

# Cleaning Data with Forbidden Itemsets

Joeri Rammelaere  
University of Antwerp  
Middelheimlaan 1  
Antwerp, Belgium  
joeri.rammelaere@uantwerp.be

Floris Geerts  
University of Antwerp  
Middelheimlaan 1  
Antwerp, Belgium  
floris.geerts@uantwerp.be

Bart Goethals  
University of Antwerp  
Middelheimlaan 1  
Antwerp, Belgium  
bart.goethals@uantwerp.be

**Abstract**—Methods for cleaning dirty data typically rely on additional information about the data, such as user-specified constraints that specify when a database is dirty. These constraints often involve domain restrictions and illegal value combinations. Traditionally, a database is considered clean if all constraints are satisfied. However, many real-world scenarios only have a dirty database available. In such a context, we adopt a dynamic notion of data quality, in which the data is clean if an error discovery algorithm does not find any errors. We introduce *forbidden itemsets* which capture unlikely value co-occurrences in dirty data, and we derive properties of the lift measure to provide an efficient algorithm for mining low lift forbidden itemsets. We further introduce a repair method which guarantees that the repaired database does not contain any low lift forbidden itemsets. The algorithm uses nearest neighbor imputation to suggest possible repairs. Optional user interaction can easily be integrated into the proposed cleaning method. Evaluation on real-world data shows that errors are typically discovered with high precision, while the suggested repairs are of good quality and do not introduce new forbidden itemsets, as desired.

## I. INTRODUCTION

In recent years, research on detecting inconsistencies in data has focused on a constraint-based data quality approach: a set of constraints in some logical formalism is associated with a database, and the data is considered consistent or clean if and only if all constraints are satisfied. Many such formalisms exist, capturing a wide variety of inconsistencies, and systematic ways of repairing the detected inconsistencies are in place. A frequently asked question is, “where do these constraints come from?”. The common answer is that they are either supplied by experts, or automatically discovered from the data [1], [2]. In most real-world scenarios, however, the underlying data is dirty. This raises concerns about the reliability of the discovered constraints.

Assume for the moment that the constraints are reliable and are used to repair (clean) a dirty database. For the sake of the argument, what if one re-runs the constraint discovery algorithm on the repair and finds other constraints? Does this imply that the repair is not clean after all, or that the discovery algorithm may in fact find unreliable constraints? It is a typical chicken or the egg dilemma. The problem is that constraints are considered to be *static*: once found, they are treated as a gold standard. To remedy this, we propose a *dynamic* notion of data quality. The idea is simple:

“We consider a given database to be clean if a constraint discovery algorithm does not detect any violated constraints on that data.”

Constraints thus reflect inconsistencies which depend on the actual data. Clearly, this dynamic notion presents a new challenge when repairing, since the constraints may shift during repairs. Indeed, it does not suffice to only resolve inconsistencies for constraints found on the original dirty data, one also has to ensure that no new constraints (and thus new inconsistencies) are found on the repaired data. Note that we focus on repairing data under dynamic constraints, as opposed to cleaning dynamic data. To our knowledge, this is a new view on data quality, raising many interesting challenges.

The main contribution of this paper is to illustrate this dynamic view on data quality for a particular class of constraints. In particular, we consider errors that can be caught by so-called *edits*, which is “the” constraint formalism used by census bureaus worldwide [3], [4] and can be seen as a simple class of denial constraints [5], [6]. Intuitively, an edit specifies *forbidden value combinations*. For example, an age attribute cannot take a value higher than 130, a city can only have certain zip codes, and people of young age cannot have a driver’s license. These edits are typically designed by domain experts and are generally accepted to be a good constraint formalism for detecting errors that occur in single records. However, we aim to automatically discover such edits in an unsupervised way by using pattern mining techniques. To make the link to pattern mining more explicit and to differentiate from edits, we call our patterns *forbidden itemsets*. Although our technique is designed such that human supervision is not compulsory, we show that optional user interaction can be readily integrated to improve the reliability of the methods.

Pattern mining techniques are typically used to uncover positive associations between items, measured by different interestingness measures such as frequency, confidence, lift, and many others. Based on experience, *discovered patterns often reveal errors in the data* in addition to interesting associations. For example, an association rule which holds for 99% of the data could be interesting in itself, but might also represent a well-known dependency in the data which should hold for 100%. The fact that the association doesn’t hold for 1% of the data is then more interesting. This 1% of the data often points to unlikely co-occurrences of values in the data, which forbidden itemsets aim to capture. In order to reliably detect unlikely co-occurrences, it is clear that a large body of clean data is needed. We therefore focus on low error rate data, such as census data or data that underwent some curation [7].

Apart from detecting errors we also aim to provide suggestions for how to repair them, i.e., suggest modifications to the data such that after these modifications, no new forbidden

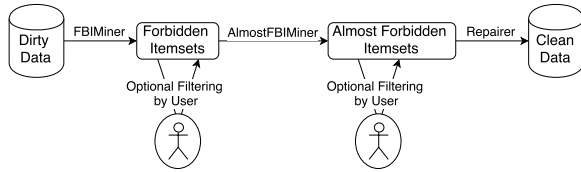


Figure 1. Schematic overview of the proposed cleaning mechanism.

itemsets are found. Here we again take inspiration from census data imputation methods that assume the presence of enough clean data [4] and take suggested modifications from similar, clean objects. The availability of clean data is commonly assumed in repairing methods, either as a large part of the input data or, for example, in the form of master data [8], [9].

Figure 1 gives a schematic overview of the proposed cleaning mechanism in its entirety. We capture unlikely co-occurrences by means of forbidden itemsets. Our algorithm FBIMINER employs pruning strategies to optimize the discovery of forbidden itemsets. Linking back to the beginning of the introduction, we will regard data to be dirty if FBIMINER finds forbidden itemsets in the data. Users may optionally filter out forbidden itemsets by declaring them as valid. Furthermore, we also devise a *repairing algorithm* that repairs a dirty dataset and ensures that no forbidden itemsets exist in the repair, hence it is indeed clean. To achieve this, we consider so-called *almost forbidden itemsets* and present an algorithm A-FBIMINER for mining them. Again, users can optionally filter out such itemsets. All algorithms are experimentally validated.

**Organization of the paper.** The paper is organized as follows: In Sect. II we discuss the most relevant related work. Notations and basic concepts are presented in Sect. III, followed by a formal problem statement in Sect. IV. Section V introduces forbidden itemsets, their properties, and the FBIMINER algorithm. Section VI presents the repair algorithm, focusing on our strategy to avoid new inconsistencies. The possibility of user interaction is discussed in Sect. VII. In Sect. VIII we show experimental results, before we pose a conclusion in Sect. IX. Proofs and additional plots can be found in the appendix.

## II. RELATED WORK

There has been a substantial amount of work on constraint-based data quality in the database community (see [1], [2] for recent surveys). Most relevant to our work are constraints that concern single tuples such as constant conditional functional dependencies [1] and constant denial constraints [5], [6]. Algorithms are in place to (i) discover these constraints from data [10], [11], [1], [12]; and (ii) repair the errors, once the constraints are fixed [13], [14], [15], [5]. Moreover, user interaction is often used to guide the repairing algorithms [16], [17], [18] and statistical methods are in place to measure the impact of repairing [19]. As previously mentioned, all these methods use a static notion of cleanliness. Our notion of forbidden itemsets is closest in spirit to edits, used by census bureaus [3] and our repairing method is similar to hot-deck imputation methods [4]. Capturing and detecting inconsistencies is also closely related to anomaly and outlier detection (see [20], [21], [22] for recent surveys). A recent study [23] provides a comparison of detection methods. In the pattern

mining community many different interestingness measures exist. We mention [24] in which outliers are discovered using a measure that is similar to ours. Furthermore, Error-Tolerant Itemsets [25], [26] can be regarded as the inverse of the forbidden itemsets. Compared to these methods, we identify new properties of the lift measure to speed up the discovery, and use forbidden itemsets to both detect *and repair* data.

## III. PRELIMINARIES

We consider datasets  $\mathcal{D}$  consisting of a finite collection of objects. An *object*  $o$  is a pair  $\langle \text{oid}, I \rangle$  where  $\text{oid}$  is an *object identifier*, e.g., a natural number, and  $I$  is an itemset. An *itemset* is a set of *items* of the form  $(A, v)$ , where  $A$  comes from a set  $\mathcal{A}$  of *attributes* and  $v$  is a *value* from a finite domain of categorical values  $\text{dom}(A)$  of  $A$ . An itemset contains at most one item for each attribute in  $\mathcal{A}$ . We define the value of an object  $o = \langle \text{oid}, I \rangle$  in attribute  $A$ , denoted by  $o[A]$ , as the value  $v$  where  $(A, v) \in I$ , and let  $o[A]$  be undefined otherwise. We denote by  $\mathcal{I}$  the set of all attribute/value pairs  $(A, v)$ .

An object  $o = \langle \text{oid}, I \rangle$  is said to *support* an itemset  $J$  if  $J \subseteq I$ , i.e.,  $J$  is contained in  $I$ . The *cover* of an itemset  $J$  in  $\mathcal{D}$ , denoted by  $\text{cov}(J, \mathcal{D})$ , is the set of  $\text{oid}$ 's of objects in  $\mathcal{D}$  that support  $J$ . The *support* of  $J$  in  $\mathcal{D}$ , denoted by  $\text{supp}(J, \mathcal{D})$ , is the number of  $\text{oid}$ 's in its cover in  $\mathcal{D}$ . Similarly, the *frequency* of an itemset  $J$  in  $\mathcal{D}$  is the fraction of  $\text{oid}$ 's in its cover:  $\text{freq}(J, \mathcal{D}) = \text{supp}(J, \mathcal{D}) / |\mathcal{D}|$ , where  $|\mathcal{D}|$  is the number of objects in  $\mathcal{D}$ . We sometimes represent  $\mathcal{D}$  in a *vertical data layout* denoted by  $\mathcal{D}\downarrow$ . Formally,  $\mathcal{D}\downarrow = \{(i, \text{cov}(\{i\}, \mathcal{D})) \mid i \in \mathcal{I}, \text{cov}(\{i\}, \mathcal{D}) \neq \emptyset\}$ . Clearly, one can freely go from  $\mathcal{D}$  to  $\mathcal{D}\downarrow$ , and vice versa.

We assume that a similarity measure between objects is given, and denote by  $\text{sim}(o, o')$  the similarity between objects  $o$  and  $o'$ . The similarity between two datasets  $\mathcal{D}$  and  $\mathcal{D}'$  is denoted by  $\text{sim}(\mathcal{D}, \mathcal{D}')$ . Any similarity measure can be used in our setting. We describe the similarity function used in our experiments in Sect. VIII.

