

---

## THE SHAPLEY VALUE OF INCONSISTENCY MEASURES FOR FUNCTIONAL DEPENDENCIES

ESTER LIVSHITS AND BENNY KIMELFELD

Technion, Haifa, Israel  
*e-mail address*: {esterliv,bennyk}@cs.technion.ac.il

---

**ABSTRACT.** Quantifying the inconsistency of a database is motivated by various goals including reliability estimation for new datasets and progress indication in data cleaning. Another goal is to attribute to individual tuples a level of responsibility to the overall inconsistency, and thereby prioritize tuples in the explanation or inspection of errors. Therefore, inconsistency quantification and attribution have been a subject of much research in knowledge representation and, more recently, in databases. As in many other fields, a conventional responsibility sharing mechanism is the Shapley value from cooperative game theory. In this article, we carry out a systematic investigation of the complexity of the Shapley value in common inconsistency measures for functional-dependency (FD) violations. For several measures we establish a full classification of the FD sets into tractable and intractable classes with respect to Shapley-value computation. We also study the complexity of approximation in intractable cases.

### 1. INTRODUCTION

Inconsistency measures for knowledge bases have received considerable attention from the Knowledge Representation (KR) and Logic communities [KLM03, Kni03, HK06, GH06, HK08, HK10, GH17, Thi17]. More recently, inconsistency measures have also been studied from the database viewpoint [Ber18, LKT<sup>+</sup>21]. Such measures quantify the extent to which the database violates a set of integrity constraints. There are multiple reasons why one might be using such measures. For one, the measure can be used for estimating the usefulness or reliability of new datasets for data-centric applications such as business intelligence [CPRT15]. Inconsistency measures have also been proposed as the basis of progress indicators for data-cleaning systems [LKT<sup>+</sup>21]. Finally, the measure can be used for attributing to individual tuples a level of responsibility to the overall inconsistency [MLJ11, Thi09], thereby prioritize tuples in the explanation/inspection/correction of errors.

**Example 1.1.** Figure 1 depicts an inconsistent database that stores a train schedule. For example, the tuple  $f_1$  states that train number 16 will depart from the New York Penn Station at time 1030 and arrive at the Boston Back Bay Station after 315 minutes. Assume that we have the functional dependency stating that the train number and departure time determine the departure station. All tuples in the database are involved in violations of this constraint, as they all agree on the train number and departure time, but there is some disagreement on the departure station. Hence, one can argue that every fact in the database affects the overall level of inconsistency in the database. But how should we measure the

---

*Key words and phrases*: Shapley value, inconsistent databases, functional dependencies, database repairs.

fact	train	departs	arrives	time	duration
$f_1$	16	NYP	BBY	1030	315
$f_2$	16	NYP	PVD	1030	250
$f_3$	16	PHL	WIL	1030	20
$f_4$	16	PHL	BAL	1030	70
$f_5$	16	PHL	WAS	1030	120
$f_6$	16	BBY	PHL	1030	260
$f_7$	16	BBY	NYP	1030	260
$f_8$	16	BBY	WAS	1030	420
$f_9$	16	WAS	PVD	1030	390

FIGURE 1. The inconsistent database of our running example.

*responsibility* of the tuples to this inconsistency? For example, which of the tuples  $f_1$  and  $f_3$  has a greater contribution to inconsistency? To this end, we can adopt some conventional concepts for responsibility sharing, and in this article we study the computational aspects involved in the measurement of those.  $\diamond$

A conventional approach to dividing the responsibility for a quantitative property (here an inconsistency measure) among entities (here the database tuples) is the *Shapley value* [Sha53], which is a game-theoretic formula for wealth distribution in a cooperative game. The Shapley value has been applied in a plethora of domains, including economics [Gul89], law [Nen03], environmental science [PZ03, LZS15], social network analysis [NN11], physical network analysis [MCL<sup>+</sup>10], and advertisement [BDG<sup>+</sup>19]. In data management, the Shapley value has been used for determining the relative contribution of features in machine-learning predictions [LF18, LL17], the responsibility of tuples to database queries [RKL20, LBKS20, BG20], and the reliability of data sources [CPRT15].

The Shapley value has also been studied in a context similar to the one we adopt in this article—assigning a level of inconsistency to statements in an inconsistent knowledge base [HK10, YVCB18, MLJ11, Thi09]. Hunter and Konieczny [HK06, HK10, HK08] use the maximal Shapley value of one inconsistency measure in order to define a new inconsistency measure. Grant and Hunter [GH15] considered information systems distributed along data sources of different reliabilities, and apply the Shapley value to determine the expected blame of each statement to the overall inconsistency. Yet, with all the investigation that has been conducted on the Shapley value of inconsistency, we are not aware of any results or efforts regarding the computational complexity of calculating this value.

**Example 1.2.** Let us define the following cooperative game over the database of Figure 1. We have nine players—the tuples of the database. One of the measures that we consider for quantifying the level of inconsistency of a coalition of players is the number of tuple pairs in this group that violate the constraints. For example, consider the constraint defined in Example 1.1. The inconsistency level of the group  $\{f_1, f_3, f_5\}$  is 2, as there are two conflicting tuple pairs:  $\{f_1, f_3\}$  and  $\{f_1, f_5\}$ . The inconsistency level of the entire database is 29, as this is the total number of conflicting pairs in the database. The Shapley value allows us to measure the contribution of each individual tuple to the overall inconsistency level. For example, the Shapley value of the tuple  $f_1$ , in this case, will be lower than the Shapley

TABLE 1. The complexity of the (exact ; approximate) Shapley value of different inconsistency measures.

	lhs chain	no lhs chain, PTime c-repair	other
$\mathcal{I}_d$	PTime	FP <sup>#P</sup> -complete ; FPRAS	
$\mathcal{I}_{MI}$	PTime		
$\mathcal{I}_P$	PTime		
$\mathcal{I}_R$	PTime	? ; FPRAS	NP-hard [LKR20] ; no FPRAS
$\mathcal{I}_{MC}$	PTime	FP <sup>#P</sup> -complete [LK17] ; ?	

value of the tuple  $f_3$  (we will later show how this value is computed), which indicates that  $f_3$  has a higher impact on the inconsistency than  $f_1$ .

In this work, we embark on a systematic analysis of the complexity of the Shapley value of database tuples relative to inconsistency measures, where the goal is to calculate the contribution of a tuple to inconsistency. Our main results are summarized in Table 1. We consider inconsistent databases with respect to functional dependencies (FDs), and basic measures of inconsistency following Bertossi [Ber19] and Livshits, Ilyas, Kimelfeld and Roy [LKT<sup>+</sup>21]. We note that these measures are all adopted from the measures studied in the aforementioned KR research. In our setting, an individual tuple affects the inconsistency of only its containing relation, since the constraints are FDs. Hence, our analysis focuses on databases with a single relation; in the end of each relevant section, we discuss the generalization to multiple relations. While most of our results easily extend to multiple relations, some extensions require a more subtle proof.

More formally, we investigate the following computational problem for any fixed combination of a relational signature, a set of FDs, and an inconsistency measure: given a database and a tuple, compute the Shapley value of the tuple with respect to the inconsistency measure. As Table 1 shows, two of these measures are computable in polynomial time:  $\mathcal{I}_{MI}$  (number of FD violations) and  $\mathcal{I}_P$  (number of problematic facts that participate in violations). For two other measures, we establish a full dichotomy in the complexity of the Shapley value:  $\mathcal{I}_d$  (the drastic measure—0 for consistency and 1 for inconsistency) and  $\mathcal{I}_{MC}$  (number of maximal consistent subsets, a.k.a. repairs). The dichotomy in both cases is the same: when the FD set has, up to equivalence, an lhs chain (i.e., the left-hand sides form a chain w.r.t. inclusion [LK17]), the Shapley value can be computed in polynomial time; in any other case, it is FP<sup>#P</sup>-hard (hence, requires at least exponential time under conventional complexity assumptions). In the case of  $\mathcal{I}_R$  (the minimal number of tuples to delete for consistency), the problem is solvable in polynomial time in the case of an lhs chain, and NP-hard whenever it is intractable to find a cardinality repair [LKR20]; however, the problem is open for every FD set in between, for example, the bipartite matching constraint  $\{A \rightarrow B, B \rightarrow A\}$ .

We also study the complexity of approximating the Shapley value and show the following (as described in Table 1). First, in the case of  $\mathcal{I}_d$ , there is a (multiplicative) fully polynomial-time approximation scheme (FPRAS) for every set of FDs. In the case of  $\mathcal{I}_{MC}$ , approximating the Shapley value of *any* intractable (non-lhs-chain) FD set is at least as hard as approximating the number of maximal matchings of a bipartite graph—a long standing open problem [JR18]. In the case of  $\mathcal{I}_R$ , we establish a full dichotomy, namely FPRAS

vs. hardness of approximation, that has the same separation as the problem of finding a cardinality repair.

This article is the full version of a conference publication [LK21]. We have added all of the proofs, intermediate results and algorithms that were excluded from the conference version. In particular, we have included in this version the proofs of Observation 3.2, Lemma 5.4, Lemma 6.3, Lemma 7.2, and Lemma 7.3, and the algorithms of Figures 5, 8, and 10. Furthermore, the results of the conference publication have been restricted to schemas with a single relation symbol. While some of the results (e.g., all of the lower bounds) immediately generalize to schemas with multiple relation symbols, some generalizations (in particular, the upper bounds for  $\mathcal{I}_d$  and  $\mathcal{I}_{MC}$ ) require a more subtle analysis that we provide in this article. We generalize the upper bounds for all the measures to schemas with multiple relation symbols, in the corresponding sections.

The rest of the article is organized as follows. After presenting the basic notation and terminology in Section 2, we formally define the studied problem and give initial observations in Section 3. In Section 4, we describe polynomial-time algorithms for  $\mathcal{I}_{MI}$  and  $\mathcal{I}_P$ . Then, we explore the measures  $\mathcal{I}_d$ ,  $\mathcal{I}_R$  and  $\mathcal{I}_{MC}$  in Sections 5, 6 and 7, respectively. We conclude and discuss future directions in Section 8.

## 2. PRELIMINARIES

We begin with preliminary concepts and notation that we use throughout the article.

**2.1. Database Concepts.** By a *relational schema* we refer to a sequence  $(A_1, \dots, A_n)$  of attributes. A database  $D$  over  $(A_1, \dots, A_n)$  is a finite set of tuples, or *facts*, of the form  $(c_1, \dots, c_n)$ , where each  $c_i$  is a constant from a countably infinite domain. For a fact  $f$  and an attribute  $A_i$ , we denote by  $f[A_i]$  the value associated by  $f$  with the attribute  $A_i$  (that is,  $f[A_i] = c_i$ ). Similarly, for a sequence  $X = (A_{j_1}, \dots, A_{j_m})$  of attributes, we denote by  $f[X]$  the tuple  $(f[A_{j_1}], \dots, f[A_{j_m}])$ . Generally, we use letters from the beginning of the English alphabet (i.e.,  $A, B, C, \dots$ ) to denote single attributes and letters from the end of the alphabet (i.e.,  $X, Y, Z, \dots$ ) to denote sets of attributes. We may omit stating the relational schema of a database  $D$  when it is clear from the context or irrelevant.

A *Functional Dependency* (FD for short) over  $(A_1, \dots, A_n)$  is an expression of the form  $X \rightarrow Y$ , where  $X, Y \subseteq \{A_1, \dots, A_n\}$ . We may also write the attribute sets  $X$  and  $Y$  by concatenating the attributes (e.g.,  $AB \rightarrow C$  instead of  $\{A, B\} \rightarrow \{C\}$ ). A database  $D$  satisfies  $X \rightarrow Y$  if every two facts  $f, g \in D$  that agree on the values of the attributes of  $X$  also agree on the values of the attributes of  $Y$  (that is, if  $f[X] = g[X]$  then  $f[Y] = g[Y]$ ). A database  $D$  *satisfies* a set  $\Delta$  of FDs, denoted by  $D \models \Delta$ , if  $D$  satisfies every FD of  $\Delta$ . Otherwise,  $D$  *violates*  $\Delta$  (denoted by  $D \not\models \Delta$ ). Two FD sets over the same relational schema are *equivalent* if every database that satisfies one of them also satisfies the other.

Let  $\Delta$  be a set of FDs and  $D$  a database (which may violate  $\Delta$ ). A *repair* (of  $D$  w.r.t.  $\Delta$ ) is a maximal consistent subset of  $D$ ; that is,  $E \subseteq D$  is a repair if  $E \models \Delta$  but  $E' \not\models \Delta$  for every  $E \subsetneq E'$ . A *cardinality repair* (or *c-repair* for short) is a repair of maximum cardinality; that is, it is a repair  $E$  such that  $|E| \geq |E'|$  for every repair  $E'$ .

**Example 2.1.** Consider again the database of Figure 1 over the relational schema

(train, departs, arrives, time, duration).

The FD set  $\Delta$  consists of the two FDs:

- train time  $\rightarrow$  departs
- train time duration  $\rightarrow$  arrives

The first FD states that the departure station is determined by the train number and departure time, and the second FD states that the arrival station is determined by the train number, the departure time, and the duration of the ride.

Observe that the database of Figure 1 violates the FDs as all the facts refer to the same train number and departure time, but there is no agreement on the departure station. Moreover, the facts  $f_6$  and  $f_7$  also agree on the duration, but disagree on the arrival station. The database has five repairs: (a)  $\{f_1, f_2\}$ , (b)  $\{f_3, f_4, f_5\}$ , (c)  $\{f_6, f_8\}$ , (d)  $\{f_7, f_8\}$ , and (e)  $\{f_9\}$ ; only the second one is a cardinality repair.  $\diamond$

**2.2. Shapley Value.** A *cooperative game* of a set  $A$  of players is a function  $v : \mathcal{P}(A) \rightarrow \mathbb{R}$ , where  $\mathcal{P}(A)$  is the power set of  $A$ , such that  $v(\emptyset) = 0$ . The value  $v(B)$  should be thought of as the joint wealth obtained by the players of  $B$  when they cooperate. The *Shapley value* of a player  $a \in A$  measures the contribution of  $a$  to the total wealth  $v(A)$  of the game [Sha53], and is formally defined by

$$\text{Shapley}(A, v, a) \stackrel{\text{def}}{=} \frac{1}{|A|!} \sum_{\sigma \in \Pi_A} (v(\sigma_a \cup \{a\}) - v(\sigma_a))$$

where  $\Pi_A$  is the set of all permutations over the players of  $A$  and  $\sigma_a$  is the set of players that appear before  $a$  in the permutation  $\sigma$ . Intuitively, the Shapley value of a player  $a$  is the expected contribution of  $a$  to a subset constructed by drawing players randomly one by one (without replacement), where the contribution of  $a$  is the change to the value of  $v$  caused by the addition of  $a$ . An alternative formula for the Shapley value, that we will use in this article, is the following.

$$\text{Shapley}(A, v, a) \stackrel{\text{def}}{=} \sum_{B \subseteq A \setminus \{a\}} \frac{|B|! \cdot (|A| - |B| - 1)!}{|A|!} (v(B \cup \{a\}) - v(B))$$

Observe that  $|B|! \cdot (|A| - |B| - 1)!$  is the number of permutations where the players of  $B$  appear first, then  $a$ , and then the rest of the players.

**2.3. Complexity.** In this article, we focus on the standard notion of *data complexity*, where the relational schema and set of FDs are considered fixed and the input consists of a database and a fact. In particular, a polynomial-time algorithm may be exponential in the number of attributes or FDs. Hence, each combination of a relational schema and an FD set defines a distinct problem, and different combinations may have different computational complexities. We discuss both exact and approximate algorithms for computing Shapley values.

