}$, the queries need not check A either. (2) One can further partition $\Sigma_{\operatorname{cind}}^{(i,j)}$ such that each group has a bounded number of attributes. The price that comes with it is that more than one SQL query will be needed. A balance between the number of groups and the number of queries can be struck based on the dependencies and relational schemas by using e.q., the clustering method [16].

Example 11. An extreme case of partitioning is to generate an SQL query for each CIND^p. For example, a query Q' for validating ψ_1 of Fig. 3 alone is:

```
select R_1.* from item R_1
where R_1.type ! = 'art' and not exists (select R_2.* from tax R_2
where R_2.state = R_1.state)
```

Contrasting Q' with query Q of Example 10, one can see that Q' is significantly simpler than Q. Note that a single test R_1 type ! = 'art' is used in query Q', instead of using: **not exists** (select N.* from $enc_{\neq} N$ where N.cid = 1 and

 $N.\mathsf{pos} = \text{`LHS'}$ and $N.\mathsf{att} = \text{`type'}$ and $R_1.\mathsf{type} = N.\mathsf{val}$). This is because enc_{\neq} contains a single tuple, and thus the complex test in terms of **not exists** is unnecessary. In general, if only a small number of negation patterns are present, one can use a conjunctions of ! = tests instead of the **not exists** test.

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 in real-life data (see, e.g., [9] for a comprehensive survey). As an extension of traditional FDs, CFDs were developed in [11], for improving the quality of data. It was shown in [11] that the satisfiability and implication problems for CFDs are NP-complete and conp-complete, respectively. Along the same lines, CINDs were proposed in [7] to extend INDs. It was shown [7] that the satisfiability and implication problems for CINDs are in constant time and EXPTIME-complete, respectively. SQL techniques were developed in [11] to detect errors by using CFDs, but have not been studied for CINDs. This work extends the static analyses of conditional dependencies of [11, 7], 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 [6], and to specify patterns in terms of value ranges [15]. While CFD^ps are more powerful than the extension of [15], they cannot express disjunctions studied in [6]. 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 [8, 15, 12] and for repairing data based on either CFDs [10], or traditional FDs and INDs taken together [4]. 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 [3, 5, 17, 18]. While the languages of [3, 17] cannot express CFDs, constraint-generating dependencies (CGDs) of [3] and constrained tuple-generating dependencies (CTGDs) of [18] 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^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^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 and eCFDs [6]), 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 [8, 15, 12]. Third, we plan to extend the methods of [4, 10] to repair data based on CFD^p s and $CIND^p$ s, instead of using CFD^p s [10] and traditional FDs and INDS [4].

References

- S. Abiteboul, R. Hull, and V. Vianu. Foundations of Databases. Addison-Wesley, 1995.
- M. Arenas, L. E. Bertossi, and J. Chomicki. Consistent query answers in inconsistent databases. In PODS, 1999.
- M. Baudinet, J. Chomicki, and P. Wolper. Constraint-Generating Dependencies. J. Comput. Syst. Sci., 59(1):94-115, 1999.
- 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.
- P. D. Bra and J. Paredaens. Conditional dependencies for horizontal decompositions. In ICALP, 1983.
- L. Bravo, W. Fan, F. Geerts, and S. Ma. Increasing the expressivity of conditional functional dependencies without extra complexity. In *ICDE*, 2008.
- L. Bravo, W. Fan, and S. Ma. Extending dependencies with conditions. In VLDB, 2007
- 8. F. Chiang and R. J. Miller. Discovering data quality rules. In VLDB, 2008.
- 9. J. Chomicki. Consistent query answering: Five easy pieces. In ICDT, 2007.
- G. Cong, W. Fan, F. Geerts, X. Jia, and S. Ma. Improving data quality: Consistency and accuracy. In VLDB, 2007.
- W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies. TODS, 33(2), 2008.
- W. Fan, F. Geerts, L. V. Lakshmanan, and M. Xiong. Discovering conditional functional dependencies. In *ICDE*, 2009.
- 13. Full version. http://homepages.inf.ed.ac.uk/sma1/dexa-full.pdf.
- 14. M. Garey and D. Johnson. Computers and Intractability: A Guide to the Theory of NP-Completeness. W. H. Freeman and Company, 1979.
- L. Golab, H. J. Karloff, F. Korn, D. Srivastava, and B. Yu. On generating nearoptimal tableaux for conditional functional dependencies. In VLDB, 2008.
- J. Han and M. Kamber. Data Mining: Concepts and Techniques. Morgan Kaufmann, 2000.
- 17. M. J. Maher. Constrained dependencies. TCS, 173(1):113-149, 1997.
- M. J. Maher and D. Srivastava. Chasing Constrained Tuple-Generating Dependencies. In PODS, 1996.
- R. van der Meyden. The complexity of querying indefinite data about linearly ordered domains. JCSS, 54(1), 1997.