## IV. PROBLEM STATEMENT

We first phrase our problem in full generality (following [27]) before making things more specific in the next section. Consider a dataset  $\mathcal{D}$  and some *constraint language*  $\mathcal{L}$  for expressing properties that indicate *dirty*ness in the data, e.g.,  $\mathcal{L}$  could consist of conditional functional dependencies [1], edits [4], or association rules [28]. Furthermore, let  $q$  be a *selection predicate* (evaluating to true or false) that assesses the relevance of constraints  $\varphi \in \mathcal{L}$  in  $\mathcal{D}$ . For example, when  $\varphi$  is a conditional functional dependency,  $q(\mathcal{D}, \varphi)$  may return true if  $\varphi$  is violated in  $\mathcal{D}$ . We denote by  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  the set of all *dirty constraints*, i.e., all  $\varphi \in \mathcal{L}$  for which  $q(\mathcal{D}, \varphi)$  evaluates to true. For example,  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  may consist of all violated conditional functional dependencies, all edits that apply to an object, or all low confidence association rules.

**Definition 1.** A dataset  $\mathcal{D}$  is said to be *clean* relative to language  $\mathcal{L}$  and predicate  $q$  if  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  is empty;  $\mathcal{D}$  is called *dirty* otherwise.  $\square$

With this definition we take a *completely new view* on data quality. Indeed, existing work in this area [1] typically *fixes* the constraints up front, regardless of the data. For example,

Table I. EXAMPLE FORBIDDEN ITEMSETS FOUND IN UCI DATASETS

Forbidden Itemsets	Dataset	$\tau$
Sex=Female, Relationship=Husband Sex=Male, Relationship=Wife Age=<18, Marital-status=Married-c-s Age=<18, Relationship=Husband Relationship=Not-in-family, Marital=Married-c-s	Adult	0.01
aquatic=0, breathes=0 (clam) type=mammal, eggs=1 (platypus) milk=1, eggs=1 (platypus) type=mammal, toothed=0 (platypus) eggs=0, toothed=0 (scorpion) milk=1, toothed=0 (platypus) tail=1, backbone=0 (scorpion)	Zoo	0.1
bruises=t, habitat=l population=y, cap-shape=k cap-surface=s, odor=n, habitat=d cap-surface=s, gill-size=b, habitat=d edible=e, habitat=d, cap-shape=k	Mushroom	0.025

edits are often designed by experts and then compared with the data. In our definition, we only specify the *class* of constraints, e.g., edits. Which edits are used for declaring the data clean or dirty *depends entirely on the underlying data*. We thus introduce a *dynamic* rather than a static notion of dirtiness/cleanliness of data: when the data changes, so do the edits under consideration. With this notion at hand, we are next interested in repairs of the data. Intuitively, a repair of a dirty dataset is a modified dataset that is clean.

**Definition 2.** Given datasets  $\mathcal{D}$  and  $\mathcal{D}'$ , language  $\mathcal{L}$ , predicate  $q$  and similarity function  $\text{sim}$ , we say that  $\mathcal{D}'$  is an  $(\mathcal{L}, q)$ -*repair* if (i)  $\mathcal{D}'$  has the same set of object identifiers as  $\mathcal{D}$ ; and (ii)  $\text{dirty}(\mathcal{D}', \mathcal{L}, q)$  is empty. An  $(\mathcal{L}, q)$ -repair  $\mathcal{D}'$  is *minimal* if  $\text{sim}(\mathcal{D}, \mathcal{D}')$  is maximal amongst all  $(\mathcal{L}, q)$ -repairs of  $\mathcal{D}$ .  $\square$

## V. FORBIDDEN ITEMSETS

We first specialize constraint language  $\mathcal{L}$  and predicate  $q$  such that  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  corresponds to *inconsistencies* in  $\mathcal{D}$ . We define  $\mathcal{L}$  as the class of itemsets and define  $q$  such that  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  corresponds to what we will call *forbidden itemsets* (Sect. V-A). Intuitively, these are itemsets which do occur in the data, despite being very improbable with respect to the rest of the data. Furthermore, we show how to compute  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  for *low lift* forbidden itemsets. As such itemsets are typically *infrequent*, existing itemset mining algorithms are not optimized for this task. We therefore derive properties of the lift measure that allow for substantial pruning when mining forbidden itemsets (Sect. V-B). We conclude by presenting a version of the well-known Eclat algorithm [29], enhanced with our derived pruning strategies and optimizations specific for the task of mining low lift forbidden itemsets (Sect. V-C).

Before formally introducing forbidden itemsets as a constraint language  $\mathcal{L}$ , let us provide some additional motivation for considering invalid or unlikely value combinations (as represented by forbidden itemsets) as error detection formalism. First of all, invalid value combinations have been used for decades to detect errors in census data starting with the seminal work by Fellegi and Holt [3]. Second, although more expressive formalisms such as conditional functional dependencies (CFDs) [30] and denial constraints (DCs) [5], [6] have become popular for error detection and repairing, many constraints used in practice are very simple. As an example, most of the constraints reported in Table 4 in [23] can be

regarded as constraints that only involve constants. This is clear for “checks” that specify invalid domain values. However, even a functional dependency such as  $(\text{zip} \rightarrow \text{state})$  can be regarded as a (finite) collection of constant rules that associate specific zip codes to state names. The violations of these rules clearly are invalid value combinations. Similarly, almost half of the DCs reported in [6] only involve constants and concern single tuples. It therefore seems natural to first gain a better understanding of these simple constraints in our dynamic data quality setting. Finally, the discovery of CFDs and DCs [1], [10], [11] in their full generality is very slow due to the high expressiveness of these constraint languages. Experiments show that the discovery may take up to hours on a single machine. This makes general constraints infeasible in settings such as ours, where interactivity or quick response times are needed. For all these reasons, we believe that forbidden itemsets provide a good balance between expressiveness, efficiency of discovery, and efficacy in error detection. Furthermore, they are easily interpretable, allowing users to inspect and filter out false positives, as discussed in Sect. VII.

### A. Low Lift Itemsets

We consider  $\mathcal{L}$  consisting of the class of itemsets and want to define  $q$  such that for an itemset  $I$ ,  $q(\mathcal{D}, I)$  is true if  $I$  corresponds to a possible inconsistency in the data. In general, we can use a *likeliness function*  $L : 2^{\mathcal{I}} \times \mathcal{D} \rightarrow \mathbb{R}$  that indicates how likely the occurrence of an itemset is in  $\mathcal{D}$ . Such a likeliness function can be defined to accommodate various types of constraints on value combinations within a single tuple. If we denote by  $\tau$  a maximum likeliness threshold, then we define  $\tau$ -*forbidden itemsets* as follows.

**Definition 3.** Let  $\mathcal{D}$  be a dataset,  $L$  a likeliness function,  $I$  an itemset and  $\tau$  a maximum likeliness threshold. Then,  $I$  is called a  $\tau$ -*forbidden itemset* whenever  $L(I, \mathcal{D}) \leq \tau$ .  $\square$

Phrased in the general framework from the previous section, we thus have that  $\mathcal{L}$  is the class of itemsets and  $q(\mathcal{D}, I)$  is true if  $I$  is a  $\tau$ -forbidden itemset. Hence,  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  consists of all  $\tau$ -forbidden itemsets in  $\mathcal{D}$ .

In this paper, we propose to use the *lift* measure of an itemset as likeliness function. Intuitively, it gives an indication of how probable the co-occurrence of a set of items is given their separate frequencies. Lift is generally used as an interestingness measure in association rule mining, where rules with a high lift between antecedent and consequent are considered the most interesting [24], and has also been used for constraint discovery [10]. A straightforward extension of lift from rules to itemsets assumes “full” independence among all individual items [31, p. 310]. In other words, an itemset  $I$  is regarded as “unlikely” when  $\text{freq}(I, \mathcal{D})$  is much smaller than  $\text{freq}(\{i_1\}, \mathcal{D}) \times \dots \times \text{freq}(\{i_k\}, \mathcal{D})$ , where  $i_1, \dots, i_k$  are the items in  $I$ . However, this introduces an undesirable bias towards larger itemsets: many items with a slight negative correlation might have a lower lift than two items with a strong negative correlation.

Instead of full independence, we adopt “pairwise” independence, in which  $\text{freq}(I, \mathcal{D})$  is compared against  $\text{freq}(J, \mathcal{D}) \times \text{freq}(I \setminus J, \mathcal{D})$  for any non-empty  $J \subset I$ , and the *maximal discrepancy ratio* is taken as the *lift* of  $I$  in  $\mathcal{D}$ :

**Definition 4.** Let  $\mathcal{D}$  be a dataset and let  $I$  be an itemset. The *lift* of  $I$ , denoted by  $\text{lift}(I, \mathcal{D})$ , is defined as

$$\text{lift}(I, \mathcal{D}) := \frac{|\mathcal{D}| \times \text{supp}(I, \mathcal{D})}{\min_{\emptyset \subset J \subset I} \{\text{supp}(J, \mathcal{D}) \times \text{supp}(I \setminus J, \mathcal{D})\}} \quad \square$$

We note that this definition is conceptually related to association rules and has been used successfully in the context of redundancy [24] and outlier detection [32], two concepts very similar in spirit to our intended notion of inconsistencies.

One could further generalize the notion of lift by ranging over partitions of  $I$  consisting of more than two parts, full independence being a special case in which  $I$  is partitioned in all its items. We find, however, that pairwise independence is already effective for detecting unlikely value combinations and is more efficient to compute.

From now on, we refer to  $\tau$ -forbidden itemsets as itemsets  $I$  for which  $\text{lift}(I, \mathcal{D}) \leq \tau$  holds, following Def. 3. When using lift,  $\tau$  will typically be small, and we assume that  $\tau < 1$ .

**Example 1.** To illustrate that  $\tau$ -forbidden itemsets are an effective formalism for capturing unlikely value combinations, we show some example forbidden itemsets found in UCI datasets in Table I. In the Adult dataset, the co-occurrence of Female and Husband, as well as Male and Wife, are clearly erroneous. Other examples involve a married person under the age of 18 and people who are married, yet living in with an unrelated household. In the Zoo dataset, the first forbidden itemset shows that the animal *clam* in the dataset is neither aquatic nor breathing. To our knowledge, clams are in fact aquatic, so these values are indeed in error. The other forbidden itemsets detect animals that are in some way an exception in nature. For example, the *platypus* is famous for being one of few existing mammal species that lays eggs, the other species being anteaters. Similar exceptional combinations are encountered in the other UCI datasets, such as the Mushroom dataset, although the forbidden itemsets are more difficult to interpret for this dataset. While not all of these examples represent actual *errors*, they show that the forbidden itemsets are capable of detecting peculiar objects that require extra attention. Of course, it makes little sense to repair objects such as the platypus. Typically, user validation of the discovered errors will be preferable over fully automatic repairs. This is addressed in Sect. VII.  $\diamond$

## B. Properties of the Lift Measure

Before presenting our algorithm FBIMINER that mines  $\tau$ -forbidden itemsets, we describe some properties of the lift measure that underlie our algorithm.