Recall that a *Fully-Polynomial Randomized Approximation Scheme* (FPRAS, for short) for a function  $f$  is a randomized algorithm  $A(x, \epsilon, \delta)$  that returns an  $\epsilon$ -approximation of  $f(x)$  with probability at least  $1 - \delta$ , given an input  $x$  for  $f$  and  $\epsilon, \delta \in (0, 1)$ , in time polynomial in  $x$ ,  $1/\epsilon$ , and  $\log(1/\delta)$ . Formally, an FPRAS, satisfies:

$$\Pr [f(x)/(1 + \epsilon) \leq A(x, \epsilon, \delta) \leq (1 + \epsilon)f(x)] \geq 1 - \delta.$$

Note that this notion of FPRAS refers to a *multiplicative* approximation, and we adopt this notion implicitly unless stated otherwise. We may also write “multiplicative” explicitly for

stress. In cases where the function  $f$  has a bounded range, it also makes sense to discuss an *additive* FPRAS where  $\Pr[f(x) - \epsilon \leq A(x, \epsilon, \delta) \leq f(x) + \epsilon] \geq 1 - \delta$ . We refer to an additive FPRAS, and explicitly state so, in cases where the Shapley value is in the range  $[0, 1]$ .

### 3. THE SHAPLEY VALUE OF INCONSISTENCY MEASURES

In this article, we study the Shapley value of facts with respect to measures of database inconsistency. More precisely, the cooperative game that we consider here is determined by an inconsistency measure  $\mathcal{I}$ , and the facts of the database take the role of the players. In turn, an *inconsistency measure*  $\mathcal{I}$  is a function that maps pairs  $(D, \Delta)$  of a database  $D$  and a set  $\Delta$  of FDs to a number  $\mathcal{I}(D, \Delta) \in [0, \infty)$ . Intuitively, the higher the value  $\mathcal{I}(D, \Delta)$  is, the more inconsistent (or, the less consistent) the database  $D$  is w.r.t.  $\Delta$ . The Shapley value of a fact  $f$  of a database  $D$  w.r.t. an FD set  $\Delta$  and inconsistency measure  $\mathcal{I}$  is then defined as follows.

$$\text{Shapley}(D, \Delta, f, \mathcal{I}) \stackrel{\text{def}}{=} \sum_{E \subseteq (D \setminus \{f\})} \frac{|E|! \cdot (|D| - |E| - 1)!}{|D|!} \left( \mathcal{I}(E \cup \{f\}, \Delta) - \mathcal{I}(E, \Delta) \right) \quad (3.1)$$

We note that the definition of the Shapley value requires the cooperative game to be zero on the empty set [Sha53] and this is indeed the case for all of the inconsistency measures  $\mathcal{I}$  that we consider in this work. Next, we introduce each of these measures.

- $\mathcal{I}_d$  is the *drastic measure* that takes the value 1 if the database is inconsistent and the value 0 otherwise [Thi17].
- $\mathcal{I}_{\text{MI}}$  counts the *minimal inconsistent subsets* [HK08, HK10]; in the case of FDs, these subsets are simply the pairs of tuples that jointly violate an FD.
- $\mathcal{I}_{\text{P}}$  is the number of *problematic facts*, where a fact is problematic if it belongs to a minimal inconsistent subset [GH11]; in the case of FDs, a fact is problematic if and only if it participates in a pair of facts that jointly violate  $\Delta$ .
- $\mathcal{I}_{\text{R}}$  is the minimal number of facts that we need to delete from the database for  $\Delta$  to be satisfied (similarly to the concept of a cardinality repair and proximity in Property Testing) [GH13, GGR98, Ber19].
- $\mathcal{I}_{\text{MC}}$  is the number of *maximal consistent subsets* (i.e., repairs) [GH11, GH17].

Table 1 summarizes the complexity results for the different measures. The first column (lhs chain) refers to FD sets that have a left-hand-side chain—a notion that was introduced by Livshits et al. [LK17], and we recall in the next section. The second column (no lhs chain, PTime c-repair) refers to FD sets that do not have a left-hand-side chain, but entail a polynomial-time cardinality repair computation according to the dichotomy of Livshits et al. [LKR20] that we discuss in more details in Section 6.

**Example 3.1.** Consider again the database of our running example. Since the database is inconsistent w.r.t. the FD set defined in Example 2.1, we have that  $\mathcal{I}_d(D, \Delta) = 1$ . As for the measure  $\mathcal{I}_{\text{MI}}$ , the reader can easily verify that there are twenty nine pairs of tuples that jointly violate the FDs; hence, we have that  $\mathcal{I}_{\text{MI}}(D, \Delta) = 29$ . Since each tuple participates in at least one violation of the FDs, it holds that  $\mathcal{I}_{\text{P}}(D, \Delta) = 9$ . Finally, as we have already seen in Example 2.1, the database has five repairs and a single cardinality repair obtained by deleting six facts. Thus,  $\mathcal{I}_{\text{R}}(D, \Delta) = 6$  and  $\mathcal{I}_{\text{MC}}(D, \Delta) = 5$ . In the next sections, we discuss the computation of the Shapley value for each one of these measures.  $\diamond$

**Preliminary analysis.** We study the *data complexity* of computing  $\text{Shapley}(D, \Delta, f, \mathcal{I})$  for different inconsistency measures  $\mathcal{I}$ . To this end, we give here two important observations that we will use throughout the article. The first observation is that the computation of  $\text{Shapley}(D, \Delta, f, \mathcal{I})$  can be easily reduced to the computation of the expected value of the inconsistency measure over all Here, we denote by  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D' \cup \{f\}, \Delta))$  the expected value of  $\mathcal{I}(D' \cup \{f\}, \Delta)$  over all subsets  $D'$  of  $D \setminus \{f\}$  of a given size  $m$ , assuming a uniform distribution. Similarly,  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D', \Delta))$  is the expected value of  $\mathcal{I}(D', \Delta)$  over all such subsets  $D'$ .

**Observation 3.2.** Let  $\mathcal{I}$  be an inconsistency measure. The following holds.

$$\text{Shapley}(D, \Delta, f, \mathcal{I}) = \frac{1}{|D|} \sum_{m=0}^{|D|-1} [\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D' \cup \{f\}, \Delta)) - \mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D', \Delta))]$$

*Proof.* We have the following.

$$\begin{aligned} \text{Shapley}(D, \Delta, f, \mathcal{I}) &= \sum_{D' \subseteq (D \setminus \{f\})} \frac{|D'|!(|D| - |D'| - 1)!}{|D|!} (\mathcal{I}(D' \cup \{f\}, \Delta) - \mathcal{I}(D', \Delta)) \\ &= \sum_{m=0}^{|D|-1} \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m}} \frac{m!(|D| - m - 1)!}{|D|!} (\mathcal{I}(D' \cup \{f\}, \Delta) - \mathcal{I}(D', \Delta)) \\ &= \sum_{m=0}^{|D|-1} \frac{m!(|D| - m - 1)!}{|D|!} \binom{|D| - 1}{m} \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m}} \frac{1}{\binom{|D|-1}{m}} (\mathcal{I}(D' \cup \{f\}, \Delta)) \end{aligned} \quad (3.2)$$

$$\begin{aligned} &- \sum_{m=0}^{|D|-1} \frac{m!(|D| - m - 1)!}{|D|!} \binom{|D| - 1}{m} \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m}} \frac{1}{\binom{|D|-1}{m}} (\mathcal{I}(D', \Delta)) \\ &= \sum_{m=0}^{|D|-1} \frac{m!(|D| - m - 1)!}{|D|!} \binom{|D| - 1}{m} \mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D' \cup \{f\}, \Delta)) \\ &- \sum_{m=0}^{|D|-1} \frac{m!(|D| - m - 1)!}{|D|!} \binom{|D| - 1}{m} \mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D', \Delta)) \\ &= \frac{1}{|D|} \sum_{m=0}^{|D|-1} [\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D' \cup \{f\}, \Delta)) - \mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}(D', \Delta))] \end{aligned} \quad (3.3)$$

Note that in Equation (3.2) we multiply and divide by the value  $\binom{|D|-1}{m}$ . The expectation expression of Equation (3.3) is due to the fact that  $1/\binom{|D|-1}{m}$  is the probability of a random subset of size  $m$  of  $D \setminus \{f\}$  in the uniform distribution.  $\square$

Observation 3.2 implies that to compute the Shapley value of  $f$ , it suffices to compute the expectations of the amount of inconsistency over subsets  $D'$  and  $D' \cup \{f\}$ , where  $D'$  is drawn uniformly from the space of subsets of size  $m$ , for every  $m$ . More precisely, the

computation of the Shapley value is Cook reducible<sup>1</sup> to the computation of these expectations. Our algorithms will, indeed, compute these expectations instead of the Shapley value.

The second observation is the following. One of the basic properties of the Shapley value is one termed “efficiency”—the sum of the Shapley values over all the players equals the total wealth [Sha53]. This property implies that  $\sum_{f \in D} \text{Shapley}(D, \Delta, f, \mathcal{I}) = \mathcal{I}(D, \Delta)$ . Thus, whenever the measure itself is computationally hard, so is the Shapley value of facts.

**Fact 3.3.** Let  $\mathcal{I}$  be an inconsistency measure. The computation of  $\mathcal{I}$  is Cook reducible to the computation of the Shapley value of facts under  $\mathcal{I}$ .

This observation can be used for showing lower bounds on the complexity of the Shapley value, as we will see in the next sections.

#### 4. MEASURES $\mathcal{I}_{\text{MI}}$ AND $\mathcal{I}_{\text{P}}$ : THE TRACTABLE MEASURES

We start by discussing two tractable measures, namely  $\mathcal{I}_{\text{MI}}$  and  $\mathcal{I}_{\text{P}}$ . We first give algorithms for computing the Shapley value for these measures, and then discuss the generalization to multiple relations.

**4.1. Computation.** Recall that  $\mathcal{I}_{\text{MI}}$  counts the pairs of facts that jointly violate at least one FD. An easy observation is that a fact  $f$  increases the value of the measure  $\mathcal{I}_{\text{MI}}$  by  $i$  in a permutation  $\sigma$  if and only if  $\sigma_f$  contains exactly  $i$  facts that are in conflict with  $f$ . Hence, assuming that  $D$  contains  $N_f$  facts that conflict with  $f$ , the Shapley value for this measure can be computed in the following way:

$$\begin{aligned} \text{Shapley}(D, \Delta, f, \mathcal{I}_{\text{MI}}) &= \sum_{E \subseteq (D \setminus \{f\})} \frac{|E|! \cdot (|D| - |E| - 1)!}{|D|!} \left( \mathcal{I}(E \cup \{f\}, \Delta) - \mathcal{I}(E, \Delta) \right) \\ &= \frac{1}{|D|!} \sum_{i=1}^{N_f} \sum_{\substack{E \subseteq (D \setminus \{f\}) \\ |E \cap N_f| = i}} |E|! \cdot (|D| - |E| - 1)! \cdot i = \frac{1}{|D|!} \sum_{i=1}^{N_f} \sum_{m=i}^{|D|-1} \sum_{\substack{E \subseteq (D \setminus \{f\}) \\ |E|=m \\ |E \cap N_f|=i}} m! \cdot (|D| - m - 1)! \cdot i \\ &= \frac{1}{|D|!} \sum_{i=1}^{N_f} \sum_{m=i}^{|D|-1} \binom{N_f}{i} \binom{|D| - N_f - 1}{m - i} \cdot m! \cdot (|D| - m - 1)! \cdot i \end{aligned}$$

Therefore, we immediately obtain the following result.

**Theorem 4.1.** *Let  $\Delta$  be a set of FDs.  $\text{Shapley}(D, \Delta, f, \mathcal{I}_{\text{MI}})$  is computable in polynomial time, given  $D$  and  $f$ .*

We now move on to  $\mathcal{I}_{\text{P}}$  that counts the “problematic” facts; that is, facts that participate in a violation of  $\Delta$ . Here, a fact  $f$  increases the measure by  $i$  in a permutation  $\sigma$  if and only if  $\sigma_f$  contains precisely  $i - 1$  facts that are in conflict with  $f$ , but not in conflict with any other fact of  $\sigma_f$  (hence, all these facts and  $f$  itself are added to the group of problematic facts). We prove the following.

<sup>1</sup>Recall that a *Cook reduction* from a function  $F$  to a function  $G$  is a polynomial-time *Turing reduction* from  $F$  to  $G$ , that is, an algorithm that computes  $F$  with an oracle to a solver of  $G$ .

**Theorem 4.2.** *Let  $\Delta$  be a set of FDs.  $\text{Shapley}(D, f, \Delta, \mathcal{I}_P)$  is computable in polynomial time, given  $D$  and  $f$ .*

*Proof.* We now show how the expected values of Observation 3.2 can be computed in polynomial time. We start with  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_P(D', \Delta))$ . We consider the uniform distribution  $U_m(D \setminus \{f\})$  over the subsets of size  $m$  of  $D \setminus \{f\}$ . We denote by  $X$  the random variable holding the number of problematic facts in the random subset. We denote by  $Y_g$  the random variable that holds 1 if the fact  $g$  is in the random subset and, moreover, participates there in a violation of the FDs. In addition, we denote the expectations of these variables by  $\mathbb{E}(X)$  and  $\mathbb{E}(Y_g)$ , respectively (without explicitly stating the distribution  $D' \sim U_m(D \setminus \{f\})$  in the subscript). Due to the linearity of the expectation we have:

$$\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_P(D', \Delta)) = \mathbb{E}(X) = \mathbb{E}\left(\sum_{g \in D \setminus \{f\}} Y_g\right) = \sum_{g \in D \setminus \{f\}} \mathbb{E}(Y_g)$$

Hence, the computation of  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_P(D', \Delta))$  reduces to the computation of  $\mathbb{E}(Y_g)$ , and this value can be computed as follows.

$$\begin{aligned} \mathbb{E}(Y_g) &= \Pr[g \text{ is selected}] \times \Pr[\text{a conflicting fact is selected} \mid g \text{ is selected}] \\ &= \frac{\binom{|D|-2}{m-1}}{\binom{|D|-1}{m}} \cdot \frac{\sum_{k=1}^{N_g} \binom{N_g}{k} \cdot \binom{|D|-1-N_g}{m-k-1}}{\binom{|D|-2}{m-1}} = \frac{\sum_{k=1}^{N_g} \binom{N_g}{k} \cdot \binom{|D|-1-N_g}{m-k-1}}{\binom{|D|-1}{m}} \end{aligned}$$

where  $N_g$  is the number of facts in  $D \setminus \{f\}$  that are in conflict with  $g$ .

We can similarly consider the distribution  $U_m(D \setminus \{f\})$  and show that the expectation  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_P(D' \cup \{f\}, \Delta))$  is equal to  $\sum_{g \in D \setminus \{f\}} \mathbb{E}(Y'_g)$ , where  $Y'_g$  is a random variable that holds 1 if  $g$  is selected in the random subset and, moreover, participates in a violation of the FDs, and 0 otherwise. For a fact  $g$  that is not in conflict with  $f$  it holds that  $\mathbb{E}(Y'_g) = \mathbb{E}(Y_g)$ , while for a fact  $g$  that is in conflict with  $f$  it holds that

$$\mathbb{E}(Y'_g) = \Pr[g \text{ is selected}] = \frac{\binom{|D|-2}{m-1}}{\binom{|D|-1}{m}}. \quad \square$$

**4.2. Generalization to Multiple Relations.** The results of this section immediately generalize to schemas with multiple relation symbols. This is true since one of the basic properties of the Shapley value is linearity [Sha53]:

$$\text{Shapley}(D, f, \Delta, a \cdot \alpha + b \cdot \beta) = a \cdot \text{Shapley}(D, f, \Delta, \alpha) + b \cdot \text{Shapley}(D, f, \Delta, \beta)$$

and both measures,  $\mathcal{I}_M$  and  $\mathcal{I}_P$ , are additive over multiple relations, that is, the value of the measure on the entire database is the sum of the values over the individual relations.