APPENDIX: Proofs

Proposition 1. The satisfiability problem for CFD^p s is NP-complete.

Proof. The lower bound follows from the NP-hardness of their CFDs counterparts [11], since CFDs are a special case of CFD p s.

П

We next show that the problem is in NP by presenting an NP algorithm that, given a set Σ of CFD^ps defined on a relational schema R, determines whether Σ is satisfiable. The satisfiability problem has the following small model property: If there exists a nonempty instance I of R such that $I \models \Sigma$, then for any tuple $t \in I$, $I_t = \{t\}$ is an instance of R and $I_t \models \Sigma$. Thus it suffices to consider single-tuple instances $I = \{t\}$ for deciding whether Σ is satisfiable.

Assume w.l.o.g. that the attributes $\mathsf{attr}(R) = \{A_1, \ldots, A_m\}$ and the number of all pattern tuples of T_p in Σ is h. For each $i \in [1, m]$, define the active domain of A_i to be a set $\mathsf{adom}(A_i) = C_0 \cup C_1$, where (1) set C_0 consists of all constants in $T_p[A_i]$ of all pattern tableaux T_p in Σ and let $C_0 = \{a_1\}$, where $a_1 \in \mathsf{dom}(A_i)$, if C_0 is empty, and (2) set C_1 contains the following set of constants:

- Arrange all constants in C_0 in the increasing order, and assume that the resulting $C_0 = \{a_1, \ldots, a_k\}$ $(k \ge 1)$;
- Add a constant $b_0 \in dom(A_i)$ to C_1 such that $b_0 < a_1$ if there exists one;
- For each $j \in [1, k-1]$, add a constant $b_j \in \text{dom}(A_i)$ to C_1 such that $b_j > a_j$ and $b_j < a_{j+1}$ if there exists one;
- Add a constant $b_k \in dom(A_i)$ to C_1 such that $b_k > a_k$ if there exists one.

Moveover, the number of elements in $\mathsf{adom}(A_i)$ is at most O(2*h+1). Then one can easily verify that Σ is satisfiable iff there exists a mapping ρ from $t[A_i]$ to $\mathsf{adom}(A_i)$ $(i \in [1, m])$ such that $I = \{(\rho(t[A_1]), \ldots, \rho(t[A_m]))\}$ and $I \models \Sigma$.

Based on these, we give an NP algorithm as follows: (1) Guess an instance, which contains a single tuple t of R such that $t[A_i] \in \mathsf{adom} A_i$ 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.

Theorem 2. In the absence of finite domain attributes, the satisfiability problem for CFD^p s remains NP-complete.