While the lift measure is neither monotonic nor anti-monotonic, two properties that are typically used for pruning in pattern mining algorithms, it still has some properties that allow the pruning of large branches of the search tree: since a low lift requires that an itemset occurs *much less often* than its subsets, we can use the relation between the support of a  $\tau$ -forbidden itemset and the support of its subsets to prune. As we will explain shortly, FBIMINER performs a depth-first search for forbidden itemsets. To make this search efficient, pruning strategies should be in place that discard all supersets

of a particular itemset. For this purpose, we derive properties that must hold for all subsets of a  $\tau$ -forbidden itemset.

Given an itemset  $I$  we denote by  $\sigma_I^{max}$  the highest support of an item  $\{i\}$  in  $I$  or any of  $I$ 's supersets  $J$ , i.e.,  $\sigma_I^{max}$  is the highest support in  $I$ 's branch of the search tree. More formally,

$$\sigma_I^{max} := \max\{\text{supp}(\{i\}, \mathcal{D}) \mid i \in J, I \subseteq J\}.$$

This quantity can be used to obtain a lower bound on the support of subsets of  $J$  when  $J$  is a  $\tau$ -forbidden itemset:

**Proposition 1.** For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , if  $J$  is a  $\tau$ -forbidden itemset then  $\text{supp}(I, \mathcal{D}) \geq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\sigma_I^{max} \times \tau}$ .  $\square$

Furthermore, for any  $\tau$ -forbidden itemset  $J$  in the dataset, it trivially holds that  $\text{supp}(J, \mathcal{D}) \geq 1$  and thus any itemset  $I \subset J$  must have  $\text{supp}(I, \mathcal{D}) \geq \frac{|\mathcal{D}|}{\sigma_I^{max} \times \tau}$ . This implies that in the depth-first search, it suffices to expand itemsets  $I$  for which  $\text{supp}(I, \mathcal{D}) \geq \frac{|\mathcal{D}|}{\sigma_I^{max} \times \tau}$ .

Furthermore, Prop. 1 can be leveraged to show that a minimum reduction in support between subsets of a  $\tau$ -forbidden itemset is required.

**Proposition 2.** For any three itemsets  $I$ ,  $J$  and  $K$  such that  $I \subset J \subseteq K$  holds, if  $K$  is a  $\tau$ -forbidden itemset, then  $\text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D}) \geq \frac{1}{\tau} - \frac{\sigma_I^{max}}{|\mathcal{D}|} > 0$ .  $\square$

In particular, for  $J$  to be  $\tau$ -forbidden, we must have that  $\text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D}) \geq \frac{1}{\tau} - \frac{\sigma_I^{max}}{|\mathcal{D}|} > 0$  holds for any of its subsets  $I$ . During the depth-first search, when expanding  $I$  to  $J$ , if this condition is not satisfied, then  $J$  and all of its supersets  $K$  can be pruned. Furthermore, the proposition implies that  $\tau$ -forbidden itemsets are so-called generators [33], i.e., if  $J$  is  $\tau$ -forbidden then  $\text{supp}(I, \mathcal{D}) > \text{supp}(J, \mathcal{D})$  for any  $I \subset J$ . A known property of generators is that all their subsets are generators as well, meaning that the entire subtree can be pruned if a non-generator is encountered during the search.

Our next pruning method uses the lift of an itemset to bound the support of its supersets. Indeed, the denominator of the lift measure is in fact anti-monotonic.

**Proposition 3.** For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , it holds that:

$$\text{lift}(J, \mathcal{D}) \geq \frac{\text{supp}(J, \mathcal{D}) \times |\mathcal{D}|}{\min_{S \subset I} \{\text{supp}(S, \mathcal{D}) \times \text{supp}(I \setminus S, \mathcal{D})\}}. \quad \square$$

Clearly, this lower bound can be used to prune itemsets  $J$  that cannot lead to  $\tau$ -forbidden itemsets. Note that to use the lower bound, one needs a lower bound on  $\text{supp}(J, \mathcal{D})$ . We again use the trivial lower bound  $\text{supp}(J, \mathcal{D}) \geq 1$ .

Finally, since itemsets with low lift are obtained when they occur much less often than their subsets, it is expected that such forbidden itemsets will have a low support. In fact, one can precisely characterize the maximal frequency of a  $\tau$ -forbidden itemset.

**Proposition 4.** If  $I$  is a  $\tau$ -forbidden itemset then its frequency is bounded by  $\text{freq}(I, \mathcal{D}) \leq \frac{2}{\tau} - 2\sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}} - 1$ . Furthermore, for small  $\tau$ -values this bound converges (from above) to  $\frac{\tau}{4}$ .  $\square$

---

**Algorithm 1** An Eclat-based algorithm for mining low lift  $\tau$ -forbidden itemsets

---

```

1: procedure FBIMINER( $\mathcal{D}\downarrow, I \subseteq \mathcal{I}, \tau$ )
2:   FBI  $\leftarrow \emptyset$ 
3:   for all  $i \in \mathcal{I}$  occurring in  $\mathcal{D}$  in reverse order do
4:      $J \leftarrow I \cup \{i\}$ 
5:     if not isGenerator( $J$ ) then
6:       continue
7:     storeGenerator( $J$ )
8:     if  $|J| > 1$  &  $\text{freq}(J, \mathcal{D}) \leq \frac{2}{\tau} - 2\sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}} - 1$  then
9:       if a subset of  $J$  has been pruned then
10:        continue
11:       if  $\text{lift}(J, \mathcal{D}) \leq \tau$  then
12:         FBI  $\leftarrow \text{FBI} \cup \{J\}$ 
13:       if  $\frac{|\mathcal{D}|}{\tau} > \min_{S \subseteq J} \{\text{supp}(S, \mathcal{D}) \times \text{supp}(J \setminus S, \mathcal{D})\}$  then
14:         continue
15:       if  $\text{supp}(J, \mathcal{D}) < \frac{|\mathcal{D}|}{\sigma_J^{\max} \times \tau}$  then
16:         continue
17:        $\mathcal{D}\downarrow[i] \leftarrow \emptyset$ 
18:       for all  $j \in \mathcal{I}$  in  $\mathcal{D}$  such that  $j > i$  do
19:          $C \leftarrow \text{cov}(\{i\}, \mathcal{D}) \cap \text{cov}(\{j\}, \mathcal{D})$ 
20:         if  $\text{supp}(J, \mathcal{D}) - |C| \geq (1/\tau) - (\sigma_J^{\max}/|\mathcal{D}|)$  then
21:           if  $|C| > 0$  then
22:              $\mathcal{D}\downarrow[i] \leftarrow \mathcal{D}\downarrow[i] \cup \{(j, C)\}$ 
23:       FBI  $\leftarrow \text{FBI} \cup \text{FBIMINER}(\mathcal{D}\downarrow[i], J, \tau)$ 
24:   return FBI

```

---

The proposition tells that for small values of  $\tau$ ,  $\tau$ -forbidden itemsets are at most approximately  $\tau/4$ -frequent and that itemsets whose frequency exceeds  $\frac{2}{\tau} - 2\sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}} - 1$  cannot be  $\tau$ -forbidden. This result can be used to obtain an initial estimate for  $\tau$ : Clearly, this bound should be at least 1, yet not too high, such that frequent itemsets cannot be forbidden.

### C. Forbidden Itemset Mining Algorithm

We now present an algorithm, FBIMINER, for mining  $\tau$ -forbidden itemsets in a dataset  $\mathcal{D}$ . That is, the algorithm computes  $\text{dirty}(\mathcal{D}, \mathcal{L}, q)$  for the language and predicate defined earlier. The algorithm is based on the well-known Eclat algorithm for frequent itemset mining [29]. Here, we only describe the main differences with Eclat. The pseudo-code of FBIMINER is shown in Alg. 1. The algorithm is initially called with  $\mathcal{D}\downarrow$  ( $\mathcal{D}$  in vertical data layout),  $I = \emptyset$  and the lift threshold  $\tau$ . Just like Eclat, FBIMINER employs a depth-first search of the itemset lattice (for loop line 3 – 23, and recursive call on line 23) using set intersections of the covers of the items to compute the support of an itemset (line 19). When expanding an itemset  $I$  in the search tree (line 4), new itemsets are generated by extending  $I$  with all items in the dataset that occur in the objects in  $I$ 's cover (lines 18 – 22). Furthermore, these items are added according to some total ordering on the items, i.e., items are only added when they come after each item in  $I$  (line 18). Items are ordered by ascending support, as this is known to improve efficiency.

A first challenge is to tweak the Eclat algorithm such that the lift of itemsets can be computed. Observe that the lift of an itemset is dependent on the support of *all* of its subsets. For this purpose, we use the same depth-first traversal as Eclat,

but traverse it in reverse order (line 3). Indeed, such a reverse pre-order traversal of the search space visits all subsets of an itemset  $J$  before visiting  $J$  itself [34]. This is exactly what is required to compute the lift measure, provided that the support of each processed itemset is stored. However, Eclat generates a candidate itemset based on only two of its subsets [29]. Hence, the supports of all subsets of an itemset are not immediately available in the search tree. To remedy this, we store the support of the processed itemsets using a prefix tree, for time- and memory-efficient lookup during lift computation.

Having integrated lift computation in the algorithm, we next turn to our pruning and optimization strategies. We deploy four pruning strategies (lines 9, 13, 15 and 20). The first strategy (line 9) applies to itemsets  $J$  for which the lift cannot be computed. This happens when some of its subsets are pruned away in an earlier step. Since itemsets are only pruned when *none* of their supersets can be  $\tau$ -forbidden, this implies that  $J$  cannot be  $\tau$ -forbidden. The absence of subsets is detected when the lift computation requests the support of a subset that is not stored. Our pruning then implies that  $J$  and its supersets cannot be  $\tau$ -forbidden and thus can be pruned (line 7).

The second pruning strategy (line 13) applies to itemsets  $J$  for which we have been able to compute their lift. Indeed, when  $\frac{|\mathcal{D}|}{\tau} > \min_{S \subseteq J} \{\text{supp}(S, \mathcal{D}) \times \text{supp}(J \setminus S, \mathcal{D})\}$  then Prop. 3 tells that  $J$  cannot be a subset of a  $\tau$ -forbidden itemset. Hence, all itemsets in the tree below  $J$  are pruned (line 14).

By contrast, the third strategy (line 15) leverages Prop. 1 and skips supersets of itemsets  $J$ , regardless of whether its lift was computed. Indeed, when  $\text{supp}(J, \mathcal{D}) < \frac{|\mathcal{D}|}{\sigma_J^{\max} \times \tau}$  holds then  $J$  cannot be part of a  $\tau$ -forbidden itemset, resulting in a further pruning of the search space.