## 5. MEASURE $\mathcal{I}_d$ : THE DRASTIC MEASURE

In this section, we consider the drastic measure  $\mathcal{I}_d$ . While the measure itself is extremely simple and, in particular, computable in polynomial time (testing whether  $\Delta$  is satisfied), it might be intractable to compute the Shapley value of a fact. In particular, we prove a dichotomy for this measure, classifying FD sets into ones where the Shapley value can be computed in polynomial time and the rest where the problem is  $\text{FP}^{\#\text{P}}$ -complete.<sup>2</sup>

**5.1. Dichotomy.** Before giving our dichotomy, we recall the definition of a *left-hand-side chain* (lhs chain, for short), introduced by Livshits et al. [LK17].

**Definition 5.1** [LK17]. An FD set  $\Delta$  has a left-hand-side chain if for every two FDs  $X \rightarrow Y$  and  $X' \rightarrow Y'$  in  $\Delta$ , either  $X \subseteq X'$  or  $X' \subseteq X$ .

**Example 5.2.** The FD set of our running example (Example 2.1) has an lhs chain. We could also define  $\Delta$  with redundancy by adding the following FD: train time arrives  $\rightarrow$  departs. The resulting FD set does not have an lhs chain, but it is *equivalent* to an FD set with an lhs chain. An example of an FD set that does not have an lhs chain, not even up to equivalence, is  $\{\text{train time} \rightarrow \text{departs}, \text{train departs} \rightarrow \text{time}\}$ .  $\diamond$

We prove the following.

**Theorem 5.3.** *Let  $\Delta$  be a set of FDs. If  $\Delta$  is equivalent to an FD set with an lhs chain, then  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$  is computable in polynomial time, given  $D$  and  $f$ . Otherwise, the problem is  $\text{FP}^{\#\text{P}}$ -complete.*

Interestingly, this is the exact same dichotomy that we obtained in prior work [LK17] for the problem of counting subset repairs. We also showed that this tractability criterion is decidable in polynomial time by computing a minimal cover: if  $\Delta$  is equivalent to an FD set with an lhs chain, then every minimal cover of  $\Delta$  has an lhs chain. In the remainder of this section, we prove Theorem 5.3.

**5.1.1. Hardness Side.** The proof of the hardness side of Theorem 5.3 has two steps. We first show hardness for the matching constraint  $\{A \rightarrow B, B \rightarrow A\}$  over the schema  $(A, B)$ , and this proof is similar to the proof of Livshits et al. [LBKS20] for the problem of computing the Shapley contribution of facts to the result of the query  $q() :- R(x), S(x, y), T(y)$ . Then, from this case to the remaining cases we apply the *fact-wise reductions* that have been devised in prior work [LK17]. We start by proving hardness for  $\{A \rightarrow B, B \rightarrow A\}$ .

**Lemma 5.4.** *Computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$  for the FD set  $\Delta = \{A \rightarrow B, B \rightarrow A\}$  over the relational schema  $(A, B)$  is  $\text{FP}^{\#\text{P}}$ -complete.*

*Proof.* We construct a reduction from the problem of computing the number  $|\mathbf{M}(g)|$  of matchings in a bipartite graph  $g$  [Val79a]. Note that we consider partial matchings; that is, subsets of edges that consist of mutually-exclusive edges. Given an input bipartite graph  $g$ , we construct  $m + 1$  input instances  $(D_1, f_1), \dots, (D_{m+1}, f_{m+1})$  to our problem, where  $m$  is the number of edges in  $g$ , in the following way. For every  $r \in \{1, \dots, m + 1\}$ , we add one vertex  $v_1$  to the left-hand side of  $g$  and  $r + 1$  vertices  $u_1, \dots, u_r, v_2$  to the right-hand side

<sup>2</sup>Recall that  $\text{FP}^{\#\text{P}}$  is the class of polynomial-time functions with an oracle to a problem in  $\#\text{P}$  (e.g., count the satisfying assignments of a propositional formula).

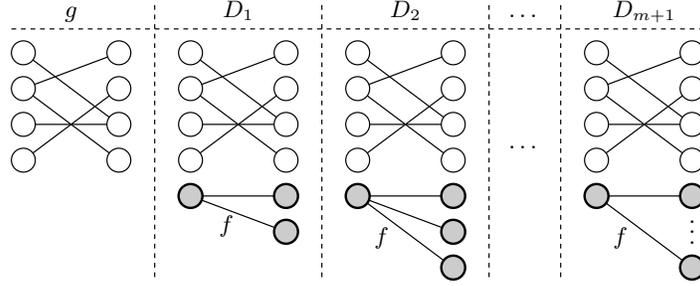


FIGURE 2. The databases constructed in the reduction of the proof of Lemma 5.4.

of  $g$ . Then, we connect the vertex  $v_1$  to every new vertex on the right-hand side of  $g$ . We construct the instance  $D_r$  from the resulting graph by adding a fact  $(u, v)$  for every edge  $(u, v)$  in the graph. We will compute the Shapley value of the fact  $f$  corresponding to the edge  $(v_1, v_2)$ . The reduction is illustrated in Figure 2.

In every instance  $D_r$ , the fact  $f$  will increase the value of the measure by one in a permutation  $\sigma$  if and only if  $\sigma_f$  satisfies two properties: (1) the facts of  $\sigma_f$  jointly satisfy the FDs in  $\Delta$ , and (2)  $\sigma_f$  contains at least one fact that is in conflict with  $f$ . Hence, for  $f$  to affect the value of the measure in a permutation, we have to select a set of facts corresponding to a matching from the original graph  $g$ , as well as exactly one of the facts corresponding to an edge  $(v_1, u_i)$  (since the facts  $(v_1, u_i)$  and  $(v_1, u_j)$  for  $i \neq j$  jointly violate the FD  $A \rightarrow B$ ). We have the following.

$$\text{Shapley}(D_r, f, \Delta, \mathcal{I}_d) = \sum_{k=0}^m |\mathbf{M}(g, k)| \cdot r \cdot (k+1)! \cdot (m-k+r-1)!$$

where  $\mathbf{M}(g, k)$  is the set of matchings of  $g$  containing precisely  $k$  edges.

Hence, we obtain  $m+1$  equations from the  $m+1$  constructed instances, and get the following system of equations.

$$\begin{pmatrix} 1 \cdot 1!m! & 1 \cdot 2!(m-1)! & \dots & 1 \cdot (m+1)!0! \\ 2 \cdot 1!(m+1)! & 2 \cdot 2!m! & \dots & 2 \cdot (m+1)!1! \\ \vdots & \vdots & \vdots & \vdots \\ (m+1) \cdot 1!2m! & (m+1) \cdot 2!(m-1)! & \dots & (m+1) \cdot (m+1)!m! \end{pmatrix} \begin{pmatrix} |\mathbf{M}(g, 0)| \\ |\mathbf{M}(g, 1)| \\ \vdots \\ |\mathbf{M}(g, m)| \end{pmatrix} \\ = \begin{pmatrix} \text{Shapley}(D_1, f, \Delta, \mathcal{I}_d) \\ \text{Shapley}(D_2, f, \Delta, \mathcal{I}_d) \\ \vdots \\ \text{Shapley}(D_{m+1}, f, \Delta, \mathcal{I}_d) \end{pmatrix}$$

Let us divide each column in the above matrix by the constant  $(j+1)!$  (where  $j$  is the column number, starting from 0) and each row by  $i+1$  (where  $i$  is the row number, starting from 0), and reverse the order of the columns. We then get the following matrix.

$$A = \begin{pmatrix} 0! & 1! & \dots & m! \\ 1! & 2! & \dots & (m+1)! \\ \vdots & \vdots & \vdots & \vdots \\ m! & (m+1)! & \dots & 2m! \end{pmatrix}$$

This matrix has coefficients  $a_{i,j} = (i+j)!$ , and the determinant of  $A$  is  $\det(A) = \prod_{i=0}^m i!i! \neq 0$ ; hence, the matrix is non-singular [Bac02]. Since dividing a column by a constant divides the determinant by a constant, and reversing the order of the columns can only change the sign of the determinant, the determinant of the original matrix is not zero as well, and the matrix is non-singular. Therefore, we can solve the system of equations, and compute the value  $\sum_{k=0}^m M(g, k)$ , which is precisely the number of matchings in  $g$ .  $\square$

**Generalization via Fact-Wise Reductions.** Using the concept of a fact-wise reduction [Kim12], we can prove hardness for any FD set that is not equivalent to an FD set with an lhs chain. We first give the formal definition of a fact-wise reduction. Let  $(R, \Delta)$  and  $(R', \Delta')$  be two pairs of a relational schema and an FD set. A *mapping* from  $R$  to  $R'$  is a function  $\mu$  that maps facts over  $R$  to facts over  $R'$ . (We say that  $f$  is a fact *over*  $R$  if  $f$  is a fact of some database  $D$  over  $R$ .) We extend a mapping  $\mu$  to map databases  $D$  over  $R$  to databases over  $R'$  by defining  $\mu(D)$  to be  $\{\mu(f) \mid f \in D\}$ . A *fact-wise reduction* from  $(R, \Delta)$  to  $(R', \Delta')$  is a mapping  $\Pi$  from  $R$  to  $R'$  with the following properties.

- (1)  $\Pi$  is injective; that is, for all facts  $f$  and  $g$  over  $R$ , if  $\Pi(f) = \Pi(g)$  then  $f = g$ .
- (2)  $\Pi$  preserves consistency and inconsistency; that is, for all facts  $f$  and  $g$  over  $R$ ,  $\{f, g\}$  satisfies  $\Delta$  if and only if  $\{\Pi(f), \Pi(g)\}$  satisfies  $\Delta'$ .
- (3)  $\Pi$  is computable in polynomial time.

We have previously shown a fact-wise reduction from  $((A, B), \{A \rightarrow B, B \rightarrow A\})$  to any  $(R, \Delta)$ , where  $\Delta$  is not equivalent to an FD set with an lhs chain [LK17]. Clearly, fact-wise reductions preserve the Shapley value of facts, that is,  $\text{Shapley}(D, f, \mathcal{I}, \Delta) = \text{Shapley}(\Pi(D), \Pi(f), \mathcal{I}, \Delta')$ . It thus follows that there is a polynomial-time reduction from the problem of computing the Shapley value over  $\{A \rightarrow B, B \rightarrow A\}$  to the problem of computing the Shapley value over any  $\Delta$  that has no lhs chain (even up to equivalence), and that concludes our proof of hardness.

**5.1.2. Tractability Side.** For the tractability side of Theorem 5.3, we present a polynomial-time algorithm to compute the Shapley value. As stated in Observation 3.2, the computation of  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$  reduces in polynomial time to the computation of the expected value of the measure over all subsets of the database of a given size  $m$ . In this case it holds that  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_d(D' \cup \{f\}, \Delta))$  and  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_d(D', \Delta))$  are the probabilities that a uniformly chosen  $D' \subseteq D \setminus \{f\}$  of size  $m$  is such that  $(D' \cup \{f\}) \not\models \Delta$  and  $D' \not\models \Delta$ , respectively. Due to the structure of FD sets with an lhs chain, we can compute these probabilities efficiently, as we explain next.

Our main observation is that for an FD  $X \rightarrow Y$ , if we group the facts of  $D$  by  $X$  (i.e., split  $D$  into maximal subsets of facts that agree on the values of all attributes in  $X$ ), then this FD and the FDs that appear later in the chain may be violated only among facts from the same group. Moreover, when we group by  $XY$  (i.e., further split each group of  $X$  into maximal subsets of facts that agree on the values of all attributes in  $Y$ ), facts from different

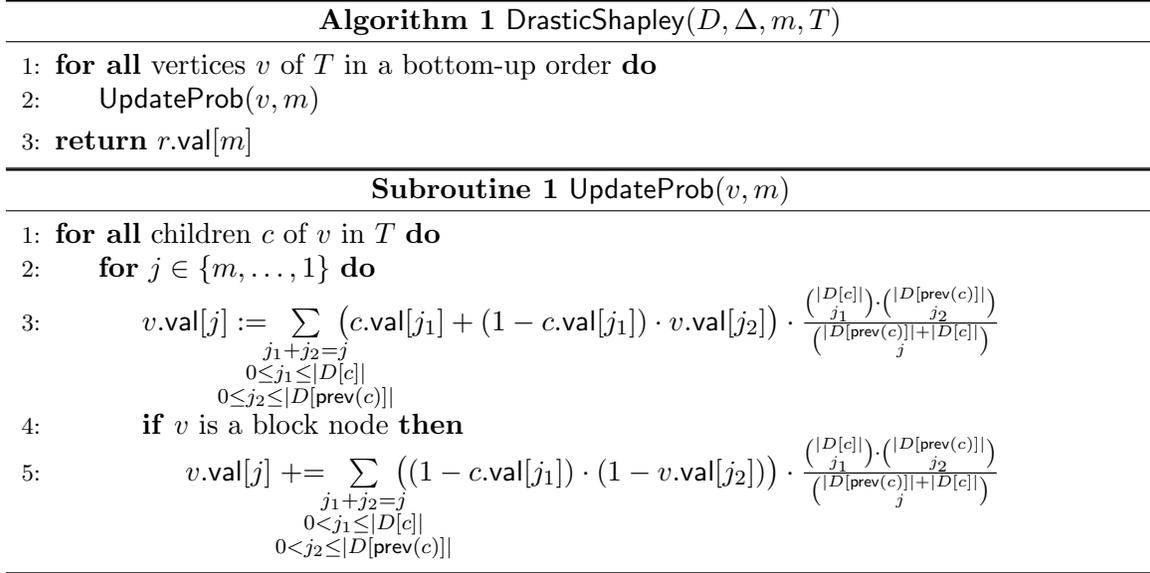


FIGURE 3. An algorithm for computing  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_d(D', \Delta))$  for  $\Delta$  with an lhs chain.

groups always violate this FD, and hence, violate  $\Delta$ . We refer to the former groups as *blocks* and the latter groups as *subblocks*. This special structure allows us to split the problem into smaller problems, solve each one of them separately, and then combine the solutions via dynamic programming.

We define a data structure  $T$  where each vertex  $v$  is associated with a subset of  $D$  that we denote by  $D[v]$ . The root  $r$  is associated with  $D$  itself, that is,  $D[r] = D$ . At the first level, each child  $c$  of  $r$  is associated with a block of  $D[r]$  w.r.t.  $X_1 \rightarrow Y_1$ , and each child  $c'$  of  $c$  is associated with a subblock of  $D[c]$  w.r.t.  $X_1 \rightarrow Y_1$ . At the second level, each child  $c''$  of  $c'$  is associated with a block of  $D[c']$  w.r.t.  $X_2 \rightarrow Y_2$ , and each child  $c'''$  of  $c''$  is associated with a subblock of  $D[c'']$  w.r.t.  $X_2 \rightarrow Y_2$ . This continues all the way to the  $n$ th FD, where at the  $i$ th level, each child  $u$  of an  $(i-1)$ th level subblock vertex  $v$  is associated with a block of  $D[v]$  w.r.t.  $X_i \rightarrow Y_i$  and each child  $u'$  of  $u$  is associated with a subblock of  $D[u]$  w.r.t.  $X_i \rightarrow Y_i$ .