Proof. The problem is in NP following from Proposition 1. We next show that the problem is NP-hard by reduction from the 3SAT problem, which is NP-complete (cf. [14]). Consider an instance $\phi = C_1 \wedge \cdots \wedge C_n$ of 3SAT, where all the variables in ϕ are x_1, \ldots, x_m, C_j is of the form $y_{j_1} \vee y_{j_2} \vee y_{j_3}$, and moreover, for $i \in [1, 3]$, y_{j_i} is either $x_{p_{j_i}}$ or $\overline{x_{p_{j_i}}}$ for $p_{j_i} \in [1, m]$. Given ϕ , we construct a relation schema R and a set Σ of CFD^{Ps} defined on R such that ϕ is satisfiable iff Σ is satisfiable. (1) Define relation $R(X_1, \ldots, X_m, C_1, \ldots, C_n, Z)$, where all attributes share an infinite domain dom and constant $a \in$ dom. Intuitively, for each R tuple t, $t[X_1, \ldots, X_m]$ specifies a truth assignment ξ for variables x_1, \ldots, x_m of ϕ , and $t[C_i]$ and t[Z] are the truth values of clause C_i and sentence ϕ w.r.t. ξ respectively.

(2) Let the set Σ of CFD^ps be $\Sigma_0 \cup \Sigma_1 \cup \ldots \cup \Sigma_n \cup \Sigma_{n+1}$.

- Σ_0 contains n+1 CFD^ps. Intuitively, Σ_0 encodes the relationships between the truth values of clauses C_1, \ldots, C_n and the truth value of sentence ϕ . For each clause C_i ($i \in [1, n]$), we define CFD^p $\varphi_i = (C_1, \ldots, C_n \to Z, T_{pi})$ such that $T_{pi} = \{t_{pi}\}$ and $t_{pi}[C_i, Z] = (\neq a \parallel \neq a)$ and $t_{pi}[C_j] = `_' (j \neq a)$

- $i, j \in [1, n]$). Moreover, CFD^p $\varphi_0 = (C_1, \dots, C_n \to Z, T_{p0})$, where $T_{p0} = \{(=a, \dots, =a \parallel =a)\}$. Intuitively, we use $\neq a$ and =a to represent false and true for the truth values of the clauses and the sentence.
- For $i \in [1, n]$, Σ_i contains 8 CFD^ps. Intuitively, Σ_i encodes the relationships between the truth values of the three variables in clause C_i and clause C_i . Consider clause $C_i = x_{j_1} \vee \overline{x_{j_2}} \vee x_{j_3}$ of ϕ , where $1 \leq j_1, j_2, j_3 \leq m$, we define CFD^ps $\varphi_{i,0} = (X_{j_1}, X_{j_2}, X_{j_3} \to C_i, T_{pi,0})$, where $T_{pi,0} = \{(< a, < a, < a \parallel = a)\}$, and $\varphi_{i,2} = (X_{j_1}, X_{j_2}, X_{j_3} \to C_i, T_{pi,2})$, where $T_{pi,2} = \{(< a, \geq a, < a \parallel \neq a)\}$. Intuitively, we use < a and $\geq a$ to represent false and true for a variable. Similarly, we can define the rest 6 CFD^ps $\varphi_{i,1}, \varphi_{i,3}, \varphi_{i,4}, \varphi_{i,5}, \varphi_{i,6}$ and $\varphi_{i,7}$.
- Similarly, we can define the rest 6 CFD^ps $\varphi_{i,1}$, $\varphi_{i,3}$, $\varphi_{i,4}$, $\varphi_{i,5}$, $\varphi_{i,6}$ and $\varphi_{i,7}$. $-\Sigma_{n+1}$ contains a single CFD^p $\varphi_{n+1} = (Z \to Z, T_{p(n+1)})$, where $T_{p(n+1)} = \{(-\parallel = a)\}$. Intuitively, φ_{n+1} requires that for any tuple t of R, t[Z] = a.

Observe that Σ contains 8*m+n+2 CFD^ps. Thus the reduction is in PTIME. We now show that ϕ is satisfiable if and only if Σ is satisfiable.