The fourth strategy employs Prop. 2 to prune extensions of  $J$  that do not cause a sufficient reduction in support. This check is performed prior to the recursive call (line 20).

Finally, we also implement an optimization that avoids certain lift computations (line 8). Only when the algorithm encounters an itemset  $J$  with at least two items and a frequency lower than the bound from Prop. 4, the lift of  $J$  is computed. All other itemsets cannot be  $\tau$ -forbidden by Prop. 4. Note that this only eliminates the need for checking the lift of certain itemsets but by itself does not lead to a pruning of its supersets.

A careful reader may have spotted the optimized pruning of non-generators on lines 5-7. Recall that as a direct consequence of Prop. 2, any  $\tau$ -forbidden itemset must be a generator, i.e., have a support which is strictly lower than that of all of its subsets. The Talky-G algorithm [35] implements specific optimizations for mining such generators, using a hash-based method that was introduced in the Charm algorithm for closed itemset mining [36]. We use the same technique in FBIMINER to efficiently prune non-generators.

During the mining process, all encountered generators are stored in a hashmap (procedure `storeGenerator` on line 7). As hash value we use, just like the Charm algorithm, the sum of the oid's of all objects in which an itemset occurs. If an itemset has the same support as one of its subsets, it is clear that this itemset must occur in exactly the same objects, and

will map to the same hash value. Moreover, the probability of unrelated itemsets having the same *sum of oids* is lower than the probability of them having the same support. Therefore this sum is taken as hash value instead of the support of itemsets.

Procedure `isGenerator` on line 5 checks all stored itemsets with the same hash value as  $J$ . If any of these itemsets are a subset of  $J$  with the same support,  $J$  is discarded as a non-generator. Furthermore, since all supersets of a non-generator are also non-generators, the entire subtree can be pruned. If no subset with identical support is discovered for an itemset  $J$ , then  $J$  is either a generator, or a subset with identical support has previously been pruned, in which case  $J$  will eventually be pruned on line 9.

## VI. REPAIRING INCONSISTENCIES

The algorithm FBIMINER discovers a set of  $\tau$ -forbidden itemsets that describe inconsistencies in  $\mathcal{D}$ . When this set is non-empty,  $\mathcal{D}$  is regarded as dirty. We next want to clean  $\mathcal{D}$ . Following the general framework outlined in section IV, we wish to compute a repair  $\mathcal{D}'$  of  $\mathcal{D}$  such that (i)  $\mathcal{D}'$  is clean, i.e., no  $\tau$ -forbidden itemsets should be found in  $\mathcal{D}'$ ; and (ii)  $\mathcal{D}'$  differs minimally from  $\mathcal{D}$ . Due to the dynamic notion of data quality, (i) becomes more challenging than in a traditional repair setting. We choose not to optimize (ii) directly by computing minimal modifications to dirty objects, but instead impute values from clean objects that differ minimally from a dirty object, in line with common practice in data imputation. We start by showing how the creation of new  $\tau$ -forbidden itemsets can be avoided by means of so-called almost forbidden itemsets (Sect. VI-A) and explain how these can be mined (Sect. VI-B), before describing the repair algorithm itself (Sect. VI-C).

### A. Ensuring a Clean Repair

Given a dataset  $\mathcal{D}$  and its set of  $\tau$ -forbidden itemsets  $\text{FBI}(\mathcal{D}, \tau)$ , we define the *dirty* objects in  $\mathcal{D}$ , denoted by  $\mathcal{D}_{\text{dirty}}$ , as those objects in  $\mathcal{D}$  that appear in the cover of an itemset in  $\text{FBI}(\mathcal{D}, \tau)$ . In other words,  $\mathcal{D}_{\text{dirty}}$  consists of all objects that support a  $\tau$ -forbidden itemset. The remaining set of clean objects in  $\mathcal{D}$  is denoted by  $\mathcal{D}_{\text{clean}}$ . The repair algorithm will produce a dataset  $\mathcal{D}'$  by modifying all objects in  $\mathcal{D}_{\text{dirty}}$  to remove the forbidden itemsets. We consider value modifications where values come from clean objects. Recall however that we need to obtain a repair  $\mathcal{D}'$  of  $\mathcal{D}$  such that  $\text{FBI}(\mathcal{D}', \tau)$  is empty, i.e., such that  $\mathcal{D}'$  is clean. We first present an example to show that this is not a trivial problem.

**Example 2.** People typically graduate High School in the year they turn 18. Depending on the timing of a census, there may be graduates who are still only 17 years old. In the Adult Census dataset, the itemset (AGE=<18, EDUCATION=HS-GRAD) is rare, with a lift  $\approx 0.072$ . Assume that  $\tau$ -forbidden itemsets were mined with  $\tau = 0.07$ . The itemset is thus not considered forbidden, and rightly so. However, an object containing (AGE=<18, EDUCATION=HS-GRAD) could, for example, contain the forbidden itemset (AGE=<18, MARITALSTATUS=DIVORCED), where MaritalStatus is in error. If the repair algorithm were to change *Age* instead, the lift of (AGE=<18, EDUCATION=HS-GRAD) would drop to  $\approx 0.068$ ! This itemset will then become  $\tau$ -forbidden in  $\mathcal{D}'$ , yielding again a dirty dataset.  $\diamond$

A naive approach for avoiding new inconsistencies would be to run FBIMINER on  $\mathcal{D}'$  for each candidate modification, and reject the modification in case  $\text{FBI}(\mathcal{D}', \tau)$  is non-empty. In view of the possibly exponential number of candidate modifications, this approach is not feasible for all but the smallest datasets. As an alternative, we propose to compute up front enough information to ensure cleanliness of multiple repairs. More specifically, the procedure A-FBIMINER computes a set  $\mathbb{A}$  of *almost  $\tau$ -forbidden itemsets*, i.e., itemsets that could become  $\tau$ -forbidden after a given number of modifications  $k$ . This computation relies only on the dataset  $\mathcal{D}$  and its dirty objects, and not on the particular modifications made during repairing.

$\triangleright$  **Almost Forbidden Itemsets.** Almost forbidden itemsets are mined by algorithm A-FBIMINER. It mines itemsets  $J$  similarly to FBIMINER, but uses a relaxed notion of lift, called the *minimal possible lift of  $J$  after  $k$  modifications*, to be explained below. More specifically, by using the minimal possible lift measure, A-FBIMINER returns a set  $\mathbb{A}$  of itemsets such that for *any* dataset  $\mathcal{D}'$  obtained from  $\mathcal{D}$  by at most  $k$  modifications,

$$\text{FBI}(\mathcal{D}', \tau) \subseteq \mathbb{A}. \quad (\dagger)$$

We next explain what precisely  $\mathbb{A}$  consists of and how it can be used to avoid new inconsistencies whilst repairing. It is crucial for our approach that  $\mathbb{A}$  accommodates for *any* repair  $\mathcal{D}'$  of  $\mathcal{D}$  obtained from at most  $k$  modifications as it eliminates the need for considering all possible repairs one by one.

$\triangleright$  **Minimal Possible Lift.** To define the relaxed lift measure used by A-FBIMINER, we consider the following problem: *Given an itemset  $J$  in  $\mathcal{D}$  and its lift( $J, \mathcal{D}$ ), can  $J$  become  $\tau$ -forbidden after  $k$  modifications to  $\mathcal{D}$  have been made?* Let us first analyze the minimal possible lift in case a single modification is made. Suppose that  $\text{lift}(J, \mathcal{D}) = |\mathcal{D}| \times \text{supp}(J, \mathcal{D}) / (\text{supp}(I, \mathcal{D}) \times \text{supp}(J \setminus I, \mathcal{D}))$  for some  $I \subset J$ , with  $\text{supp}(I, \mathcal{D}) \leq \text{supp}(J \setminus I, \mathcal{D})$ . It can be shown that, after one modification, the minimal possible lift is either

$$\frac{|\mathcal{D}| \times (\text{supp}(J, \mathcal{D}) - 1)}{\text{supp}(I, \mathcal{D}) \times (\text{supp}(J \setminus I, \mathcal{D}) - 1)}$$

or

$$\frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{(\text{supp}(I, \mathcal{D}) + 1) \times \text{supp}(J \setminus I, \mathcal{D})}.$$

Which case is minimal depends on the ratio of the supports of the itemsets  $I$ ,  $J \setminus I$  and  $J$ . Nevertheless, given these supports we can return the smallest of the two as minimal possible lift after one modification and denote the result by  $\text{mpl}(J, I, 1)$ . To generalize this to an arbitrary number  $k$  of modifications and obtain  $\text{mpl}(J, I, k)$ , we recursively repeat this computation  $k$  times, with updated supports of the itemsets  $I$ ,  $J$  and  $J \setminus I$ . The crucial property of minimal possible lift is the following:

**Proposition 5.** *Let  $J$ ,  $I$  as above. If  $J \in \text{FBI}(\mathcal{D}', \tau)$  for some  $\mathcal{D}'$  obtained from  $\mathcal{D}$  by at most  $k$  modifications, then  $\text{mpl}(J, I, k) \leq \tau$ .  $\square$*

Hence, if  $\mathbb{A}$  consists of all itemsets  $J$  for which  $\text{mpl}(J, I, k) \leq \tau$ , for some  $I \subset J$  such that  $\text{lift}(J, \mathcal{D}) = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\text{supp}(I, \mathcal{D}) \times \text{supp}(J \setminus I, \mathcal{D})}$ , then it is guaranteed that all itemsets in  $\text{FBI}(\mathcal{D}', \tau)$  are returned, as desired by property  $(\dagger)$ . Algorithm

A-FBIMINER mines all itemsets  $J$  for which  $\text{mpl}(J, I, k) \leq \tau$ , for  $I$  as above, and thus returns  $\mathbb{A}$ .

▷ **Avoiding New Inconsistencies.** Before explaining how A-FBIMINER works, we first explain how the set  $\mathbb{A}$  of almost forbidden itemsets can be used to guarantee clean repairs. We come back to the repair algorithm in Sect. VI-C. Consider a modification  $o_{\text{rep},i}$  of a dirty object  $o_i$ . Our repair algorithm rejects this modification in the following cases:

— *Old inconsistency:* These are itemsets which should not be present in the repaired dataset, as they are already known to be inconsistent in  $\mathcal{D}$ . This happens if  $o_{\text{rep},i}$  covers an itemset  $C \in \mathbb{A} \cap \text{FBI}(\mathcal{D}, \tau)$ . It also happens if  $o_{\text{rep},i}$  covers an itemset  $C \in \mathbb{A}$  with  $\text{supp}(C, \mathcal{D}) = 0$ , which are itemsets that do not occur in the original dataset, but would be forbidden if they did occur. Indeed, an incautious repair could introduce such an itemset in the “repaired” dataset. In other words, no itemset should be present in  $o_{\text{rep},i}$  that is already known to be forbidden. We denote this set as  $\mathbb{A}_{\text{old}}$ .