We assume that the data structure  $T$  is constructed in a preprocessing phase. Clearly, the number of vertices in  $T$  is polynomial in  $|D|$  and  $n$  (recall that  $n$  is the number of FDs in  $\Delta$ ) as the height of the tree is  $2n$ , and each level contains at most  $|D|$  vertices; hence, this preprocessing phase requires polynomial time (even under combined complexity). Then, we compute both  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_d(D', \Delta))$  and  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_d(D' \cup \{f\}, \Delta))$  by going over the vertices of  $T$  from bottom to top, as we will explain later. Note that for the computation of these values, we construct  $T$  from the database  $D \setminus \{f\}$ . Figure 4 depicts the data structure  $T$  used for the computation of Shapley( $D, f_9, \Delta, \mathcal{I}_d$ ) for the database  $D$  and fact  $f_9$  of our running example. Next, we explain the meaning of the values stored in each vertex.

Each vertex  $v$  in  $T$  stores an array  $v.\text{val}$  with  $|D[v]| + 1$  entries (that is initialized with zeros) such that  $v.\text{val}[j] = \mathbb{E}_{D' \sim U_j(D[v])}(\mathcal{I}_d(D', \Delta))$  for all  $j \in \{0, \dots, |D[v]|\}$  at the end of

the execution. For this measure, we have that:

$$v.\text{val}[j] \stackrel{\text{def}}{=} \Pr[\text{a random subset of size } j \text{ of } D[v] \text{ violates } \Delta]$$

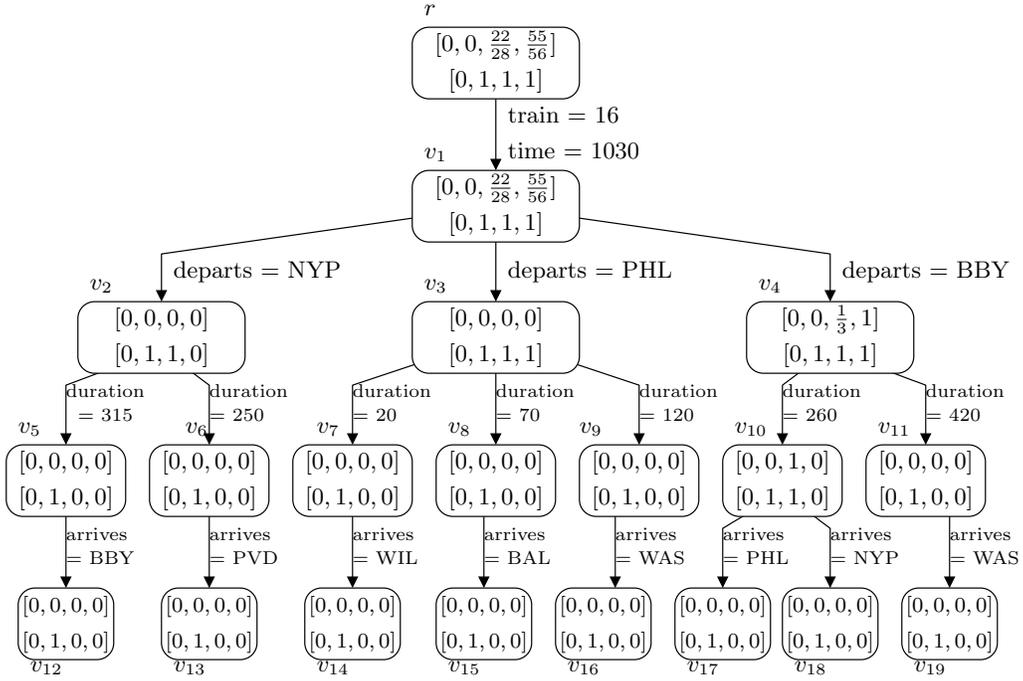
Our final goal is to compute  $r.\text{val}[m]$ , where  $r$  is the root of  $T$ . For that purpose, in the algorithm `DrasticShapley`, depicted in Figure 3, we go over the vertices of  $T$  in a bottom-up order and compute the values of  $v.\text{val}$  for every vertex  $v$  in the `UpdateProb` subroutine. Observe that we only need one execution of `DrasticShapley` with  $m = |D| - 1$  to compute the required values for all  $m \in \{1, \dots, |D| - 1\}$ , as we calculate all these values in our intermediate computations.

To compute  $v.\text{val}$  for a subblock vertex  $v$ , we iterate over its children in  $T$  (which are the  $(i + 1)$ th level blocks) according to an arbitrary order defined in the construction of  $T$ . For a child  $c$  of  $v$ , we denote by  $\text{prev}(c)$  the set of children of  $v$  that occur before  $c$  in that order, and by  $D[\text{prev}(c)]$  the database  $\bigcup_{c' \in \text{prev}(c)} D[c']$ . When considering  $c$  in the for loop of line 1, we compute the expected value of the measure on a subset of  $D[\text{prev}(c)] \cup D[c]$ . Hence, when we consider the last child of  $v$  in the for loop of line 1, we compute the expected value of the measure on a subset of the entire database  $D[v]$ .

For a child  $c$  of  $v$ , there are  $N_1 = \binom{|D[\text{prev}(c)]| + |D[c]|}{j}$  subsets of size  $j$  of all the children of  $v$  considered so far (including  $c$  itself). Each such subset consists of  $j_1$  facts of the current  $c$  (there are  $N_2 = \binom{|D[c]|}{j_1}$  possibilities) and  $j_2$  facts of the previously considered children (there are  $N_3 = \binom{|D[\text{prev}(c)]|}{j_2}$  possibilities), for some  $j_1, j_2$  such that  $j_1 + j_2 = j$ , with probability  $N_2 N_3 / N_1$ . Moreover, such a subset violates  $\Delta$  if either the facts of the current  $c$  violate  $\Delta$  (with probability  $c.\text{val}[j_1]$  that was computed in a previous iteration) or these facts satisfy  $\Delta$ , but the facts of the previous children violate  $\Delta$  (with probability  $(1 - c.\text{val}[j_1]) \cdot v.\text{val}[j_2]$ ). Observe that since we go over the values  $j$  in reverse order in the for loop of line 2 (i.e., from  $m$  to 1), at each iteration of this loop, we have that  $v.\text{val}[j_2]$  (for all considered  $j_2 \leq j$ ) still holds the expected value of  $\mathcal{I}_d$  over subsets of size  $j_2$  of the previous children of  $v$ , which is indeed the value that we need for our computation.

This computation of  $v.\text{val}$  also applies to the block vertices. However, the addition of line 5 only applies to blocks. Since the children of a block belong to different subblocks, and two facts from the same  $i$ th level block but different  $i$ th level subblocks always jointly violate  $X_i \rightarrow Y_i$ , a subset of size  $j$  of a block also violates the constraints if we select a non-empty subset of the current child  $c$  and a non-empty subset of the previous children, even if each of these subsets by itself is consistent w.r.t.  $\Delta$ . Hence, we add this probability in line 5. Note that all the three cases that we consider are disjoint, so we sum the probabilities. Observe also that the leaves of  $T$  have no children and we do not update their probabilities, and, indeed the probability to select a subset from a leaf  $v$  that violates the constraints is zero, as all the facts of  $D[v]$  agree on the values of all the attributes that occur in  $\Delta$ .

**Example 5.5.** We now illustrate the computation of  $\mathbb{E}_{D' \sim U_m(D \setminus \{f_9\})}(\mathcal{I}_d(D', \Delta))$  on the database  $D$  and the fact  $f_9$  of our running example for  $m = 3$ . Inside each node of the data structure  $T$  of Figure 4, we show the values  $[v.\text{val}[0], v.\text{val}[1], v.\text{val}[2], v.\text{val}[3]]$  used for this computation. Below them, we present the corresponding values used in the computation of  $\mathbb{E}_{D' \sim U_m(D \setminus \{f_9\})}(\mathcal{I}_d(D' \cup \{f\}, \Delta))$ . For the leaves  $v$  and each vertex  $v \in \{v_5, \dots, v_9, v_{11}\}$ , we have that  $v.\text{val}[j] = 0$  for every  $j \in \{0, 1, 2, 3\}$ , as  $D[v]$  has a single fact. As for  $v_{10}$ , when we consider its first child  $v_{17}$  in the for loop of line 1 of `UpdateProb`, all the values in  $v_{10}.\text{val}$  remain zero (since  $v_{17}.\text{val}[j_1] = v_{10}.\text{val}[j_2] = 0$  for any  $j_1, j_2$ , and  $|D[\text{prev}(c)]| = 0$ ). However, when we consider its second child  $v_{18}$ , while the computation of line 3 again has

FIGURE 4. The data structure  $T$  of our running example.

no impact on  $v_{10}.\text{val}$ , after the computation of line 5 we have that  $v_{10}.\text{val}[2] = 1$ . And, indeed, there is a single subset of size two of  $D[v_{10}]$ , which is  $\{f_6, f_7\}$ , and it violates the FD  $\text{train time duration} \rightarrow \text{arrives}$ . This also affects the values of  $v_4.\text{val}$ . In particular, when we consider the first child  $v_{10}$  of  $v_4$ , we have that  $v_4.\text{val}[j] = 1$  for  $j = 2$  and  $v_4.\text{val}[j] = 0$  for any other  $j$ . Then, when we consider the second child  $v_{11}$  of  $v_4$ , it holds that  $v_4.\text{val}[2] = \frac{1}{3}$  (as the only subset of size two of  $D[v_4]$  that violates the FDs is  $\{f_6, f_7\}$ , and there are three subsets in total) and  $v_4.\text{val}[3] = 1$  (as every subset of size three contains both  $f_6$  and  $f_7$ ). Finally, we have that  $\mathbb{E}_{D' \sim U_3(D \setminus \{f_9\})}(\mathcal{I}_d(D', \Delta)) = \frac{55}{56}$  and  $\mathbb{E}_{D' \sim U_3(D \setminus \{f_9\})}(\mathcal{I}_d(D' \cup \{f_9\}, \Delta)) = 1$ .  $\diamond$

To compute  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_d(D' \cup \{f\}, \Delta))$ , we use the algorithm `DrasticShapleyF` in Figure 5. There, we distinguish between several types of vertices w.r.t.  $f$ , and show how this expectation can be computed for each one of these types. Before elaborating on the algorithm, we give some non-standard definitions. Recall that all the facts in  $D[v]$ , for an  $i$ th level block vertex  $v$ , agree on the values of all the attributes in  $X_1 Y_1 \dots X_i$ . Moreover, all the facts in  $D[u]$ , for an  $i$ th level subblock vertex  $u$ , agree on the values of all the attributes in  $X_1 Y_1 \dots X_i Y_i$ . We say that  $f$  conflicts with an  $i$ th level block vertex  $v$  if for some  $X_j \rightarrow Y_j$  such that  $j \in \{1, \dots, i-1\}$  it holds that  $f$  agrees with the facts of  $D[v]$  on all the values of the attributes in  $X_j$  but disagrees with them on the attributes of  $Y_j$ . Note that in this case, *every* fact of  $D[v]$  conflicts with  $f$ . Similarly, we say that  $f$  conflicts with an  $i$ th level subblock vertex  $u$  if it violates an FD  $X_j \rightarrow Y_j$  for some  $j \in \{1, \dots, i\}$  with the facts of  $D[u]$ . We also say that  $f$  matches an  $i$ th level block or subblock vertex  $v$  if it agrees with the facts of  $D[v]$  on the values of all the attributes in  $X_1 Y_1 \dots X_i$ .

In `DrasticShapleyF`, we define  $v.\text{val}'[j]$  to be the probability that a random subset of size  $j$  of  $D[v]$  violates  $\Delta$  when  $f$  is added. We first compute  $v.\text{val}$  for all vertices  $v$  of  $T$

---

**Algorithm 2** DrasticShapleyF( $D, \Delta, m, T, f$ )

---

```

1: DrasticShapley( $D, \Delta, m, T$ )
2: for all vertices  $v$  of  $T$  in a bottom-up order do
3:   UpdateProbF( $v, m, f$ )
4: return  $r.\text{val}[m]$ 

```

---

**Subroutine 2** UpdateProbF( $v, m, f$ )

---

```

1: if  $f$  conflicts with  $v$  then
2:    $v.\text{val}'[j] = 1$  for all  $1 \leq j \leq |D[v]|$ 
3:   return
4: if  $f$  does not match  $v$  then
5:    $v.\text{val}'[j] = v.\text{val}[j]$  for all  $1 \leq j \leq m$ 
6:   return
7: for all children  $c$  of  $v$  in  $T$  do
8:   for  $j \in \{m, \dots, 1\}$  do
9:      $v.\text{val}'[j] := \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} (c.\text{val}'[j_1] + (1 - c.\text{val}'[j_1]) \cdot v.\text{val}'[j_2]) \cdot \frac{\binom{|D[c]|}{j_1} \cdot \binom{|D[\text{prev}(c)]|}{j_2}}{\binom{|D[\text{prev}(c)]| + |D[c]|}{j}}$ 
10:    if  $v$  is a block node then
11:       $v.\text{val}'[j] += \sum_{\substack{j_1+j_2=j \\ 0 < j_1 \leq |D[c]| \\ 0 < j_2 \leq |D[\text{prev}(c)]|}} ((1 - c.\text{val}'[j_1]) \cdot (1 - v.\text{val}'[j_2])) \cdot \frac{\binom{|D[c]|}{j_1} \cdot \binom{|D[\text{prev}(c)]|}{j_2}}{\binom{|D[\text{prev}(c)]| + |D[c]|}{j}}$ 

```

---

FIGURE 5. An algorithm for computing  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_d(D' \cup \{f\}, \Delta))$  for  $\Delta$  with an lhs chain.

using DrasticShapley, and then we use these values to compute  $v.\text{val}'$  for some vertices  $v$ . First, we observe that for vertices  $v$  that conflict with  $f$  we have that  $v.\text{val}'[j] = 1$  for every  $1 \leq j \leq |D[v]|$ , as every non-empty subset of  $D[v]$  violates the FDs with  $f$ . Note that this computation also applies to the leaves of  $T$  that are in conflict with  $f$ . For vertices  $v$  that do not conflict with  $f$  but also do not match with  $f$ , we have that  $v.\text{val}'[j] = v.\text{val}[j]$  for every  $1 \leq j \leq m$ , as no fact of  $D[v]$  agrees with  $f$  on the left-hand side of an FD in  $\Delta$  (for  $j = 0$  we clearly have that  $v.\text{val}'[0] = v.\text{val}[0] = 0$ ).

For the rest of the vertices, the arguments given in Section 5 for the computation of  $v.\text{val}$  still hold in this case; hence, the computation of  $v.\text{val}'$  is similar. In particular, for a child  $c$  of  $v$ , a subset  $E$  of size  $j$  of  $D[\text{prev}(c)] \cup D[c]$  is such that  $E \cup \{f\}$  violates  $\Delta$  if either  $(E \cap D[c]) \cup \{f\}$  violates  $\Delta$  or  $(E \cap D[c]) \cup \{f\}$  satisfies  $\Delta$  but  $(E \cap D[\text{prev}(c)]) \cup \{f\}$  violates  $\Delta$ . If  $v$  is a block vertex, then  $E \cup \{f\}$  also violates  $\Delta$  if we choose a non-empty subset from both  $(E \cap D[c])$  and  $(E \cap D[\text{prev}(c)])$ . Therefore, the main difference between the computation of  $v.\text{val}$  in DrasticShapley and the computation of  $v.\text{val}'$  in DrasticShapleyF is the use of the value  $c.\text{val}'$  instead of the value  $c.\text{val}$  in lines 9 and 11.