Suppose first Σ is satisfiable, then there exists a nonempty instance I of R such that $I \models \Sigma$. For any tuple $t \in I$, (1) Σ_{n+1} forces t[Z] = a, (2) Σ_0 forces $t[C_1, \ldots, C_n] = (a, \ldots, a)$, and for each clause C_i ($i \in [1, n]$) with variables $X_{j_1}, X_{j_2}, X_{j_3}$, Σ_j forces that $t[X_{j_1}, X_{j_2}, X_{j_3}]$ does not match the LHS of the CFD^p that forces $t[C_i] \neq a$. From t, we can construct a truth assignment ξ of ϕ such that $\xi(x_i) = false$ if $t[X_i] < a$ and $\xi(x_i) = true$ 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) $t[C_1, \ldots, C_n, Z] = (a, \ldots, a)$ and (b) for each $i \in [1, m]$, $t[X_i] = a_i$ such that $a_i \in \text{dom}$, $a_i \geq a$ if $\xi(x_i) = true$ and $a_i < a$ otherwise. Let $I = \{t\}$, then one can easily verify that $I \models \Sigma$. \square

Proposition 3. A set Σ of CIND^ps is always satisfiable.

Proof. It suffices to show that given a set Σ of $CIND^p$ s over a database schema $\mathcal{R}(R_1,\ldots,R_n)$, we can 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 $\mathsf{adom}(A)$, which consists of certain data values in $\mathsf{dom}(A)$. Second, using these active domains, we construct D.

We start with the construction of active domains. (1) For each attribute A, initialize $\mathsf{adom}(A)$ along the same lines as the one for CFD^p s in Proposition 1; (2) For each CIND^p $(R_a[A_1,A_2,\ldots,A_m;X_p]\subseteq R_b[B_1,\ldots,B_m;Y_p],T_p)$ in Σ , let $\mathsf{adom}(B_i)=\mathsf{adom}(B_i)\cup \mathsf{adom}(A_i)$ for each $i\in[1,m]$; (3) This rule is repeatedly applied until a fixpoint of $\mathsf{adom}(A)$ is reached for every attribute A. It is easy to verify that this process terminates since we start with a finite set of data values.

We next construct D. For each relation $R_i(A_1, \ldots, A_k) \in \mathcal{R}$, we define $I_i = \mathsf{adom}(A_1) \times \cdots \times \mathsf{adom}(A_k)$, where \times is the *Cartesian Product* operation. Let $D = \{I_1, \ldots, I_n\}$, then it is easy to verify that D is nonempty and $D \models \Sigma$. \square

Proposition 5. The implication problem for CFD^p s is conp-complete.

Proof. The lower bound follows from the conp-hardness of their CFDs counterparts [11], since CFDs are a special case of CFD p s.

We show that the problem is in conp by presenting an NP algorithm for its complement problem, *i.e.*, for determining whether $\Sigma \not\models \varphi$. The algorithm is

based on a small model property: if $\varphi = (R: X \to Y, T_p)$ 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 $t_1[X] = t_2[X] \times T_p[X]$, but either $t_1[Y] \not\neq t_2[Y]$ 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 for deciding whether $\Sigma \not\models \varphi$.

Assume w.l.o.g. that the attributes $\mathsf{attr}(R) = \{A_1, \dots, A_m\}$. For each $i \in [1, m]$, define the active domain of A_i to be a set $\mathsf{adom}(A_i)$ in the same way as Proposition 1. Then one can easily verify that $\Sigma \not\models \varphi$ iff there exist two mappings ρ_1, ρ_2 from $t[A_i]$ to $\mathsf{adom}(A_i)$ $(i \in [1, m])$ such that $I = \{(\rho_1(t[A_1]), \dots, \rho_1(t[A_m])), (\rho_2(t[A_1]), \dots, \rho_2(t[A_m]))\}$, $I \models \Sigma$ and $I \not\models \varphi$.

Based on these, we give an NP algorithm as follows: (1) Guess two tuples t_1, t_2 of R such that $t_1[A_i], t_2[A_i] \in \mathsf{adom}(A_i)$ 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 Σ . Hence the algorithm is in NP, and the problem is in conp. \square

Theorem 6. In the absence of finite domain attributes, the implication problem for CFD^ps remains conp-complete.

Proof. The problem is in conp following from Proposition 5. We next show that the problem is conp-hard by reduction from the 3SAT problem to the complement problem of the implication problem, where 3SAT is NP-complete (cf. Proposition 2). Given an instance ϕ of 3SAT, we construct a relation schema R, a set $\Sigma \cup \{\varphi\}$ of CFD^ps defined on R, such that ϕ is satisfiable iff $\Sigma \not\models \varphi$.