**Repair Safety:** Old inconsistencies are avoided when a repair  $o_{\text{rep},i}$  does not support any itemsets in  $\mathbb{A}_{\text{old}}$ .

— *Potential inconsistency:* Object  $o_{\text{rep},i}$  covers an itemset  $C \in \mathbb{A}$  with  $\text{lift}(C, \mathcal{D}) > \tau$  and  $\text{lift}(C, \mathcal{D}'_i) < \text{lift}(C, \mathcal{D})$ , where  $\mathcal{D}'_i$  is  $\mathcal{D}$  in which only  $o_i$  is replaced by  $o_{\text{rep},i}$  and all other objects in  $\mathcal{D}$  are preserved. In other words, when  $o_{\text{rep},i}$  covers an almost  $\tau$ -forbidden itemset that is not yet  $\tau$ -forbidden in  $\mathcal{D}$ , modifications that reduce the lift of this itemset should be prevented. We denote this set as  $\mathbb{A}_{\text{pot}}$ .

**Repair Safety:** Since it is infeasible to recompute the lift of all itemsets in  $\mathbb{A}_{\text{pot}}$ , we opt for a more efficient method which suffices to ensure that the lift of the itemsets  $I \in \mathbb{A}_{\text{pot}}$  doesn't drop, by asserting that (i) no occurrence of  $I$  has been removed (which would decrease the numerator in its lift) and (ii) no occurrence of a strict subset of  $I$  has been added (which would increase the denominator in its lift). By guaranteeing that  $\text{supp}(I, \mathcal{D}) \geq \text{supp}(I, \mathcal{D}')$  and for all  $J \subsetneq I : \text{supp}(J, \mathcal{D}) \leq \text{supp}(J, \mathcal{D}')$ , it follows that  $\text{lift}(I, \mathcal{D}') \geq \text{lift}(I, \mathcal{D})$ .

It can be shown that these two safety checks suffice to guarantee that no forbidden itemsets will be found in accepted repairs. We declare a candidate repair to be safe if it passes both checks.

## B. Mining Almost Forbidden Itemsets

We now describe algorithm A-FBIMINER. It is similar to FBIMINER, except for the relaxed lift measure and different pruning strategies. Recall that algorithm A-FBIMINER is to mine almost forbidden itemsets without looking at repairs, i.e., only  $\mathcal{D}$  and an upper bound  $k$  on the number of modified objects is available. Clearly  $k$  is at most  $|\mathcal{D}_{\text{dirty}}|$ . To adapt the pruning strategies from FBIMINER, the underlying properties must be revised to take possible modifications into account. Since clean objects are never modified, a tight bound on the support of an itemset in any  $\mathcal{D}'$ , obtained from  $\mathcal{D}$  by at most  $k$  modifications, can be obtained. Indeed, observe that for any itemset  $I$ , the following holds:

$$\text{supp}(I, \mathcal{D}_{\text{clean}}) \leq \text{supp}(I, \mathcal{D}') \leq \text{supp}(I, \mathcal{D}_{\text{clean}}) + k$$

An immediate consequence is that A-FBIMINER must also consider itemsets  $I$  with  $\text{supp}(I, \mathcal{D}_{\text{clean}}) = 0$  as these can become supported in repairs. Furthermore, we can now modify Propositions 1 and 2:

**Proposition 6.** *For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , if  $J$  is a  $\tau$ -forbidden itemset in  $\mathcal{D}'$ , then we have that  $\text{supp}(I, \mathcal{D}_{\text{clean}}) \geq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{\text{max}} \times \tau} - k$ .  $\square$*

**Proposition 7.** *For any three itemsets  $I$ ,  $J$  and  $K$  such that  $I \subset J \subset K$ , if  $K$  is a  $\tau$ -forbidden itemset in  $\mathcal{D}'$ , then  $\text{supp}(I, \mathcal{D}_{\text{clean}}) - \text{supp}(J, \mathcal{D}_{\text{clean}}) \geq \frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}^{\text{max}}}{|\mathcal{D}|} - k$ .  $\square$*

Similarly to the pruning strategies for FBIMINER, we again use the trivial lower bound  $\text{supp}(J, \mathcal{D}') = 1$ . Furthermore, to make use of these propositions for pruning, note that we do not know  $\sigma_{I, \mathcal{D}'}^{\text{max}}$ . Indeed, recall that we do not consider any particular repair  $\mathcal{D}'$ . Instead,  $\sigma_{I, \mathcal{D}_{\text{clean}}}^{\text{max}}$  can be computed, and it holds that  $\sigma_{I, \mathcal{D}'}^{\text{max}} \leq \sigma_{I, \mathcal{D}_{\text{clean}}}^{\text{max}} + k$ . As a consequence, Prop. 6 is used to prune supersets of  $I$  whenever  $\text{supp}(I, \mathcal{D}_{\text{clean}}) < \frac{|\mathcal{D}|}{(\sigma_{I, \mathcal{D}_{\text{clean}}}^{\text{max}} + k) \times \tau} - k$ .

From Prop. 7, it follows that non-generator pruning can be applied to an itemset  $I$  if  $\frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}^{\text{max}}}{|\mathcal{D}|} > k$ . Observe that  $0 < \frac{\sigma_{I, \mathcal{D}'}^{\text{max}}}{|\mathcal{D}|} \leq 1$  and hence the impact of this ratio is almost negligible. Therefore, instead of computing  $\sigma_{I, \mathcal{D}'}^{\text{max}}$  for every  $I$ , we use an estimate for  $\sigma_{\emptyset, \mathcal{D}'}^{\text{max}}$  instead, i.e.,  $\sigma_{\emptyset, \mathcal{D}_{\text{clean}}}^{\text{max}} + k$ . Prop. 7 implies that  $\text{supp}(I, \mathcal{D}_{\text{clean}}) - \text{supp}(J, \mathcal{D}_{\text{clean}}) \geq \frac{1}{\tau} - \frac{(\sigma_{\emptyset, \mathcal{D}_{\text{clean}}}^{\text{max}} + k)}{|\mathcal{D}|} - k$  must hold for every  $I \subseteq J$  for which  $J$  is  $\tau$ -forbidden in any  $\mathcal{D}'$  obtained by using  $k$  modifications.

Finally, Prop. 3 also needs to be adapted to account for possible modifications. Since the  $\text{mpl}$  measure depends on the ratio of  $\text{supp}(J, \mathcal{D}')$  and its partitions, it does not preserve the anti-monotonicity of the lift denominator. Instead, we need to compute the worst case increase in the denominator of the lift measure:

**Proposition 8.** *For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , it holds that  $\text{lift}(J, \mathcal{D}') \geq$*

$$\frac{\text{supp}(J, \mathcal{D}') \times |\mathcal{D}|}{\min_{S \subset I} \{(\text{supp}(S, \mathcal{D}_{\text{clean}}) + k) \times (\text{supp}(I \setminus S, \mathcal{D}_{\text{clean}}) + k)\}}. \quad \square$$

As before, we use this proposition for  $\text{supp}(J, \mathcal{D}') = 1$ . These three properties allow substantial pruning when mining almost forbidden itemsets similarly as explained for FBIMINER.

▷ **Batch execution.** Propositions 6 and 7 also show that the number  $k$  of modifications considered has a direct impact on the pruning capabilities of algorithm A-FBIMINER. Indeed, suppose that  $k$  takes the maximal value, i.e.,  $k = |\mathcal{D}_{\text{dirty}}|$ . Then, it is likely that the bounds given in these propositions do not allow for any pruning. Worse still, the minimal possible lift, which also depends on  $k$ , will declare many itemsets to be almost forbidden. Although A-FBIMINER will need to be run only once to obtain the set  $\mathbb{A}$  and all dirty objects can be repaired based on  $\mathbb{A}$ , this set will be big and inefficient to compute (due to lack of pruning). On the other hand, when  $k = 1$ , pruning will be possible and  $\mathcal{A}$  will be of reasonable size. However, to clean all dirty objects, we need to deal with

them one-at-a-time, re-running A-FBIMINER for  $k = 1$  after a single dirty object is repaired.

Between these two extremes, we propose to process the dirty objects in batches. More specifically, we partition  $\mathcal{D}_{dirty}$  into blocks of a size  $r$ , optimizing the trade-off between the runtime of A-FBIMINER and the number of runs. The question, of course, is what block size to select. We already described block size  $r = 1$  and  $r = |\mathcal{D}_{dirty}|$ . We next use Prop. 6 and Prop. 7 to identify ranges for  $r$  for which pruning may still be possible.

Let  $I$  and  $J$  be itemsets such that  $I \subset J$ . Proposition 6 is only applicable if  $\frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{max} \times \tau} > r$ , and Prop. 7 is only applicable if  $\frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}^{max}}{|\mathcal{D}|} > r$ , where  $\mathcal{D}'$  is now any repair obtained from  $\mathcal{D}$  by at most  $r$  modifications. It is easy to see that  $\frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{max} \times \tau} \geq \frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}^{max}}{|\mathcal{D}|}$  and hence  $r = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{max} \times \tau}$  is the maximal block size for which one can expect pruning. As before letting  $\text{supp}(J, \mathcal{D}') = 1$  and  $I = \emptyset$ , we identify  $r = \frac{|\mathcal{D}|}{(\sigma_{\emptyset, \mathcal{D}_{clean}}^{max} + r) \times \tau}$  as a maximal block size that allows pruning based on Prop. 6. Additional pruning based on Prop. 7 is possible for lower block sizes, i.e., when  $\frac{1}{\tau} - 1 > r$ . Here we again use that  $\frac{|\mathcal{D}|}{\sigma_{I, \mathcal{D}'}^{max}} \geq 1$ . We hence also identify  $r = \frac{1}{\tau} - 1$  as block size that allows substantial pruning.

Of course, there is no universally optimal block size  $r$ . The specifics of the data and even the choice of which objects to include in a partition of  $r$  objects all impact the pruning power. It is also important to note that the block sizes derived above guarantee that the associated propositions are *applicable*, but still offer *reduced pruning power* in comparison to smaller block sizes. In the experimental section, we show that  $r = \frac{1}{2\tau}$  provides a sensible default value of  $r$ , whereas  $r = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{max} \times \tau}$  improves the runtime on datasets with a small number of attributes.