**5.2. Approximation.** We now consider an approximate computation of the Shapley value. Using the Chernoff-Hoeffding bound, we can easily obtain an additive FPRAS of the value  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$ , by sampling  $O(\log(1/\delta)/\epsilon^2)$  permutations and computing the

average contribution of  $f$  in a permutation. As observed by Livshits et al. [LBKS20], a multiplicative FPRAS can be obtained using the same algorithm (possibly with a different number of samples) if the “gap” property holds: nonzero Shapley values are guaranteed to be large enough compared to the utility value (which is at most 1 in the case of the drastic measure). This is indeed the case here, as we now prove the following gap property of  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$ .

**Proposition 5.6.** *There is a polynomial  $p$  such that for all databases  $D$  and facts  $f$  of  $D$  the value  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$  is either zero or at least  $1/(p(|D|))$ .*

*Proof.* If no fact of  $D$  is in conflict with  $f$ , then  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d) = 0$ . Otherwise, let  $g$  be a fact that violates an FD of  $\Delta$  jointly with  $f$ . Clearly, it holds that  $\{g\} \models \Delta$ , while  $\{g, f\} \not\models \Delta$ . The probability to choose a permutation  $\sigma$ , such that  $\sigma_f$  is exactly  $\{g\}$  is  $\frac{(|D|-2)!}{|D|!} = \frac{1}{|D| \cdot (|D|-1)}$  (recall that  $\sigma_f$  is the set of facts that appear before  $f$  in  $\sigma$ ). Therefore, we have that  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d) \geq \frac{1}{|D| \cdot (|D|-1)}$ , and that concludes our proof.  $\square$

From Proposition 5.6 we conclude that we can obtain an upper bound on the multiplicative error  $\epsilon$  for  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$  by requiring an additive gap of  $\epsilon$  divided by a polynomial. Hence, we get the following.

**Corollary 5.7.**  *$\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$  has both an additive and a multiplicative FPRAS.*

**5.3. Generalization to Multiple Relations.** We now generalize our results to schemas with multiple relation symbols. More formally, we consider (relational) schemas  $\mathbf{S}$  that consists of a finite set  $\{R_1, \dots, R_n\}$  of relation symbols, each associated with a sequence of attributes. For a set  $\Delta$  of FDs over  $\mathbf{S}$  and a relation symbol  $R_j$  of  $\mathbf{S}$ , we denote by  $\Delta_{R_j}$  the restriction of  $\Delta$  to the FDs over  $R_j$ . Similarly, for a database  $D$  over  $\mathbf{S}$ , we denote by  $D_{R_j}$  the restriction of  $D$  to the facts over  $R_j$ . Finally, we denote  $\Delta_{R_1} \cup \dots \cup \Delta_{R_j}$  by  $\Delta^j$  and  $D_{R_1} \cup \dots \cup D_{R_j}$  by  $D^j$ .

It is straightforward that the lower bound provided in this section also holds for schemas with multiple relation symbols. That is, given an FD set  $\Delta$  over a schema  $\mathbf{S}$ , if for at least one relation symbol  $R$  of  $\mathbf{S}$ , the FD set  $\Delta_R$  is not equivalent to an FD set with an lhs chain, then the problem of computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_d)$  is  $\text{FP}^{\#\text{P}}$ -complete. We now generalize our upper bound to schemas with multiple relations; that is, we focus on the case where the FD set  $\Delta_R$  of *every* relation symbol  $R$  of the schema has an lhs chain (up to equivalence), and show that the Shapley value can be computed in polynomial time.

The formula given in Observation 3.2 for computing the Shapley value is general and also applies to databases over schemas with multiple relation symbols. As aforementioned, for the drastic measure, this computation boils down to computing two probabilities—the probability that a uniformly chosen subset of  $D \setminus \{f\}$  of size  $m$  violates the constraints, and the probability that such a subset  $D'$  satisfies  $D' \cup \{f\} \not\models \Delta$ . Since we consider FDs, there are no violations among facts over different relation symbols; hence, we can compute these probabilities separately for each one of the relation symbols (i.e., for every pair  $(D_{R_j}, \Delta_{R_j})$  of a database and its corresponding FD set), and then we combine these results using dynamic programming, as we explain next.

Let  $R_1, \dots, R_n$  be an arbitrary order of the relation symbols. For each  $j \in \{1, \dots, n\}$  we denote by  $T_j^m$  the probability that a uniformly chosen subset of size  $m$  of  $D_{R_j} \setminus \{f\}$  violates  $\Delta_{R_j}$ . This value can be computed in polynomial time for every relation symbol,

using the algorithm of Figure 3, as we assume that  $\Delta_{R_j}$  has an lhs chain. Next, we denote by  $P_j^m$  the probability that a uniformly chosen subset of size  $m$  of  $D^j \setminus \{f\}$  violates the constraints of  $\Delta_{R_1} \cup \dots \cup \Delta_{R_j}$ . Hence, the value  $P_n^m$  is needed for the computation of the Shapley value. We compute this value using dynamic programming. Clearly, we have that:

$$P_1^m = T_1^m$$

and for every  $j > 1$  we prove the following.

**Lemma 5.8.** *For  $j \in \{2, \dots, n\}$  we have that:*

$$P_j^m = \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \left[ \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \right] \times \left( 1 - \left( 1 - T_j^{m_1} \right) \times \left( 1 - P_{j-1}^{m_2} \right) \right)$$

*Proof.* Each subset  $D'$  of size  $m$  of  $D^j \setminus \{f\}$  contains a subset  $E_1$  of size  $m_1$  of  $D_{R_j} \setminus \{f\}$  and a subset  $E_2$  of size  $m_2$  of  $D^{j-1} \setminus \{f\}$ , for some  $m_1, m_2$  such that  $m_1 + m_2 = m$ . Clearly,  $D'$  violates the constraints if and only if at least one of  $E_1$  or  $E_2$  violates the constraints. That is,  $\mathcal{I}_d(D', \Delta^j) = 1$  if either  $\mathcal{I}_d(E_1, \Delta_{R_j}) = 1$  or  $\mathcal{I}_d(E_2, \Delta^{j-1}) = 1$  (or both). Therefore,

$$\mathcal{I}_d(D', \Delta^j) = 1 - \left( 1 - \mathcal{I}_d(E_1, \Delta_{R_j}) \right) \times \left( 1 - \mathcal{I}_d(E_2, \Delta^{j-1}) \right)$$

Then, we have the following:

$$\begin{aligned} P_j^m &= \mathbb{E}_{D' \sim U_m(D^j \setminus \{f\})} (\mathcal{I}_d(D', \Delta^j)) = \sum_{\substack{D' \subseteq D^j \setminus \{f\} \\ |D'| = m}} \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \mathcal{I}_d(D', \Delta^j) \\ &= \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_1| = m_1, |E_2| = m_2}} \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \times \left( 1 - \left( 1 - \mathcal{I}_d(E_1, \Delta_{R_j}) \right) \times \left( 1 - \mathcal{I}_d(E_2, \Delta^{j-1}) \right) \right) \\ &= \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \left[ \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \times \right. \\ &\quad \left. \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_1| = m_1, |E_2| = m_2}} \left[ \frac{1}{\binom{|D_{R_j} \setminus \{f\}|}{m_1}} \times \frac{1}{\binom{|D^{j-1} \setminus \{f\}|}{m_2}} \times \left( 1 - \left( 1 - \mathcal{I}_d(E_1, \Delta_{R_j}) \right) \times \left( 1 - \mathcal{I}_d(E_2, \Delta^{j-1}) \right) \right) \right] \right] \\ &= \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \left[ \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \times \right. \\ &\quad \left. \left( \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_1| = m_1, |E_2| = m_2}} \left[ \frac{1}{\binom{|D_{R_j} \setminus \{f\}|}{m_1}} \times \frac{1}{\binom{|D^{j-1} \setminus \{f\}|}{m_2}} \times 1 \right] \right) \right] \end{aligned}$$

---

**Algorithm 3** Simplify( $\Delta$ )

---

- 1: Remove trivial FDs from  $\Delta$
  - 2: **if**  $\Delta$  is not empty **then**
  - 3:     find a removable pair  $(X, Y)$  of attribute sets
  - 4:      $\Delta := \Delta - XY$
- return**  $\Delta$
- 

FIGURE 6. A simplification algorithm used for deciding whether a cardinality repair w.r.t.  $\Delta$  can be computed in polynomial time [LKR20].

$$\begin{aligned}
& \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_1|=m_1, |E_2|=m_2}} \left[ \frac{1}{\binom{|D_{R_j} \setminus \{f\}|}{m_1}} \times \frac{1}{\binom{|D^{j-1} \setminus \{f\}|}{m_2}} \times (1 - \mathcal{I}_d(E_1, \Delta_{R_j})) \times (1 - \mathcal{I}_d(E_2, \Delta^{j-1})) \right] \\
&= \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \left[ \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \times \right. \\
& \quad \left( 1 - \left( \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ |E_1|=m_1}} \frac{1}{\binom{|D_{R_j} \setminus \{f\}|}{m_1}} \times (1 - \mathcal{I}_d(E_1, \Delta_{R_j})) \right) \times \right. \\
& \quad \left. \left. \left( \sum_{\substack{E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_2|=m_2}} \frac{1}{\binom{|D^{j-1} \setminus \{f\}|}{m_2}} \times (1 - \mathcal{I}_d(E_2, \Delta^{j-1})) \right) \right) \right] \\
&= \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \left[ \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \times (1 - (1 - T_j^{m_1}) \times (1 - P_{j-1}^{m_2})) \right]
\end{aligned}$$

This concludes our proof.  $\square$

We can similarly compute the second probability required for the Shapley value computation. The only difference is that if the fact  $f$  that we consider is over the relation symbol  $R_j$ , then  $T_j^m$  will be the probability that a uniformly chosen  $D' \subset D_{R_j}$  of size  $m$  is such that  $D' \cup \{f\}$  violates  $\Delta_{R_j}$ . This value can be computed in polynomial time using the algorithm of Figure 5. Note that the results of Section 5.2 on the approximate computation of the Shapley value trivially generalize to schemas with multiple relation symbols; hence, there is an additive FPRAS and a multiplicative FPRAS for any set of FDs.

## 6. MEASURE $\mathcal{I}_R$ : THE COST OF A CARDINALITY REPAIR

In this section, we study the measure  $\mathcal{I}_R$  that is based on the cost of a *cardinality repair*, that is, the minimal number of facts that should be deleted from the database in order to obtain a consistent subset. Unlike the other inconsistency measures considered in this article, we do not have a full dichotomy for the measure  $\mathcal{I}_R$ .

**6.1. Complexity Results.** Livshits et al. [LKR20] established a dichotomy for the problem of computing a cardinality repair, classifying FD sets into those for which the problem is solvable in polynomial time, and those for which it is NP-hard. They presented a polynomial-time algorithm, which we refer to as **Simplify**, that takes as input an FD set  $\Delta$ , finds a *removable* pair  $(X, Y)$  of attribute sets (if such a pair exists), and removes every attribute of  $X \cup Y$  from every FD in  $\Delta$  (we denote the result by  $\Delta - XY$ ). A pair  $(X, Y)$  of attribute sets is considered removable if it satisfies the following three conditions:

- $\text{Closure}_\Delta(X) = \text{Closure}_\Delta(Y)$ ,
- $XY$  is nonempty,
- every FD in  $\Delta$  contains either  $X$  or  $Y$  on the left-hand side.

Note that it may be the case that  $X = Y$ , and then the conditions imply that every FD of  $\Delta$  contains  $X$  on the left-hand side. The algorithm is depicted in Figure 6.

Livshits et al. [LKR20] have shown that if it is possible to transform  $\Delta$  to an empty set by repeatedly applying **Simplify**( $\Delta$ ), then a cardinality repair can be computed in polynomial time. Otherwise, the problem is NP-hard (and, in fact, APX-complete).

Fact 3.3 implies that computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  is hard whenever computing  $\mathcal{I}_R(D, \Delta)$  is hard. Hence, we immediately obtain the following.

**Theorem 6.1.** *Let  $\Delta$  be a set of FDs. If  $\Delta$  cannot be emptied by repeatedly applying **Simplify**( $\Delta$ ), then computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  is NP-hard.*

In the remainder of this section, we focus on the tractable cases of the dichotomy of Livshits et al. [LKR20]. In particular, we start by proving that the Shapley value can again be computed in polynomial time for an FD set that has an lhs chain. Note that FD sets with an lhs chain are a special case of FD sets that can be emptied via **Simplify** steps. This holds since every FD set with an lhs chain has either an FD of the form  $\emptyset \rightarrow X$  or a set  $X$  of attributes that occurs on the left-hand side of every FD. In the first case,  $(\emptyset, X)$  is a removable pair, while in the second case,  $(X, X)$  is a removable pair.

**Theorem 6.2.** *Let  $\Delta$  be a set of FDs. If  $\Delta$  is equivalent to an FD set with an lhs chain, then computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  can be done in polynomial time, given  $D$  and  $f$ .*

Our polynomial-time algorithm **RShapley**, depicted in Figure 7, is very similar in structure to **DrasticShapley**. However, to compute the expected value of  $\mathcal{I}_R$ , we take the reduction of Observation 3.2 a step further, and show, that the problem of computing the expected value of the measure over subsets of size  $m$  can be reduced to the problem of computing the number of subsets of size  $m$  of  $D$  that have a cardinality repair of cost  $k$ , given  $m$  and  $k$ . Recall that we refer to the number of facts that are removed from  $D$  to obtain a cardinality repair  $E$  as the *cost* of  $E$ . In the subroutine **UpdateCount**, we compute this number. In what follows, we denote by  $\text{MR}(D, \Delta)$  the cost of a cardinality repair of  $D$  w.r.t.  $\Delta$ .