The relation schema R and the set Σ of CFD^ps are the same as the corresponding ones in Proposition 2. Moreover, φ is defined as $(Z \to Z, T_p)$, where $T_p = \{(_ || \neq a)\}$. Intuitively, φ requires that for any tuple t of R, $t[Z] \neq a$. 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.

Proposition 7. The implication problem for $CIND^p$ s is EXPTIME-complete. \Box

Proof. The implication problem for CINDs is EXPTIME-hard [7]. The lower bound carries over to CIND^ps, which subsume CINDs. We next show that the problem is in EXPTIME by presenting an EXPTIME algorithm that, given a set $\Sigma \cup \{\psi\}$ of CIND^ps defined on a database schema \mathcal{R} , determines whether $\Sigma \models \psi$.

Consider $\mathcal{R} = (R_1, \dots, R_n)$ and $\psi = (R_a[X; X_p] \subseteq R_b[Y; Y_p], T_p)$. Moreover, for each attribute A, define the active domain $\mathsf{adom}(A)$ of A based on $\Sigma \cup \{\psi\}$ along the same way as Proposition 3. One can easily verify that if $\Sigma \not\models \psi$, there must exist an instance D of \mathcal{R} such that $D \models \Sigma$ and $D \not\models \psi$, and moreover, the instance D only consists of constants in those predefined active domains.

Based on these, we give the EXPTIME algorithm as follows: (1) Initialize instance $I_a = \{t_a\}$ of R_a such that $t_a[A] \in \mathsf{adom}(A)$ for each attribute $A \in \mathsf{attr}(R_a)$, and $I_i = \{\}$ for any other relation R_i $(1 \le i \le n)$ in \mathcal{R} . (2) For each CIND^p $(R_i[X'; X'_p] \subseteq R_j[Y'; Y'_p], T'_p)$ in Σ and each tuple $t_i \in I_i$ such that $t_i \asymp T'_p[X'_p]$, add the following set of tuples to I_j . Let $\mathsf{attr}(R_j) \setminus Y' = E_1, \ldots, E_k$ and $S = \mathsf{adom}(E_1) \times \ldots \times \mathsf{adom}(E_k)$, where \times is the Cartesian Product operation. For each tuple $t \in S$, if $t \asymp T'_p[Y'_p]$, add a tuple t_j to I_j such that $t_j[\mathsf{attr}(R_j) \setminus Y'] = t[\mathsf{attr}(R_j) \setminus Y']$ and $t_j[Y'] = t_i[X']$. This rule is repeatedly applied until all

instances reach their fixpoints. (3) Check whether $D = \{I_1, \ldots, I_n\}$ satisfies Σ , but not ψ . If so, the algorithm returns "no", and otherwise it repeats steps (1) (2) and (3). If all cases in step (1) are tested and no instance D is found such that $D \models \Sigma$ and $D \not\models \psi$, the algorithm returns "yes".

Let $|\Sigma \cup \{\psi\}|$ be the size of Σ and ψ , and $|\mathcal{R}|$ be the sum of arities of relations in \mathcal{R} . Then it is easy to know that (a) the instance D in step (3) contains at most $O(|\Sigma \cup \{\psi\}|^{|\mathcal{R}|})$ tuples, and (b) there are at most $O(|\Sigma \cup \{\psi\}|^{|\mathcal{R}|})$ cases for step (1). Following from these, we can conclude that the algorithm is indeed in EXPTIME, and thus the problem is in EXPTIME.

Theorem 8. In the absence of finite domain attributes, the implication problem for $CIND^p$ s remains EXPTIME-complete.