### C. Repair Algorithm

We are finally ready to describe our algorithm REPAIR, shown in Alg. 2. It takes as input the dirty and clean objects,  $\mathcal{D}_{dirty}$  and  $\mathcal{D}_{clean}$  respectively, a similarity function, and the lift threshold  $\tau$ . The dirty objects  $\mathcal{D}_{dirty}$  are partitioned in blocks  $R_i$  of size  $r$ . Per block, the set of almost forbidden itemsets  $\mathbb{A}_i$  is computed. At this point, the repair process depends only on the two sets  $R_i$  and  $\mathbb{A}_i$ . After each processed block,  $\mathcal{D}$  is updated with the found repairs in  $\mathcal{D}'$ , denoted as  $\mathcal{D} \oplus \mathcal{D}'$ . By default, the set  $\mathcal{D}_{clean}$  is not altered during the entire repair process, ensuring that subsequent repairs are not based on previously imputed values, to avoid the propagation of undesirable repair choices.

For each dirty object in  $R_i$ , a candidate repair is generated by replacing some of its items by items from a similar, but clean object in  $\mathcal{D}_{clean}$ . By using the most similar objects to produce candidate repairs, we try to minimize the difference between  $\mathcal{D}'$  and  $\mathcal{D}$ . This approach is in line with the commonly used *hot deck imputation* in statistical survey editing [4].

If the candidate repair is safe (as explained before) with respect to the almost forbidden itemsets in  $\mathbb{A}_i$ , then it is added to the set  $\mathcal{D}'$  (line 12). Otherwise, the next candidate clean

---

### Algorithm 2 Repairing dirty objects

---

```

1: procedure REPAIR( $\mathcal{D}_{dirty}$ ,  $\mathcal{D}_{clean}$ ,  $\text{linsim}$ ,  $\tau$ )
2:   for all  $R_i \in \text{blocks}(\mathcal{D}_{dirty}, r)$  do
3:      $r := |R_i|$ 
4:      $\mathcal{D}' := \emptyset$ ;  $\mathcal{D}'' = \emptyset$ 
5:      $\mathbb{A}_i = \text{A-FBIMINER}(\mathcal{D} \oplus \mathcal{D}', r, \tau)$ 
6:     for all  $o_i \in R_i$  do
7:       success := false
8:       for all  $o_c \in \mathcal{D}_{clean}$  in  $\text{sim}(o_c, o_i)$  desc. order do
9:          $o_{rep, i} = \text{MODIFY}(o_i, o_c)$ 
10:        if SAFE( $o_i, o_{rep, i}, \mathbb{A}_i$ ) then
11:          success := true
12:           $\mathcal{D}' := \mathcal{D}' \cup o_{rep, i}$ 
13:          break
14:        if not success then  $\mathcal{D}'' = \mathcal{D}'' \cup o_i$ 
15:     return ( $\mathcal{D}', \mathcal{D}''$ )

```

---

object is considered (loop line 8–13) until a repair is found (line 13) or all candidate repairs have been rejected. In this case,  $o_i$  is added to the set of unrepaired objects  $\mathcal{D}''$  (line 14), for further user inspection.

It remains to explain how candidate repairs are generated. For each dirty object  $o_i$ , the algorithm consecutively processes the clean objects  $o_c$  in order of their similarity to  $o_i$ . The algorithm subsequently modifies the dirty object  $o_i$  by means of the procedure  $\text{MODIFY}(o_i, o_c)$  (line 8). The resulting object is denoted by  $o_{rep, i}$ . In our implementation,  $\text{MODIFY}(o_i, o_c)$  replaces items  $(A, v)$  in  $o_i$  by  $(A, o_c[A])$  that occur in the  $\tau$ -forbidden itemsets covered by  $o_i$ , i.e., only the items that are part of inconsistencies are modified.

## VII. USER INTERACTION

The cleaning methods outlined in this paper have been designed such that they can be run fully automatically, without any user input or interaction. Of course, in practice, optional user input is often desired. Such a mechanism can readily be integrated into our method. Indeed, the repairing process relies only on the sets of forbidden and almost forbidden itemsets,  $\text{FBI}(\mathcal{D}, \tau)$  and  $\mathbb{A}$ , respectively. A user can validate the discovered itemsets, answering the simple question “are these items allowed to co-occur?”. Itemsets that are rejected by the user can be discarded from the respective sets. The algorithms will then work as desired, considering the user-rejected (almost) forbidden itemsets to be semantically correct. Likewise, a user could be shown the top-k lowest lift forbidden itemsets and asked to confirm which itemsets to remove. The efficient discovery of such top-k itemsets and experimental evaluation of user interactions are left for future work.

## VIII. EXPERIMENTS

In this section, we experimentally validate our proposed techniques, by answering the following questions:

- Does FBIMINER find a manageable set of errors efficiently? Are the forbidden itemsets actually errors? What is the impact of pruning?
- Can almost forbidden itemsets be mined efficiently? How does the block size impact runtime?



- Is the repair algorithm able to find low-cost repairs efficiently? How often is it impossible to repair an object?

### A. Experimental Settings

The experiments were conducted on real-life datasets from the UCI repository (<http://archive.ics.uci.edu/ml/>). We show results for six datasets, their descriptions are given in Table II. The Adult database was preprocessed by discretizing ages and removing other continuous attributes. The algorithms have been implemented in C++, the source code and used datasets are available for research purposes<sup>1</sup>. The program has been tested on an Intel Xeon Processor (2.9GHZ) with 32GB of memory running Linux. Our algorithms run in main memory. Reported runtimes are an average over five independent runs.

Table II. STATISTICS OF THE UCI DATASETS USED IN THE EXPERIMENTS. WE REPORT THE NUMBER OF OBJECTS, DISTINCT ITEMS, AND ATTRIBUTES.

Dataset	$ \mathcal{D} $	$ \mathcal{I} $	$ \mathcal{A} $
Adult	48842	202	11
CensusIncome	199524	235	12
CreditCard	30000	216	12
Ipums	70187	364	32
LetterRecognition	20000	282	17
Mushroom	8124	119	23

### B. Forbidden Itemset Mining

We ran the forbidden itemset mining algorithm FBMINER with full pruning, reporting the total runtime, the number of forbidden itemsets, and the number of objects containing a forbidden itemset, for increasing values of  $\tau$ . For the larger datasets, Ipums and CensusIncome, a smaller  $\tau$  range was considered. This prevents an explosion in the number of forbidden itemsets and the associated high runtime. The results are shown in Fig. 2, Fig. 3 and Fig. 4, respectively.

The results show that the runtime of the algorithm (Fig. 2) scales linearly with  $\tau$ . As a result of the depth-first search, the runtime is strongly influenced by the number of distinct items. As a consequence, the algorithm runs slowest on the Ipums dataset. The runtime on the LetterRecognition dataset is explained by its relatively high number of items, and the fact that it contains many forbidden itemsets.

The number of forbidden itemsets (Fig. 3) is typically small, although there is a stronger than linear increase as the lift threshold increases, illustrating that  $\tau$  should indeed be chosen very small. Especially for the LetterRecognition database, the number of forbidden itemsets increases exponentially. This occurs because the dataset is very noisy, since the contained letters were randomly distorted. In contrast, the less noisy Adult and CensusIncome datasets have relatively few dirty objects. The number of dirty objects (Fig. 4) naturally follows a similar pattern to the number of forbidden itemsets, with an occasionally big increase if a forbidden itemset with a relatively high support is discovered.

To answer the question “Are the forbidden itemsets actually errors?”, a gold standard for the subjective task of data cleaning

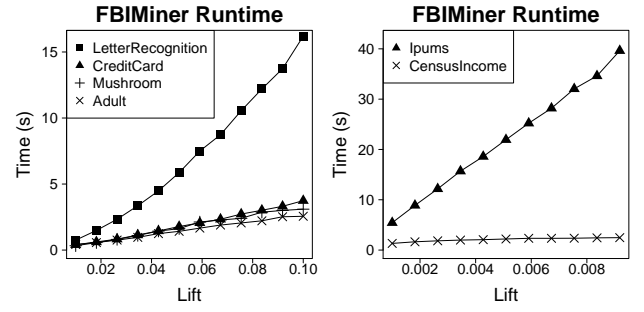


Figure 2. Runtime of FBMINER in function of maximum lift threshold  $\tau$ .

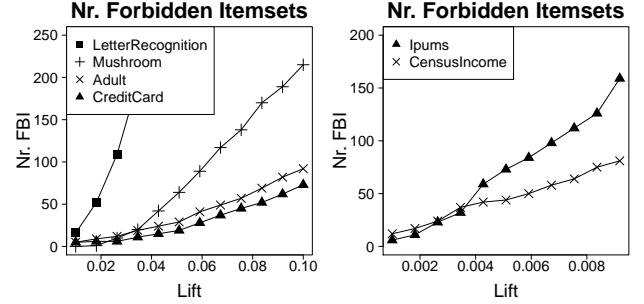


Figure 3. Number of Forbidden Itemsets in function of maximum lift threshold  $\tau$ .

would be needed. Although synthetic error generators exist [37], [38], they require the constraints to be known up front. In line with [10], we therefore evaluate forbidden itemsets manually for usefulness, obtaining only a precision score. The results for the Adult dataset, which is the most readily interpretable, are shown in Table III. Precision is very high for small  $\tau$ -values, and keeps up around 50% in the  $\tau$ -range, which is a high score for an uninformed method. We believe a high precision is important to instill faith in an eventual user. Note that we report the precision of the *forbidden itemsets*, the number of erroneous objects is a multiple of this value, as illustrated by Fig. 4.

In order to evaluate the influence of the pruning strategies, we report the number of itemsets processed with only one type of pruning enabled, and contrast these with the number of itemsets processed when all pruning is enabled. We distinguish between *Min. Supp* pruning, using Prop. 1 on line 15 of Alg. 1; *Lift Denominator* pruning, using Prop. 3 on line 13; and *Support Diff* pruning, using Prop. 2 on line 20.

The results are shown in Fig. 5 for the Adult and Credit-Card datasets; results for CensusIncome and LetterRecognition were similar. On the Mushroom and Ipums datasets, which have many attributes, FBMINER became infeasible for larger values of  $\tau$  without Support Diff pruning. Clearly, Support Diff pruning is dominant in most cases. Since this strategy also entails non-generator pruning, it is definitely crucial for the runtime of FBMINER. Especially as  $\tau$  increases, the other strategies also improve the overall result, indicating that all are beneficial and complementary to each other. We do not show the results when all pruning is disabled since all itemsets are then considered (independently of  $\tau$ ), leading to a high number of processed itemsets and running time.