**Lemma 6.3.** *The following holds.*

$$\text{Shapley}(D, f, \Delta, \mathcal{I}_R) = \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{k=0}^m \frac{k}{\binom{|D|-1}{m}} |S_{m,k}^f| - |S_{m,k}|$$

where:

$$S_{m,k} = \{D' \subseteq D \setminus \{f\} \mid |D'| = m, \text{MR}(D' \cup \{f\}, \Delta) = k\}$$

$$S_{m,k}^f = \{D' \subseteq D \setminus \{f\} \mid |D'| = m, \text{MR}(D', \Delta) = k\}$$

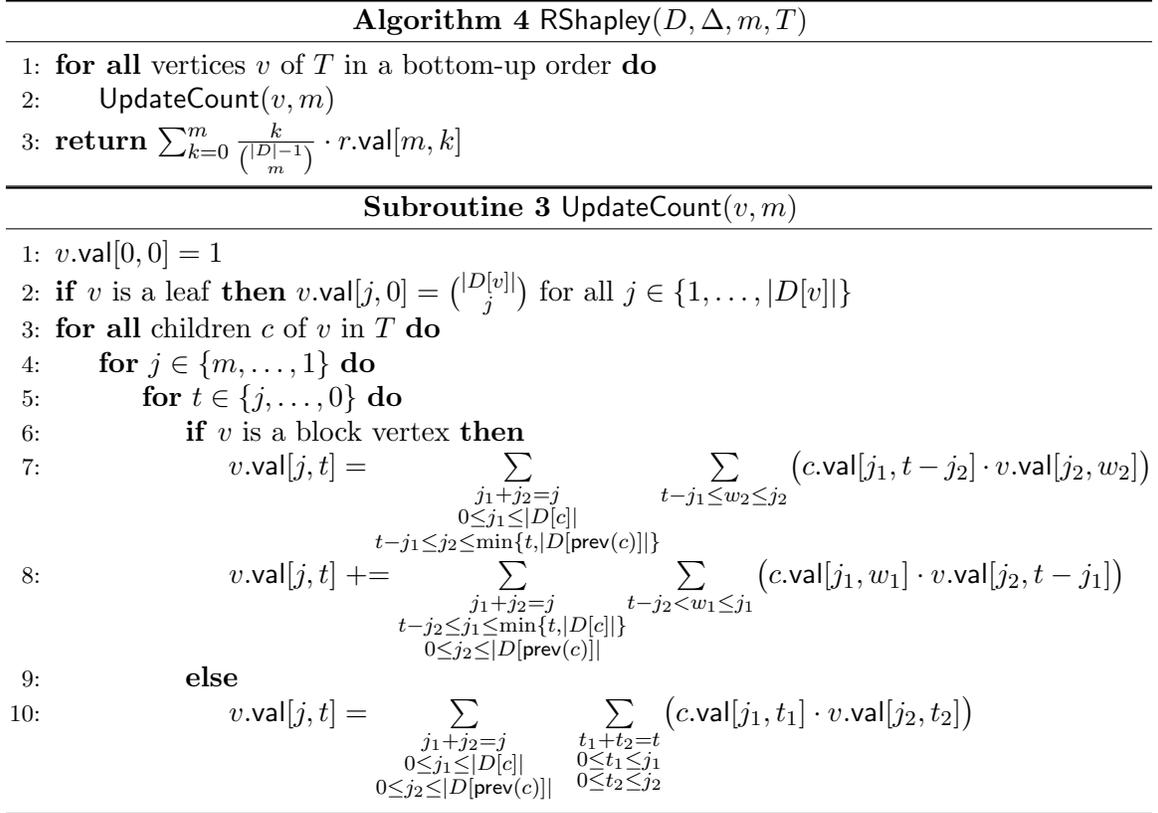


FIGURE 7. An algorithm for computing  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_R(D', \Delta))$  for  $\Delta$  with an lhs chain.

*Proof.* We further develop the reduction of Observation 3.2.

$$\begin{aligned}
\text{Shapley}(D, f, \Delta, \mathcal{I}_R) &= \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m}} \frac{1}{\binom{|D|-1}{m}} \left( \mathcal{I}(D' \cup \{f\}, \Delta) - \mathcal{I}(D', \Delta) \right) \\
&= \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{k=0}^m \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m \\ \text{MR}(D' \cup \{f\}, \Delta) = k}} \frac{1}{\binom{|D|-1}{m}} \left( \mathcal{I}(D' \cup \{f\}, \Delta) \right) \\
&\quad - \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{k=0}^m \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m \\ \text{MR}(D', \Delta) = k}} \frac{1}{\binom{|D|-1}{m}} \left( \mathcal{I}(D', \Delta) \right) \\
&= \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{k=0}^m \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m \\ \text{MR}(D' \cup \{f\}, \Delta) = k}} \frac{k}{\binom{|D|-1}{m}} - \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{k=0}^m \sum_{\substack{D' \subseteq (D \setminus \{f\}) \\ |D'|=m \\ \text{MR}(D', \Delta) = k}} \frac{k}{\binom{|D|-1}{m}}
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{k=0}^m \frac{k}{\binom{|D|-1}{m}} |\{D' \subseteq D \setminus \{f\} \mid |D'| = m, \text{MR}(D' \cup \{f\}, \Delta) = k\}| \\
&- \frac{1}{|D|} \sum_{m=0}^{|D|-1} \sum_{k=0}^m \frac{k}{\binom{|D|-1}{m}} |\{D' \subseteq D \setminus \{f\} \mid |D'| = m, \text{MR}(D', \Delta) = k\}| \quad \square
\end{aligned}$$

We again use the data structure  $T$  defined in the previous section. For each vertex  $v$  in  $T$ , we define:

$$v.\text{val}[j, t] \stackrel{\text{def}}{=} \text{number of subsets of size } j \text{ of } D[v] \text{ with a cardinality repair of cost } t$$

For the leaves  $v$  of  $T$ , we set  $v.\text{val}[j, 0] = \binom{|D[v]|}{j}$  for  $0 \leq j \leq |D[v]|$ , as every subset of  $D[v]$  is consistent, and the cost of a cardinality repair is zero. We also set  $v.\text{val}[0, 0] = 1$  for each  $v$  in  $T$  for the same reason. Since the size of the cardinality repair is bounded by the size of the database, in  $\text{UpdateCount}(v, m)$ , we compute the value  $v.\text{val}[j, t]$  for every  $1 \leq j \leq m$  and  $0 \leq t \leq j$ . To compute this number, we again go over the children of  $v$ , one by one. When we consider a child  $c$  in the for loop of line 1, the value  $v.\text{val}[j, t]$  is the number of subsets of size  $j$  of  $D[\text{prev}(c)] \cup D[c]$  that have a cardinality repair of cost  $t$ .

The children of a block  $v$  are subblocks that jointly violate an FD of  $\Delta$ ; hence, when we consider a child  $c$  of  $v$ , a cardinality repair of a subset  $E$  of  $D[\text{prev}(c)] \cup D[c]$  is either a cardinality repair of  $E \cap D[c]$  (in which case we remove every fact of  $E \cap D[\text{prev}(c)]$ ) or a cardinality repair of  $E \cap D[\text{prev}(c)]$  (in which case we remove every fact of  $E \cap D[c]$ ). The decision regarding which of these cases holds is based on the following four parameters: (1) the number  $j_1$  of facts in  $E \cap D[c]$ , (2) the number  $j_2$  of facts in  $E \cap D[\text{prev}(c)]$ , (3) the cost  $w_1$  of a cardinality repair of  $E \cap D[c]$ , and (4) the cost  $w_2$  of a cardinality repair of  $E \cap D[\text{prev}(c)]$ . In particular:

- If  $w_1 + j_2 \leq w_2 + j_1$ , then a cardinality repair of  $E \cap D[c]$  is preferred over a cardinality repair of  $E \cap D[\text{prev}(c)]$ , as it requires removing less facts from the database.
- If  $w_1 + j_2 > w_2 + j_1$ , then a cardinality repair of  $E \cap D[\text{prev}(c)]$  is preferred over a cardinality repair of  $E \cap D[c]$ .

In fact, since we fix  $t$  in the computation of  $v.\text{val}[j, t]$ , we do not need to go over all  $w_1$  and  $w_2$ . In the first case, we have that  $w_1 = t - j_2$  (hence, the total number of removed facts is  $t - j_2 + j_2 = t$ ), and in the second case we have that  $w_2 = t - j_1$  for the same reason. Hence, in line 7 we consider the first case where  $t \leq w_2 + j_1$ , and in line 8 we consider the second case where  $w_1 + j_2 > t$ . To avoid negative costs, we add a lower bound of  $t - j_1$  on  $j_2$  and  $w_2$  in line 7, and, similarly, a lower bound of  $t - j_2$  on  $j_1$  and  $w_1$  in line 8.

For a subblock vertex  $v$ , a cardinality repair of  $D[v]$  is the union of cardinality repairs of the children of  $v$ , as facts corresponding to different children of  $v$  do not jointly violate any FD. Therefore, for such vertices, in line 10, we compute  $v.\text{val}$  by going over all  $j_1, j_2$  such that  $j_1 + j_2 = j$  and all  $t_1, t_2$  such that  $t_1 + t_2 = t$  and multiply the number of subsets of size  $j_1$  of the current child for which the cost of a cardinality repair is  $t_1$  by the number of subsets of size  $j_2$  of the previously considered children for which the cost of a cardinality repair is  $t_2$ .

Next, we give the algorithm  $\text{RShapleyF}$  for computing  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_R(D' \cup \{f\}, \Delta))$ , that again involves a special treatment for vertices that conflict with  $f$ . The algorithm is

---

**Algorithm 5** RShapleyF( $D, \Delta, m, T, f$ )

---

```

1: RShapley( $D, \Delta, m, T$ )
2: for all vertices  $v$  of  $T$  in a bottom-up order do
3:   UpdateCount( $v, m, f$ )
4: return  $\sum_{k=0}^m \frac{k}{\binom{|D|-1}{m}} \cdot r.\text{val}[m, k]$ 

```

---

**Subroutine 4** UpdateCountF( $v, m, f$ )

---

```

1:  $v.\text{val}'[0, 0] = 1$ 
2: if  $f$  conflict with  $v$  then
3:    $v.\text{val}'[j, t] = v.\text{val}[j, t - 1]$  for all  $1 \leq j \leq |D[v]|$  and  $1 \leq t \leq j$ 
4:   return
5: if  $f$  does not match  $v$  or  $v$  is a leaf then
6:    $v.\text{val}'[j, t] = v.\text{val}[j, t]$  for all  $1 \leq j \leq m$  and  $0 \leq t \leq j$ 
7:   return
8: for all children  $c$  of  $v$  in  $T$  do
9:   for  $j \in \{m, \dots, 1\}$  do
10:    for  $t \in \{j, \dots, 1\}$  do
11:      if  $v$  is a block vertex then
12:         $v.\text{val}'[j, t] = \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ t-j_1 \leq j_2 \leq \min\{t, |D[\text{prev}(c)]\}}} \sum_{t-j_1 \leq w_2 \leq j_2} (c.\text{val}'[j_1, t-j_2] \cdot v.\text{val}'[j_2, w_2])$ 
13:         $v.\text{val}'[j, t] += \sum_{\substack{j_1+j_2=j \\ t-j_2 \leq j_1 \leq \min\{t, |D[c]|\} \\ 0 \leq j_2 \leq |D[\text{prev}(c)]}} \sum_{t-j_2 < w_1 \leq j_1} (c.\text{val}'[j_1, w_1] \cdot v.\text{val}'[j_2, t-j_1])$ 
14:      else
15:         $v.\text{val}'[j, t] = \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]}} \sum_{\substack{t_1+t_2=t \\ 0 \leq t_1 \leq j_1 \\ 0 \leq t_2 \leq j_2}} (c.\text{val}'[j_1, t_1] \cdot v.\text{val}'[j_2, t_2])$ 

```

---

FIGURE 8. An algorithm for computing  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_R(D' \cup \{f\}, \Delta))$  for  $\Delta$  with an lhs chain.

depicted in Figure 8. We define:

$$v.\text{val}'[j, t] \stackrel{\text{def}}{=} \text{number of subsets of size } j \text{ of } D[v] \text{ that, jointly with } f, \\ \text{have a cardinality repair of cost } t$$

As in the case of DrasticShapleyF, we start with the execution of RShapley, which allows us to reuse some of the values computed in this execution. For every vertex, we set  $v.\text{val}'[0, 0] = 1$ , as the empty set has a single cardinality repair of cost zero. Then, we consider three types of vertices. For vertices  $v$  that conflict with  $f$  we have that  $v.\text{val}'[j, t] = v.\text{val}[j, t - 1]$  for all  $1 \leq j \leq |D[v]|$  and  $1 \leq t \leq j$ , as every non-empty subset of  $v$  conflicts with  $f$ ; hence, we have to remove  $f$  in a cardinality repair, and the cost of a cardinality repair increases by one. For vertices  $v$  that do not match  $f$ , we have that  $v.\text{val}'[j, t] = v.\text{val}[j, t]$ , as  $f$  is not in conflict with any fact of  $D[v]$ ; hence, it can be added to any cardinality repair without increasing its cost. The same holds for the leaves of  $T$  that do not conflict with  $f$ .

For a block vertex  $v$ , all the arguments given for RShapley still apply here. In particular, for a child  $c$  of  $v$ , a cardinality repair of  $E \cup \{f\}$  for a subset  $E$  of size  $j$  of  $D[\text{prev}(c)] \cup D[c]$ ,

is either a cardinality repair of  $(E \cap D[c]) \cup \{f\}$  (in which case we delete all facts of  $E \cap D[\text{prev}(c)]$ ) or a cardinality repair of  $(E \cap D[\text{prev}(c)]) \cup \{f\}$  (in which case we delete all facts of  $E \cap D[c]$ ). Therefore, the only difference in the computation of  $v.\text{val}'$  compared to the computation of  $v.\text{val}$  for such vertices is the use of  $c.\text{val}'$  (that takes  $f$  into account) rather than  $c.\text{val}$ .

For a subblock vertex  $v$  (that does not conflict with  $f$ , and, hence, matches  $f$ ), the computation of  $v.\text{val}'$  is again very similar to that of  $v.\text{val}$ , with the only difference being the use of  $c.\text{val}'$ . Observe that in this case, the children of  $v$  correspond to different blocks. Each such block that does not match  $f$  also does not violate any FD with  $f$ ; hence, when we add  $f$  to this block, a cardinality repair of the resulting group of facts does not require the removal of  $f$ . The only child of  $v$  where a cardinality repair might require the removal of  $f$  is a child that matches  $f$ , and, clearly, there is at most one such child. Therefore, we do not count the fact  $f$  twice in the computation of the value  $v.\text{val}'$ .

**6.2. Approximation.** In cases where a cardinality repair can be computed in polynomial time, we can obtain an additive FPRAS in the same way as the drastic measure. (Note that this Shapley value is also in  $[0, 1]$ .) Moreover, we can again obtain a multiplicative FPRAS using the same technique due to the following gap property (proved similarly to Proposition 5.6).

**Proposition 6.4.** *There is a polynomial  $p$  such that for all databases  $D$  and facts  $f$  of  $D$  the value  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  is either zero or at least  $1/(p(|D|))$ .*

As aforementioned, Livshits et al. [LKR20] showed that the hard cases of their dichotomy for the problem of computing a cardinality repair are, in fact, APX-complete; hence, there is a polynomial-time constant-ratio approximation, but for some  $\epsilon > 1$  there is no (randomized)  $\epsilon$ -approximation or else  $P = NP$  ( $NP \subseteq BPP$ ). Since the Shapley value of every fact w.r.t.  $\mathcal{I}_R$  is positive, the existence of a multiplicative FPRAS for  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  would imply the existence of a multiplicative FPRAS for  $\mathcal{I}_R(D, \Delta)$  (due to Fact 3.3), which is a contradiction to the APX-hardness. We conclude the following.

**Proposition 6.5.** *Let  $\Delta$  be a set of FDs. If  $\Delta$  can be emptied by repeatedly applying  $\text{Simplify}(\Delta)$ , then  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  has both an additive and a multiplicative FPRAS. Otherwise, it has neither multiplicative nor additive FPRAS, unless  $NP \subseteq BPP$ .*

**Unsolved cases for  $\mathcal{I}_R$ .** A basic open problem is the computation of  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  for  $\Delta = \{A \rightarrow B, B \rightarrow A\}$ . On the one hand, Proposition 6.5 shows that this case belongs to the tractable side if an approximation is allowed. On the other hand, our algorithm for exact  $\text{Shapley}(D, f, \Delta, \mathcal{I}_R)$  is via counting the subsets of size  $m$  that have a cardinality repair of cost  $k$ . This approach will not work here:

**Proposition 6.6.** *Let  $\Delta = \{A \rightarrow B, B \rightarrow A\}$  be an FD set over  $(A, B)$ . Counting the subsets of size  $m$  of a given database that have a cardinality repair of cost  $k$  is  $\#P$ -hard.*

*Proof.* The proof is by a reduction from the problem of computing the number of perfect matchings in a bipartite graph, known to be  $\#P$ -complete [Val79b]. Given a bipartite graph  $g = (A \cup B, E)$  (where  $|A| = |B|$ ), we construct a database  $D$  over  $(A, B)$  by adding a fact  $(a, b)$  for every edge  $(a, b) \in E$ . We then define  $m = |A|$  and  $k = 0$ . It is rather straightforward that the perfect matchings of  $g$  correspond exactly to the subsets  $D'$  of size  $|A|$  of  $D$  such that  $D'$  itself is a cardinality repair.  $\square$

Observe that the cooperative game for  $\Delta = \{A \rightarrow B, B \rightarrow A\}$  can be seen as a game on bipartite graphs where the vertices on the left-hand side represent the values of attribute  $A$ , the vertices on the right-hand side correspond to the values that occur in attribute  $B$ , and the edges represent the tuples of the database (hence, the players of the game). This game is different from the well-known matching game [AdK14] where the players are the *vertices* of the graph (and the value of the game is determined by the maximum weight matching of the subgraph induced by the coalition). In contrast, in our case the players correspond to the *edges* of the graph. It is not clear what is the connection between the two games and whether or how we can use known results on matching games to derive results for the game that we consider here.