Proof. The problem is in EXPTIME following from Proposition 7. We next show that the problem is EXPTIME-hard by reduction from the implication problem for CINDs in the general setting, which is EXPTIME-complete (cf. [7]). Given a set $\Sigma \cup \{\psi\}$ of CINDs defined on a database schema $\mathcal{R} = (R_1, \ldots, R_n)$, we construct a database schema $\mathcal{R}' = (R'_1, \ldots, R'_n)$, where 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}' . And, moreover, $\Sigma \models \psi$ iff $\Sigma' \models \psi'$.

We start with constructing \mathcal{R}' . For each $R_i(A_1,\ldots,A_k)$ of \mathcal{R} , we define $R_i'(A_1',\ldots,A_k')$ such that for each attribute A_j' $(j\in[1,k])$, $\mathsf{dom}(A_j') = \mathsf{dom}(A_j)$ if $\mathsf{dom}(A_j)$ is infinite and $\mathsf{dom}(A_j')$ is a totally ordered infinite domain, say integer or real, if $\mathsf{dom}(A_j)$ is finite. Moreover, we define a mapping $\rho_{i,j}$ from finite domain $\mathsf{dom}(A_j)$ to intervals of infinite domain $\mathsf{dom}(A_j')$: (1) Randomly choose $h = |\mathsf{dom}(A_j)| - 1$ values $\{b_1,\ldots,b_h\}$ in $\mathsf{dom}(A_j')$ such that for $i \in [1,h-1]$, $b_i < b_{i+1}$ and there exists a constant $b \in \mathsf{dom}(A_j')$ such that $b > b_i$ and $b \le b_{i+1}$. Moreover, we require that there exist two constants $b_0, b_h \in \mathsf{dom}(A_j')$ such that $b_0 \le b_1$ and $b_h > b_{h-1}$. Note that this is always doable since $\mathsf{dom}(A_j')$ is a totally ordered infinite domain. (2) Assume $\mathsf{dom}(A_j) = \{a_0, a_1, \ldots, a_h\}$, then $\rho_{i,j}(a_0) = \le b_1$, $\rho_{i,j}(a_h) = \ge b_h$, and $\rho_{i,j}(a_i) = \ge b_i \le b_{i+1}$ for $i \in [1,h-1]$.

= '≤ b_1 ', $\rho_{i,j}(a_h)$ = '> b_h ', and $\rho_{i,j}(a_i)$ = '> b_i ∧≤ b_{i+1} ' for $i \in [1, h-1]$. We next define Σ' and ψ' on \mathcal{R}' based on the mappings defined above. For each CIND ($R_a[X; A_1, \ldots, A_{m1}] \subseteq R_b[Y; B_1, \ldots, B_{m2}], T_p$) in Σ and each $t_p \in T_p$, we define a CIND^p ($R'_a[X'; A'_1, \ldots, A'_{m1}] \subseteq R'_b[Y'; B'_1, \ldots, B'_{m2}], T'_p$) such that $T'_p = \{t'_{p1}, t'_{p2}\}$ for each attribute A' in A'_1, \ldots, A'_{m1} or $B'_1, \ldots, B'_{m2}, (1)$ $t'_{p1}[A'] = t'_{p2}[A'] =$ '2' when $t_p[A] =$ '2', (2) $t'_{p1}[A'] = t'_{p2}[A'] =$ '5 a' when $t_p[A] =$ 'a' and dom(A) is infinite, (3) $t'_{p1}[A'] = t'_{p2}[A'] =$ '5 bh' when $t_p[A] =$ 'a', $\rho(a) =$ '5 bh' and dom(A) is finite, and (5) $t'_{p1}[A'] =$ '5 bh' and $t'_{p2}[A'] =$ '6 b2' when $t_p[A] =$ 'a', $\rho(a) =$ '5 b1' and dom(A) is finite, and (5) $t'_{p1}[A'] =$ 5 b1' and $t'_{p2}[A'] =$ 6 b2' when $t_p[A] =$ 'a', $\rho(a) =$ 5 b1' h and dom(A) is finite, and (5) h finite. Similarly, we can define h from h based on the mappings.

Finally, one can easily verify that $\Sigma \models \psi$ iff $\Sigma' \models \psi'$. Following from this, the problem is EXPTIME-hard.