A separate issue is the maximal frequency of a forbidden

<sup>1</sup><http://adrem.uantwerpen.be/joerirammelaere>

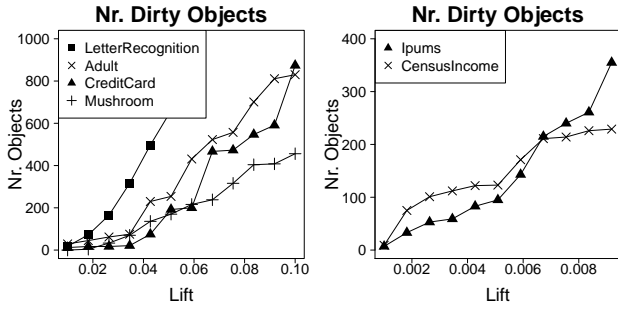


Figure 4. Number of objects containing one or more Forbidden Itemsets in function of maximum lift threshold  $\tau$ .

Table III. PRECISION OF DISCOVERED FORBIDDEN ITEMSETS ON ADULT DATASET.

	$\tau$ -value					
Dataset	0.01	0.026	0.043	0.067	0.084	0.1
Nr. FBI	5	12	24	49	69	92
Precision	100%	67%	71%	61%	55%	45%

itemset, used on line 8 of Alg. 1. Recall that this is not a pruning strategy: using the frequency bound increased the number of itemsets processed, but may reduce runtime by avoiding certain unnecessary lift computations. This bound was disabled for the previous pruning results, to prevent painting a distorted picture of the influence of each pruning strategy. Table IV shows the runtime influence of the frequency bound as a percentage of the runtime without this bound. Results range from a 20% decrease to a 10% increase, and indicate that the frequency bound is typically beneficial for the runtime of FBIMiner, especially for smaller  $\tau$ -values.

### C. Almost Forbidden Itemsets

The discovery of almost forbidden itemsets is the most computationally expensive part in our methodology. Recall that the runtime of A-FBIMINER depends both on the lift threshold  $\tau$  and the number of dirty objects as discovered by FBIMiner. Since a larger  $\tau$  automatically entails a higher number of dirty objects, clearly scalability in  $\tau$  is an issue.

For each dataset and each  $\tau$ , we first run algorithm FBIMINER to obtain the forbidden itemsets. Let  $k$  denote the number of dirty objects found. We then run algorithm A-FBIMINER a number of times with block size  $r$ , indicating the number of dirty objects to be repaired at once, until all  $k$  objects have been repaired. Firstly, we report runtimes for the extreme cases of the block size, i.e.,  $r = 1$  and  $r = k$ . These runtimes are shown in Fig. 6a-6b for the first four datasets. Results for CensusIncome and Ipums were similar, the plots are deferred to the appendix.

The difference between both block sizes is clear. For  $r = k$ , runtimes start out reasonably low, but quickly explode as the algorithm loses its pruning power and computation becomes infeasible. This is the most problematic for Mushroom, which has a larger number of attributes, and LetterRecognition, which has a very high number of dirty objects  $k$ . Block size  $r = 1$  remains feasible throughout the  $\tau$  range, but is slower overall.

Next, we focus on the optimal block size  $r$ . As outlined in Sect. VI-B, we can identify the quantities  $\frac{1}{\tau} - 1$  and

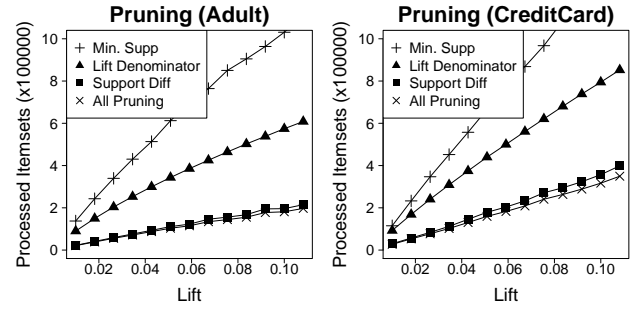


Figure 5. Number of itemsets processed with different pruning strategies in function of maximum lift threshold  $\tau$ .

Table IV. RUNTIME INFLUENCE OF MAXIMAL FREQUENCY BOUND IN FUNCTION OF  $\tau$ . VALUES SHOWN ARE THE RUNTIME WITH FREQUENCY BOUND AS A PERCENTAGE OF RUNTIME WITHOUT FREQUENCY BOUND.

	$\tau$ -value					
Dataset	0.01	0.026	0.043	0.067	0.084	0.1
Adult	95%	92%	93%	98%	98%	112%
CreditCard	95%	89%	77%	96%	110%	97%

$\frac{1}{\tau} \times \frac{|D|}{\sigma_{\theta, D}^{max}}$  as the maximal block sizes for which Prop. 6 and Prop. 7, respectively, are still applicable. Fig. 6c-6d displays the obtained runtimes using these block sizes on the first four datasets, results for CensusIncome and Ipums are again omitted but similar. Note that the chosen values for  $r$  are  $\tau$ -dependent. Consequently, for every  $\tau$ -value, a different number of dirty objects is obtained and partitioned into blocks of a *different* absolute size. As expected, the runtimes are lower than for  $k = 1$ , while feasibility is better than for  $r = k$ , proving that the right block size indeed improves the overall performance of the algorithm. Runtime on the LetterRecognition dataset naturally suffers from the exponential increase in the number of dirty objects on that dataset. However, performance on the Mushroom dataset is still problematic: the number of attributes leads to a deep search tree, and pruning power is too limited. The same holds for Ipums. Plots of the obtained runtimes are deferred to the appendix.

As an alternative, we consider the block size  $r = \frac{1}{2\tau}$ , the halfway point between  $r = 1$  and  $r = \frac{1}{\tau} - 1$ . Figure 6e-6f shows that this block size provides sufficient pruning power for the Mushroom dataset, and indeed outperforms all other considered sizes over the entire  $\tau$ -range. For higher  $\tau$ -values, the algorithm still struggles on the Ipums dataset, its high number of items proving problematic. On the other datasets, A-FBIMINER is fast for low values of  $\tau$ , and feasible across the considered  $\tau$ -range.

### D. Data Repairing

In Table V we report on the quality of repairs obtained by algorithm REPAIR for various values of  $\tau$ . The block size  $r = \frac{1}{2\tau}$  was chosen, as described in the previous paragraph. We report the minimal and maximal similarity between a dirty object and its repair (within the  $\tau$  range as above), with a similarity value of 1 indicating identical objects. The obtained repairs consistently have a high similarity in the given  $\tau$ -range.

We also report the number of objects that could not be repaired at the highest  $\tau$ -value, denoted as  $D''$ . For Adult,

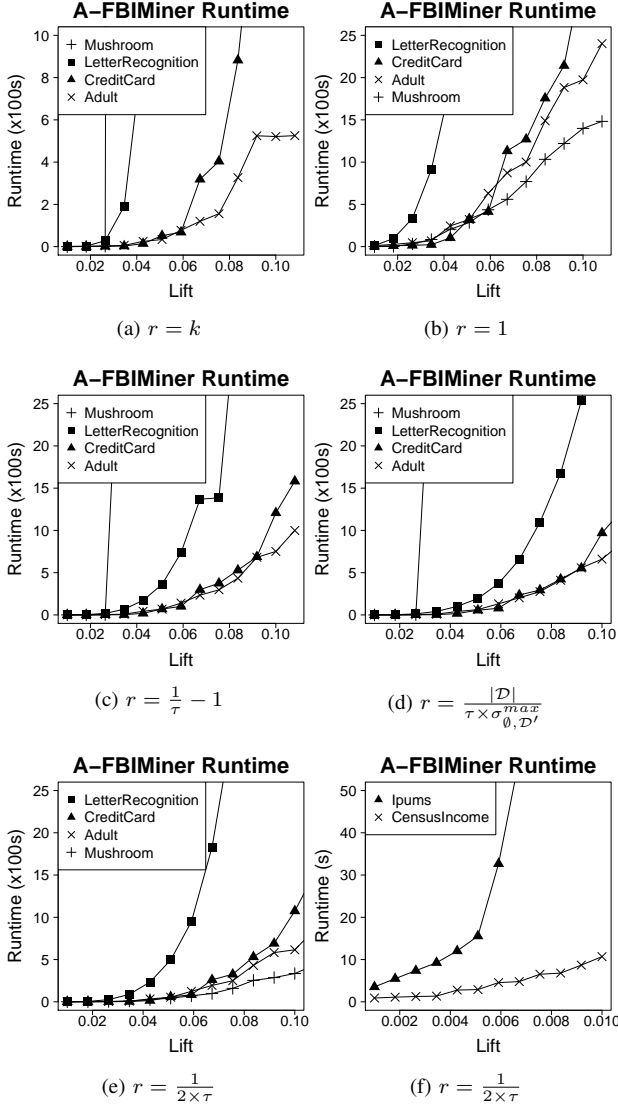


Figure 6. Runtime of A-FBIMINER in function of maximum lift threshold  $\tau$ , for various block sizes  $r$ .

Clean	Age	Marital-status	Relationship	Sex	...
X	<18	Married-civ-spouse	Husband	Male	
V	>50	Married-civ-spouse	Husband	Male	
X	<18	Married-civ-spouse	Own-child	Female	
V	<18	Never-married	Own-child	Female	
X	22-30	Married-civ-spouse	Not-in-family	Male	
V	22-30	Never-married	Own-child	Male	
X	22-30	Married-civ-spouse	Wife	Male	
V	22-30	Married-civ-spouse	Wife	Female	
X	31-50	Married-civ-spouse	Wife	Male	
V	31-50	Married-civ-spouse	Husband	Male	

Figure 7. Example repairs on Adult dataset.

CensusIncome, CreditCard and LetterRecognition, only a few objects are unrepairable and this occurs only for high values of  $\tau$ . A higher number of unrepairable objects is encountered for the Ipums and Mushroom datasets. This seems to suggest that a higher number of attributes causes problems for repairing.

Finally, Fig. 8 shows the runtime of algorithm REPAIR. The reported running times exclude the time needed for A-FBIMINER. Since the repair algorithm computes nearest

Table V. AVERAGE QUALITY OF REPAIRS.

Dataset	$\tau$ -range	Min-Max Sim.	$ \mathcal{D}' $
Adult	0.01-0.1	0.94-0.95	1
CensusIncome	0.001-0.01	0.90-0.95	0
CreditCard	0.01-0.1	0.94-0.96	10
Ipums	0.001-0.01	0.95-0.98	94
LetterRecognition	0.01-0.1	0.96-0.98	33
Mushroom	0.01-0.1	0.94-0.99	238

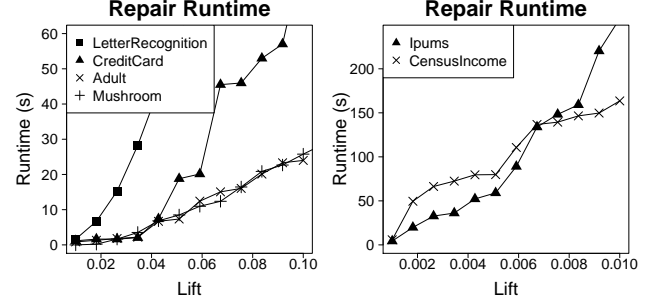


Figure 8. Runtime of the Repair algorithm in function of maximum lift threshold  $\tau$ .

neighbors for all dirty objects, the runtime plots are similar in shape to the plots in Fig. 4. The required time to repair a single dirty object depends on the number of clean objects, which is typically close to  $|\mathcal{D}|$ . Note that the repair algorithm itself is independent of  $\tau$ , which only affects FBIMINER and A-FBIMINER.