**6.3. Generalization to Multiple Relations.** As in the case of  $\mathcal{I}_{MI}$  and  $\mathcal{I}_P$ , the results of this section easily generalize to schemas with multiple relations, due to the linearity property of the Shapley value. As in the case of the drastic measure, the (positive and negative) results on the approximate computation of the Shapley value trivially generalize to schemas with multiple relation symbols.

## 7. MEASURE $\mathcal{I}_{MC}$ : THE NUMBER OF REPAIRS

The final measure that we consider is  $\mathcal{I}_{MC}$  that counts the repairs of the database.

**7.1. Dichotomy.** A dichotomy result from our previous work [LK17] states that the problem of counting repairs can be solved in polynomial time for FD sets with an lhs chain (up to equivalence), and is  $\#P$ -complete for any other FD set. The hardness side, along with Fact 3.3, implies that computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_{MC})$  is  $\text{FP}^{\#P}$ -hard whenever the FD set is not equivalent to an FD set with an lhs chain. Hence, an lhs chain is a necessary condition for tractability. We show here that it is also sufficient: if the FD set has an lhs chain, then the problem can be solved in polynomial time. Consequently, we obtain the following dichotomy.

**Theorem 7.1.** *Let  $\Delta$  be a set of FDs. If  $\Delta$  is equivalent to an FD set with an lhs chain, then computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_{MC})$  can be done in polynomial time, given  $D$  and  $f$ . Otherwise, the problem is  $\text{FP}^{\#P}$ -complete.*

The algorithm `MCSHapley`, depicted in Figure 9, for computing  $\text{Shapley}(D, f, \Delta, \mathcal{I}_{MC})$ , has the same structure as `DrasticShapley`, with the only difference being the computations in the subroutine `UpdateExpected` (that replaces `UpdateProb`).

For a vertex  $v$  in  $T$  we define:

$$v.\text{val}[j] = \mathbb{E}[\text{number of repairs of a random subset of size } j \text{ of } D[v]]$$

As the number of repairs of a consistent database  $D$  is one ( $D$  itself is a repair), we set  $v.\text{val}[0] = 1$  for every vertex  $v$  and  $v.\text{val}[j] = 1$  for  $0 \leq j \leq |D[v]|$  for every leaf  $v$ . Now, consider a block vertex  $v$  and a child  $c$  of  $v$ . Since the children of  $v$  are subblocks, each repair consists of facts of a single child. Hence, the total number of repairs is the sum of repairs of the children of  $v$ .

Using standard mathematical manipulations, we obtain the following result:

---

**Algorithm 6** MCShapley( $D, \Delta, m, T$ )

---

```

1: for all vertices  $v$  of  $T$  in a bottom-up order do
2:   UpdateExpected( $v, m$ )
3: return  $r.\text{val}[m]$ 

```

---

**Subroutine 5** UpdateExpected( $v, m$ )

---

```

1:  $v.\text{val}[0] = 1$ 
2: if  $v$  is a leaf then  $v.\text{val}[j] = 1$  for all  $j \in \{1, \dots, |D[v]|\}$ 
3: for all children  $c$  of  $v$  in  $T$  do
4:   for  $j \in \{m, \dots, 1\}$  do
5:     if  $v$  is a block vertex then
6:        $v.\text{val}[j] = \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} (c.\text{val}[j_1] + v.\text{val}[j_2])$ 
7:     else
8:        $v.\text{val}[j] = \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} (c.\text{val}[j_1] \cdot v.\text{val}[j_2])$ 

```

---

FIGURE 9. An algorithm for computing  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D', \Delta))$  for  $\Delta$  with an lhs chain.

**Lemma 7.2.** *For a block vertex  $v$  and a child  $c$  of  $v$ , we have that:*

$$\begin{aligned} & \mathbb{E}_{D' \sim U_j(D[\text{prev}(c)] \cup D[c])}(\mathcal{I}_{\text{MC}}(D', \Delta)) \\ &= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} \mathbb{E}_{D' \sim U_{j_1}(D[c])}(\mathcal{I}_{\text{MC}}(D', \Delta)) + \mathbb{E}_{D' \sim U_{j_2}(D[\text{prev}(c)])}(\mathcal{I}_{\text{MC}}(D', \Delta)) \end{aligned}$$

*Proof.* As aforementioned, each repair of a subset  $E$  of  $D[v]$  contains facts from a single child of  $v$ , and the number of repairs is the sum of repairs over the children of  $v$ . Moreover, since our choice of facts from different subblocks is independent, we have the following (where  $\text{MC}(D, \Delta)$  is the set of repairs of  $D$  w.r.t.  $\Delta$ ).

$$\begin{aligned} & \mathbb{E}_{D' \sim U_j(D[\text{prev}(c)] \cup D[c])}(\mathcal{I}_{\text{MC}}(D', \Delta)) = \sum_{\substack{D' \subseteq D[\text{prev}(c)] \cup D[c] \\ |D'|=j}} \Pr[D'] \cdot |\text{MC}(D', \Delta)| \\ &= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_1] \Pr[E_2] (|\text{MC}(E_1, \Delta)| + |\text{MC}(E_2, \Delta)|) \\ &= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_1] \Pr[E_2] |\text{MC}(E_1, \Delta)| \\ &+ \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_1] \Pr[E_2] |\text{MC}(E_2, \Delta)| \end{aligned}$$

$$\begin{aligned}
&= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \Pr[E_1] |\text{MC}(E_1, \Delta)| \sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_2] \\
&+ \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_2] |\text{MC}(E_2, \Delta)| \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \Pr[E_1] \\
&= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \mathbb{E}_{D' \sim U_{j_1}(D[c])} (\mathcal{I}_{\text{MC}}(D', \Delta)) + \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \mathbb{E}_{D' \sim U_{j_2}(D[\text{prev}(c)])} (\mathcal{I}_{\text{MC}}(D', \Delta)) \\
&= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \mathbb{E}_{D' \sim U_{j_1}(D[c])} (\mathcal{I}_{\text{MC}}(D', \Delta)) + \mathbb{E}_{D' \sim U_{j_2}(D[\text{prev}(c)])} (\mathcal{I}_{\text{MC}}(D', \Delta))
\end{aligned}$$

Recall that in our reduction from the problem of computing the Shapley value to that of computing the expected value of the measure over subsets of a given size of the database, we considered the uniform distribution where  $\Pr[E] = \frac{1}{\binom{|D|}{m}}$  for a subset  $E$  of size  $m$  of  $D$ . Therefore, we have that  $\sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_2] = \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \Pr[E_1] = 1$ .  $\square$

The result of Lemma 7.2 is reflected in line 6 of the `UpdateExpected` subroutine. Next, we show the following result for subblock vertices, that we use for the calculation of line 8.

**Lemma 7.3.** *For a subblock vertex  $v$  and a child  $c$  of  $v$ , we have that:*

$$\begin{aligned}
&\mathbb{E}_{D' \sim U_j(D[\text{prev}(c)] \cup D[c])} (\mathcal{I}_{\text{MC}}(D', \Delta)) = \\
&\sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \mathbb{E}_{D' \sim U_{j_1}(D[c])} (\mathcal{I}_{\text{MC}}(D', \Delta)) \cdot \mathbb{E}_{D' \sim U_{j_2}(D[\text{prev}(c)])} (\mathcal{I}_{\text{MC}}(D', \Delta))
\end{aligned}$$

*Proof.* Since the children of  $v$  are blocks (that do not jointly violate any FD of  $\Delta$ ), each repair of a subset  $E$  of  $D[v]$  is a union of the repairs of the children of  $v$ , and the number of repairs is the product of the number of repairs over the children of  $v$ . Hence, we have the following:

$$\begin{aligned}
&\mathbb{E}_{D' \sim U_j(D[\text{prev}(c)] \cup D[c])} (\mathcal{I}_{\text{MC}}(D', \Delta)) = \sum_{\substack{D' \subseteq D[\text{prev}(c)] \cup D[c] \\ |D'|=j}} \Pr[D'] \cdot |\text{MC}(D', \Delta)| \\
&= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_1] \Pr[E_2] (|\text{MC}(E_1, \Delta)| \cdot |\text{MC}(E_2, \Delta)|) \\
&= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq D[\text{prev}(c)]}} \sum_{\substack{E_1 \subseteq D[c] \\ |E_1|=j_1}} \Pr[E_1] |\text{MC}(E_1, \Delta)| \sum_{\substack{E_2 \subseteq D[\text{prev}(c)] \\ |E_2|=j_2}} \Pr[E_2] |\text{MC}(E_2, \Delta)|
\end{aligned}$$

---

**Algorithm 7** MCSShapleyF( $D, \Delta, m, T, f$ )

---

```

1: MCSShapley( $D, \Delta, m, T$ )
2: for all vertices  $v$  of  $T$  in a bottom-up order do
3:   UpdateExpectedF( $v, m, f$ )
4: return  $r.\text{val}[m]$ 

```

---

**Subroutine 6** UpdateExpectedF( $v, m, f$ )

---

```

1:  $v.\text{val}'[0] = 1$ 
2: if  $f$  conflict with  $v$  then
3:    $v.\text{val}'[j] = v.\text{val}[j] + 1$  for all  $1 \leq j \leq |D[v]|$ 
4:   return
5: if  $f$  does not match  $v$  or  $v$  is a leaf then
6:    $v.\text{val}'[j] = v.\text{val}[j]$  for all  $0 \leq j \leq m$ 
7:   return
8: for all children  $c$  of  $v$  in  $T$  do
9:   for  $j \in \{m, \dots, 1\}$  do
10:    if  $v$  is a block vertex then
11:      if  $c$  does not conflict with  $f$  then
12:         $v.\text{val}'[j] = \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} (c.\text{val}'[j_1] + v.\text{val}'[j_2])$ 
13:      else
14:         $v.\text{val}'[j] = \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} (c.\text{val}[j_1] + v.\text{val}'[j_2])$ 
15:      if all the children of  $v$  conflict with  $f$  then
16:         $v.\text{val}'[j] = v.\text{val}'[j] + 1$  for all  $1 \leq j \leq m$ 
17:      else
18:         $v.\text{val}'[j] = \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} (c.\text{val}'[j_1] \cdot v.\text{val}'[j_2])$ 

```

---

FIGURE 10. An algorithm for computing  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D' \cup \{f\}, \Delta))$  for  $\Delta$  with an lhs chain.

$$= \sum_{\substack{j_1+j_2=j \\ 0 \leq j_1 \leq |D[c]| \\ 0 \leq j_2 \leq |D[\text{prev}(c)]|}} \mathbb{E}_{D' \sim U_{j_1}(D[c])}(\mathcal{I}_{\text{MC}}(D', \Delta)) \cdot \mathbb{E}_{D' \sim U_{j_2}(D[\text{prev}(c)])}(\mathcal{I}_{\text{MC}}(D', \Delta)) \quad \square$$

The algorithm MCSShapleyF that computes  $\mathbb{E}_{D' \sim U_m(D \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D' \cup \{f\}, \Delta))$  is shown in Figure 10. We define:

$$v.\text{val}'[j] = \mathbb{E}[\text{number of repairs of } E \cup \{f\} \text{ for a random subset } E \text{ of size } j \text{ of } D[v]]$$

First, we set  $v.\text{val}'[0] = 1$  for every vertex  $v$ , as when  $f$  is added to the empty set we obtain a consistent database that has a single repair—the whole database. Then, we again consider three possible types of vertices. For vertices  $v$  that conflict with  $f$  we have that  $v.\text{val}'[j] = v.\text{val}[j] + 1$ , as  $f$  violates the FDs with *every* non-empty subset of  $D[v]$ ; hence,

for each such subset,  $\{f\}$  is an additional repair, and the number of repairs increases by one compared to the number of repairs without  $f$ . For a vertex  $v$  that does not match  $f$ , it holds that  $f$  does not violate the constraints with any subset of  $D[v]$ ; thus, it does not affect the number of repairs and we have that  $v.\text{val}'[j] = v.\text{val}[j]$ . The same holds for the leaves of  $T$  that do not conflict with  $f$ .

For the rest of the vertices  $v$ , the computation is similar to the one in **MCSHapley**. In particular, we go over the children  $c$  of  $v$  in the for loop of line 8, and compute  $v.\text{val}'$  using dynamic programming. We observe that if a child  $c$  of a block vertex  $v$  conflicts with  $f$ , then when  $f$  is added to a subset  $E$  of  $D[c]$ , none of the repairs of  $E \cup \{f\}$  contains  $f$ , but  $\{f\}$  is an additional repair. For such children  $c$ , we use the value  $c.\text{val}$  in the calculation of line 14, where we ignore  $f$ . Hence, if all the children of  $v$  conflict with  $f$ , we compute  $v.\text{val}'$  in the exact same way we compute  $v.\text{val}$ , while ignoring the fact  $f$ . Then, we increase the computed value by one in line 18, to reflect the additional repair  $\{f\}$ .

If one of the children  $c$  of  $v$  does not conflict with  $f$  (note that there is at most one such child), then we take  $f$  into account in the computation of line 12, where we use the value  $c.\text{val}'$  rather than  $c.\text{val}$ . In this case, there is no need to increase the computed value by one, as we have already considered the addition of  $f$ , and the fact  $f$  may appear in some of the repairs of the subset of  $D[c]$  (or, again, be a repair on its own).

For a subblock vertex  $v$ , we compute  $v.\text{val}'$  in the same way we compute  $v.\text{val}$  in **MCSHapley**, but we use the value  $c.\text{val}'$  in the computation. In this case, each repair of a subset  $E$  of  $D[c] \cup D[\text{prev}(c)]$  is a union of a repair of  $(E \cap D[c]) \cup \{f\}$  and a repair of  $(E \cap D[\text{prev}(c)]) \cup \{f\}$ . Note that in this case, for a child  $c$  of  $v$  that does not match  $f$  we have that  $c.\text{val}' = c.\text{val}$ ; hence, the fact  $f$  is again only taken into account when considering a child  $c$  of  $v$  that matches  $f$ , and there is no risk in counting the repair  $\{f\}$  twice.

**7.2. Approximation.** Repair counting for  $\Delta = \{A \rightarrow B, B \rightarrow A\}$  is the problem of counting the maximal matchings of a bipartite graph. As the values  $\text{Shapley}(D, f, \Delta, \mathcal{I}_{\text{MC}})$  are nonnegative and sum up to the number of repairs, we conclude that an FPRAS for Shapley implies an FPRAS for the number of maximal matchings. To the best of our knowledge, existence of the latter is a long-standing open problem [JR18]. This is also the case for any  $\Delta'$  that is not equivalent to an FD set with an lhs chain, since there is a fact-wise reduction from  $\Delta$  to such  $\Delta'$  [LK17].

**7.3. Generalization to Multiple Relations.** As in the case of the drastic measure, we can generalize the upper bound of this section to schemas with multiple relation symbols using dynamic programming. We again consider an arbitrary order  $R_1, \dots, R_n$  of the relation symbols of the schema, and denote:

$$T_j^m = \mathbb{E}_{D' \sim U_m(D_{R_j} \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D', \Delta_{R_j}))$$

and:

$$P_j^m = \mathbb{E}_{D' \sim U_m(D^j \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D', \Delta^j))$$

The value  $T_j^m$  can be computed in polynomial time, using the algorithm of Figure 9, as we assume that each  $\Delta_{R_j}$  has an lhs chain. As for the value  $P_j^m$ , we have that  $P_1^m = T_1^m$ , and we prove the following for  $j > 1$ . (Recall that we denote by  $\Delta^j$  the FD set  $\Delta_{R_1} \cup \dots \cup \Delta_{R_j}$  and by  $D^j$  the database  $D_{R_1} \cup \dots \cup D_{R_j}$ .)

**Lemma 7.4.** *For every  $j \in \{2, \dots, n\}$  we have that:*

$$P_j^m = \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \times T_j^{m_1} \times P_{j-1}^{m_2}$$

*Proof.* A basic observation here is that the number of repairs of  $D_{R_1} \cup \dots \cup D_{R_j}$  is a product of the number of repairs of  $D_{R_j}$  and the number of repairs of  $D_{R_1} \cup \dots \cup D_{R_{j-1}}$ , since there are no conflicts among facts over different relation symbols. Thus, we have the following:

$$\begin{aligned} P_j^m &= \mathbb{E}_{D' \sim U_m(D^j \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D', \Delta^j)) = \sum_{\substack{D' \subseteq D^j \setminus \{f\} \\ |D'|=m}} \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \mathcal{I}_{\text{MC}}(D', \Delta^j) \\ &= \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_1|=m_1, |E_2|=m_2}} \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \times (\mathcal{I}_{\text{MC}}(E_1, \Delta_{R_j}) \times \mathcal{I}_{\text{MC}}(E_2, \Delta^{j-1})) \\ &= \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \left[ \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \right. \\ &\quad \times \left. \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_1|=m_1, |E_2|=m_2}} \left[ \frac{1}{\binom{|D_{R_j} \setminus \{f\}|}{m_1}} \times \frac{1}{\binom{|D^{j-1} \setminus \{f\}|}{m_2}} \times (\mathcal{I}_{\text{MC}}(E_1, \Delta_{R_j}) \times \mathcal{I}_{\text{MC}}(E_2, \Delta^{j-1})) \right] \right] \\ &= \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \left[ \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \right. \\ &\quad \times \left( \sum_{\substack{E_1 \subseteq D_{R_j} \setminus \{f\} \\ |E_1|=m_1}} \frac{1}{\binom{|D_{R_j} \setminus \{f\}|}{m_1}} \times \mathcal{I}_{\text{MC}}(E_1, \Delta_{R_j}) \right) \\ &\quad \times \left( \sum_{\substack{E_2 \subseteq D^{j-1} \setminus \{f\} \\ |E_2|=m_2}} \frac{1}{\binom{|D^{j-1} \setminus \{f\}|}{m_2}} \times \mathcal{I}_{\text{MC}}(E_2, \Delta^{j-1}) \right) \left. \right] \\ &= \frac{1}{\binom{|D^j \setminus \{f\}|}{m}} \sum_{\substack{0 \leq m_1 \leq |D_{R_j} \setminus \{f\}| \\ 0 \leq m_2 \leq |D^{j-1} \setminus \{f\}| \\ m_1 + m_2 = m}} \binom{|D_{R_j} \setminus \{f\}|}{m_1} \times \binom{|D^{j-1} \setminus \{f\}|}{m_2} \times T_j^{m_1} \times P_{j-1}^{m_2} \quad \square \end{aligned}$$

The computation of  $\mathbb{E}_{D' \sim U_m(D^{j-1} \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D' \cup \{f\}, \Delta^j))$  is very similar, with the only difference being the fact that:

$$T_j^m = \mathbb{E}_{D' \sim U_m(D_{R_j} \setminus \{f\})}(\mathcal{I}_{\text{MC}}(D' \cup \{f\}, \Delta_{R_j}))$$

for the relation symbol  $R_j$  of  $f$ . This value can be computed in polynomial time using the algorithm of Figure 10. Finally, as in the case of the drastic measure, it is rather straightforward that the lower bound of Theorem 7.1 generalizes to the case where the FD set  $\Delta_R$  has no lhs chain (up to equivalence) for at least one relation symbol  $R$  of the schema.

## 8. CONCLUSIONS

We studied the complexity of calculating the Shapley value of database facts for basic inconsistency measures, focusing on FD constraints. We showed that two of them are computable in polynomial time: the number of violations ( $\mathcal{I}_{MI}$ ) and the number of problematic facts ( $\mathcal{I}_P$ ). In contrast, each of the drastic measure ( $\mathcal{I}_d$ ) and the number of repairs ( $\mathcal{I}_{MC}$ ) features a dichotomy in complexity, where the tractability condition is the possession of an lhs chain (up to equivalence). For the cost of a cardinality repair ( $\mathcal{I}_R$ ) we showed a tractable fragment and an intractable fragment, but a gap remains on certain FD sets—the ones that do not have an lhs chain, and yet, a cardinality repair can be computed in polynomial time. We also studied the approximability of the Shapley value and showed, among other things, an FPRAS for  $\mathcal{I}_d$  and a dichotomy in the existence of an FPRAS for  $\mathcal{I}_R$ .

Many other directions are left open for future research. First, the picture is incomplete for the measure  $\mathcal{I}_R$ . In particular, the complexity of the exact computation is open for the bipartite matching constraint  $\{A \rightarrow B, B \rightarrow A\}$  that, unlike the known FD sets in the intractable fragment, has an FPRAS. In general, we would like to complete the picture of  $\mathcal{I}_R$  towards a full dichotomy. Moreover, for the schemas where there is no FPRAS for  $\mathcal{I}_R$ , our results neither imply nor refute the existence of a constant-ratio approximation (for *some* constant). Second, the problems are immediately extendible to any type of constraints other than functional dependencies, such as denial constraints, tuple generating dependencies, and so on. Third, it would be interesting to see how the results extend to wealth distribution functions other than Shapley, for instance the Banzhaff Power Index [DS79]. The tractable cases remain tractable for the Banzhaff Power Index, but it is not clear how (and whether) our proofs for the lower bounds generalize to this function. Another direction is to investigate whether properties of the database (e.g., bounded treewidth) have an impact on the complexity of computing the Shapley value. Finally, there is the practical question of implementation: while our algorithms terminate in polynomial time, we believe that they are hardly scalable without further optimization and heuristics ad-hoc to the use case; developing those is an important challenge for future research.

## ACKNOWLEDGMENT

This work was supported by the Israel Science Foundation (ISF), Grant 768/19, and the German Research Foundation (DFG) Project 412400621 (DIP program).

## REFERENCES

- [AdK14] Haris Aziz and Bart de Keijzer. Shapley meets Shapley. In *STACS*, volume 25 of *LIPICs*, pages 99–111. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2014.
- [Bac02] Roland Bacher. Determinants of matrices related to the pascal triangle. *Journal de Théorie des Nombres de Bordeaux*, 14, 01 2002.

- [BDG<sup>+</sup>19] Omar Besbes, Antoine Désir, Vineet Goyal, Garud Iyengar, and Raghav Singal. Shapley meets uniform: An axiomatic framework for attribution in online advertising. In *WWW*, pages 1713–1723. ACM, 2019.
- [Ber18] Leopoldo E. Bertossi. Measuring and computing database inconsistency via repairs. In *SUM*, volume 11142 of *Lecture Notes in Computer Science*, pages 368–372. Springer, 2018.
- [Ber19] Leopoldo E. Bertossi. Repair-based degrees of database inconsistency. In *LPNMR*, volume 11481 of *Lecture Notes in Computer Science*, pages 195–209. Springer, 2019.
- [BG20] Leopoldo E. Bertossi and Floris Geerts. Data quality and explainable AI. *J. Data and Information Quality*, 12(2):11:1–11:9, 2020.
- [CPRT15] Laurence Cholvy, Laurent Perrussel, William Raynaut, and Jean-Marc Thévenin. Towards consistency-based reliability assessment. In *AAMAS*, pages 1643–1644. ACM, 2015.
- [DS79] Pradeep Dubey and Lloyd S. Shapley. Mathematical properties of the banzhaf power index. *Mathematics of Operations Research*, 4(2):99–131, 1979.
- [GGR98] Oded Goldreich, Shafi Goldwasser, and Dana Ron. Property testing and its connection to learning and approximation. *J. ACM*, 45(4):653–750, 1998.
- [GH06] John Grant and Anthony Hunter. Measuring inconsistency in knowledgebases. *J. Intell. Inf. Syst.*, 27(2):159–184, 2006.
- [GH11] John Grant and Anthony Hunter. Measuring consistency gain and information loss in stepwise inconsistency resolution. In *ECSQARU*, volume 6717, pages 362–373. Springer, 2011.
- [GH13] John Grant and Anthony Hunter. Distance-based measures of inconsistency. In *ECSQARU*, volume 7958 of *Lecture Notes in Computer Science*, pages 230–241. Springer, 2013.
- [GH15] John Grant and Anthony Hunter. Using Shapley inconsistency values for distributed information systems with uncertainty. In *ECSQARU*, volume 9161 of *Lecture Notes in Computer Science*, pages 235–245. Springer, 2015.
- [GH17] John Grant and Anthony Hunter. Analysing inconsistent information using distance-based measures. *Int. J. Approx. Reasoning*, 89:3–26, 2017. doi:10.1016/j.ijar.2016.04.004.
- [Gul89] Faruk Gul. Bargaining foundations of Shapley value. *Econometrica: Journal of the Econometric Society*, pages 81–95, 1989.
- [HK06] Anthony Hunter and Sébastien Konieczny. Shapley inconsistency values. In *KR*, pages 249–259. AAAI Press, 2006.
- [HK08] Anthony Hunter and Sébastien Konieczny. Measuring inconsistency through minimal inconsistent sets. In *KR*, pages 358–366. AAAI Press, 2008.
- [HK10] Anthony Hunter and Sébastien Konieczny. On the measure of conflicts: Shapley inconsistency values. *Artif. Intell.*, 174(14):1007–1026, 2010.
- [JR18] Yifan Jing and Akbar Rafiey. Counting maximal near perfect matchings in quasirandom and dense graphs. *CoRR*, abs/1807.04803, 2018.
- [Kim12] Benny Kimelfeld. A dichotomy in the complexity of deletion propagation with functional dependencies. In *PODS*, pages 191–202, 2012.
- [KLM03] Sébastien Konieczny, Jérôme Lang, and Pierre Marquis. Quantifying information and contradiction in propositional logic through test actions. In *IJCAI*, pages 106–111. Morgan Kaufmann, 2003.
- [Kni03] Kevin M. Knight. Two information measures for inconsistent sets. *Journal of Logic, Language and Information*, 12(2):227–248, 2003.
- [LBKS20] Ester Livshits, Leopoldo E. Bertossi, Benny Kimelfeld, and Moshe Sebag. The Shapley value of tuples in query answering. In *ICDT*, volume 155 of *LIPICs*, pages 20: 1–20: 19. Schloss Dagstuhl, 2020.
- [LF18] Christophe Labreuche and Simon Fossier. Explaining multi-criteria decision aiding models with an extended Shapley value. In *IJCAI*, pages 331–339. ijcai.org, 2018.
- [LK17] Ester Livshits and Benny Kimelfeld. Counting and enumerating (preferred) database repairs. In *PODS*, pages 289–301. ACM, 2017.
- [LK21] Ester Livshits and Benny Kimelfeld. The shapley value of inconsistency measures for functional dependencies. In *ICDT*, volume 186 of *LIPICs*, pages 15:1–15:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021.
- [LKR20] Ester Livshits, Benny Kimelfeld, and Sudeepa Roy. Computing optimal repairs for functional dependencies. *ACM Trans. Database Syst.*, 45(1):4: 1–4: 46, 2020.

- [LKT<sup>+</sup>21] Ester Livshits, Rina Kochirgan, Segev Tsur, Ihab F. Ilyas, Benny Kimelfeld, and Sudeepa Roy. Properties of inconsistency measures for databases. In *SIGMOD Conference*, pages 1182–1194. ACM, 2021.
- [LL17] Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *NIPS*, pages 4765–4774, 2017.
- [LZS15] Zhenliang Liao, Xiaolong Zhu, and Jiaorong Shi. Case study on initial allocation of shanghai carbon emission trading based on Shapley value. *Journal of Cleaner Production*, 103:338–344, 2015.
- [MCL<sup>+</sup>10] Richard TB Ma, Dah Ming Chiu, John Lui, Vishal Misra, and Dan Rubenstein. Internet economics: The use of Shapley value for isp settlement. *IEEE/ACM Transactions on Networking (TON)*, 18(3):775–787, 2010.
- [MLJ11] Kedian Mu, Weiru Liu, and Zhi Jin. Measuring the blame of each formula for inconsistent prioritized knowledge bases. *Journal of Logic and Computation*, 22(3):481–516, 02 2011. [arXiv: https://academic.oup.com/logcom/article-pdf/22/3/481/3177718/exr002.pdf](https://academic.oup.com/logcom/article-pdf/22/3/481/3177718/exr002.pdf).
- [Nen03] Tatiana Nenova. The value of corporate voting rights and control: A cross-country analysis. *Journal of financial economics*, 68(3):325–351, 2003.
- [NN11] Ramasuri Narayanan and Yadati Narahari. A Shapley value-based approach to discover influential nodes in social networks. *IEEE Transactions on Automation Science and Engineering*, 8(1):130–147, 2011.
- [PZ03] Leon Petrosjan and Georges Zaccour. Time-consistent Shapley value allocation of pollution cost reduction. *Journal of economic dynamics and control*, 27(3):381–398, 2003.
- [RKL20] Alon Reshef, Benny Kimelfeld, and Ester Livshits. The impact of negation on the complexity of the Shapley value in conjunctive queries. In *PODS*, pages 285–297. ACM, 2020.
- [Sha53] Lloyd S Shapley. A value for n-person games. In Harold W. Kuhn and Albert W. Tucker, editors, *Contributions to the Theory of Games II*, pages 307–317. Princeton University Press, Princeton, 1953.
- [Thi09] Matthias Thimm. Measuring inconsistency in probabilistic knowledge bases. In *UAI*, pages 530–537. AUAI Press, 2009.
- [Thi17] Matthias Thimm. On the compliance of rationality postulates for inconsistency measures: A more or less complete picture. *KI*, 31(1):31–39, 2017.
- [Val79a] Leslie G. Valiant. The complexity of enumeration and reliability problems. *SIAM J. Comput.*, 8(3):410–421, 1979.
- [Val79b] L.G. Valiant. The complexity of computing the permanent. *Theoretical Computer Science*, 8(2):189–201, 1979. doi:10.1016/0304-3975(79)90044-6.
- [YVCB18] Bruno Yun, Srdjan Vesic, Madalina Croitoru, and Pierre Bisquert. Inconsistency measures for repair semantics in OBDA. In *IJCAI*, pages 1977–1983. ijcai.org, 2018.