For illustrative purposes, Fig. 7 shows example repairs obtained on the Adult dataset ( $\tau = 0.01$ ).

For our implementation, we make use of the lin-similarity measure [39] which weights both matches and mismatches based on the frequency of the actual values:

$$\text{linsim}(o, o') =$$

$$\frac{\sum_{A \in \mathcal{A}} S(o[A], o'[A])}{\sum_{A \in \mathcal{A}} \log(\text{freq}(\{(A, o[A])\}, \mathcal{D})) + \log(\text{freq}(\{(A, o'[A])\}, \mathcal{D}))}$$

where  $S(o[A], o'[A])$  is given by

$$\begin{cases} 2 \log(\text{freq}(\{(A, o[A])\}, \mathcal{D})) & \text{if } o[A] = o'[A]; \text{ and} \\ 2 \log(\text{freq}(\{(A, o[A])\}, \mathcal{D})) + \\ \log(\text{freq}(\{(A, o'[A])\}, \mathcal{D})) & \text{otherwise.} \end{cases}$$

For example in the context of census data, a match or mismatch in gender would be more influential than a match or mismatch in the age category. Of course, any other similarity measure could be used instead. As part of future work, we intend to compare the influence of different similarity functions.

## IX. CONCLUSION

We have argued that the classical point of view on data quality is too static, and proposed a general dynamic notion of cleanliness instead. We believe that this notion is quite interesting on its own and hope that it will be adopted and explored in various data quality settings.

In this paper, we have specialized the general setting by introducing so-called forbidden itemsets, established properties of the lift measure, and provided an algorithm to mine low-lift forbidden itemsets. Our experiments show that the algorithm

is efficient, and illustrate that forbidden itemsets capture inconsistencies with high precision, while providing a concise representation of dirtiness in data.

Furthermore, we have developed an efficient repair algorithm, guaranteeing that after repairs, no new inconsistencies can be found. By first mining almost forbidden itemsets, we assure that no itemsets become forbidden during a repair. This is an essential ingredient in our dynamic notion of data quality. Experiments show high quality repairs. Crucial here are our pruning strategies for mining almost forbidden itemsets.

As part of future work, we intend to experiment with different likeliness functions for the forbidden itemsets. Intuitively, the likeliness function for any type of single-object constraints. As long as the effect of a fixed number of modifications on this likeliness function can be bounded, then our approach remains applicable for larger classes of constraints. The repair algorithm also warrants further research: can a better repairability be achieved, especially on higher dimensional data? Finally, we seek to acquire datasets with ground truth such that an in-depth comparison of error precision and repair accuracy can be performed for different likeliness functions and repair strategies.

In conclusion, the dynamic view on data quality opens the way for revisiting data quality for other kinds of constraints or patterns. It would be interesting to see how to design repair algorithms in the setting for say, standard constraints such as conditional functional dependencies, among others. In addition, the impact of user interaction on the repairing process and the quality of repairs needs to be addressed.

## REFERENCES

- [1] W. Fan and F. Geerts, *Foundations of Data Quality Management*, ser. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2012.
- [2] I. F. Ilyas and X. Chu, “Trends in cleaning relational data: Consistency and deduplication,” *Foundations and Trends in Databases*, vol. 5, no. 4, pp. 281–393, 2015.
- [3] I. P. Fellegi and D. Holt, “A systematic approach to automatic edit and imputation,” *Journal of the American Statistical Association*, vol. 71, no. 353, pp. 17–35, 1976.
- [4] T. N. Herzog, F. J. Scheuren, and W. E. Winkler, *Data Quality and Record Linkage Techniques*. Springer, 2007.
- [5] X. Chu, I. F. Ilyas, and P. Papotti, “Holistic data cleaning: Putting violations into context,” 2013, pp. 458–469.
- [6] —, “Discovering denial constraints,” *PVLDB*, vol. 6, no. 13, pp. 1498–1509, 2013.
- [7] P. Buneman, J. Cheney, W.-C. Tan, and S. Vansummeren, “Curated databases,” in *PODS*, 2008, pp. 1–12.
- [8] W. Fan, J. Li, S. Ma, N. Tang, and W. Yu, “Towards certain fixes with editing rules and master data,” *VLDB*, vol. 3, no. 1-2, pp. 173–184, 2010.
- [9] F. Geerts, G. Mecca, P. Papotti, and D. Santoro, “The llunatic data-cleaning framework,” *PVLDB*, vol. 6, no. 9, pp. 625–636, 2013.
- [10] F. Chiang and R. J. Miller, “Discovering data quality rules,” *PVLDB*, vol. 1, no. 1, pp. 1166–1177, 2008.
- [11] X. Chu, I. F. Ilyas, P. Papotti, and Y. Ye, “Ruleminer: Data quality rules discovery,” in *ICDE*, 2014, pp. 1222–1225.
- [12] Y. Huhtala, J. Kärkkäinen, P. Porkka, and H. Toivonen, “TANE: an efficient algorithm for discovering functional and approximate dependencies,” *Comput. J.*, vol. 42, no. 2, pp. 100–111, 1999.
- [13] F. Chiang and R. J. Miller, “A unified model for data and constraint repair,” in *ICDE*. IEEE, 2011, pp. 446–457.
- [14] G. Beskales, I. F. Ilyas, L. Golab, and A. Galiullin, “On the relative trust between inconsistent data and inaccurate constraints,” in *ICDE*. IEEE, 2013, pp. 541–552.
- [15] M. Mazuran, E. Quintarelli, L. Tanca, and S. Ugolini, “Semi-automatic support for evolving functional dependencies,” 2016.
- [16] J. Wang, T. Kraska, M. J. Franklin, and J. Feng, “Crowder: Crowdsourcing entity resolution,” *PVLDB*, vol. 5, no. 11, pp. 1483–1494, 2012.
- [17] M. Volkovs, F. Chiang, J. Szlichta, and R. J. Miller, “Continuous data cleaning,” in *ICDE*. IEEE, 2014, pp. 244–255.
- [18] J. He, E. Veltri, D. Santoro, G. Li, G. Mecca, P. Papotti, and N. Tang, “Interactive and deterministic data cleaning: A tossed stone raises a thousand ripples,” in *SIGMOD*. ACM, 2016, pp. 893–907.
- [19] T. Dasu and J. M. Loh, “Statistical distortion: Consequences of data cleaning,” *PVLDB*, vol. 5, no. 11, pp. 1674–1683, 2012.
- [20] C. C. Aggarwal, *Outlier analysis*. Springer, 2013.
- [21] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly detection: A survey,” *ACM computing surveys*, vol. 41, no. 3, 2009.
- [22] M. Markou and S. Singh, “Novelty detection: a review,” *Signal processing*, vol. 83, no. 12, pp. 2481–2497, 2003.
- [23] Z. Abedjan, X. Chu, D. Deng, R. C. Fernandez, I. F. Ilyas, M. Ouzzani, P. Papotti, M. Stonebraker, and N. Tang, “Detecting data errors: Where are we and what needs to be done?” *PVLDB*, vol. 9, no. 12, pp. 993–1004, 2016.
- [24] G. Webb and J. Vreeken, “Efficient discovery of the most interesting associations,” *ACM Transactions on Knowledge Discovery from Data*, vol. 8, no. 3, pp. 1–31, 2014.
- [25] J. Pei, A. K. Tung, and J. Han, “Fault-tolerant frequent pattern mining: Problems and challenges,” *DMKD*, vol. 1, p. 42, 2001.
- [26] R. Gupta, G. Fang, B. Field, M. Steinbach, and V. Kumar, “Quantitative evaluation of approximate frequent pattern mining algorithms,” in *SIGKDD*. ACM, 2008, pp. 301–309.
- [27] H. Mannila and H. Toivonen, “Levelwise search and borders of theories in knowledge discovery,” *DMKD*, vol. 1, no. 3, pp. 241–258, 1997.
- [28] R. Agrawal, T. Imieliński, and A. Swami, “Mining association rules between sets of items in large databases,” in *SIGMOD*, 1993, pp. 207–216.
- [29] M. J. Zaki, S. Parthasarathy, M. Ogihara, W. Li *et al.*, “New algorithms for fast discovery of association rules,” *KDD*, pp. 283–286, 1997.
- [30] W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis, “Conditional functional dependencies for capturing data inconsistencies,” *ACM Transactions on Database Systems*, vol. 33, no. 2, 2008.
- [31] M. J. Zaki and W. Meira Jr, *Data mining and analysis: fundamental concepts and algorithms*. Cambridge University Press, 2014.
- [32] R. Bertens, J. Vreeken, and A. Siebes, “Beauty and brains: Detecting anomalous pattern co-occurrences,” *arXiv preprint arXiv:1512.07048*, 2015.
- [33] N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal, “Discovering frequent closed itemsets for association rules,” in *ICDT*, 1999, pp. 398–416.
- [34] T. Calders and B. Goethals, “Depth-first non-derivable itemset mining,” in *SDM*, 2005, pp. 250–261.
- [35] L. Szathmary, P. Valtchev, A. Napoli, and R. Godin, “Efficient vertical mining of frequent closures and generators,” in *Advances in Intelligent Data Analysis VIII*. Springer, 2009, pp. 393–404.
- [36] M. J. Zaki and C.-J. Hsiao, “Charm: An efficient algorithm for closed itemset mining,” in *SDM*, 2002, pp. 457–473.
- [37] P. C. Arocena, B. Glavic, G. Mecca, R. J. Miller, P. Papotti, and D. Santoro, “Messing up with bart: error generation for evaluating data-cleaning algorithms,” *PVLDB*, vol. 9, no. 2, pp. 36–47, 2015.
- [38] A. Arasu, R. Kaushik, and J. Li, “Datasynt: Generating synthetic data using declarative constraints,” *PVLDB*, vol. 4, no. 12, 2011.
- [39] S. Boriah, V. Chandola, and V. Kumar, “Similarity measures for categorical data: A comparative evaluation,” in *SDM*, 2008, pp. 243–254.

APPENDIX

PROOFS OF SECT. V-B

PROPOSITION 1: For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , if  $J$  is a  $\tau$ -forbidden itemset then  $\text{supp}(I, \mathcal{D}) \geq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\sigma_I^{max} \times \tau}$ .

We show this by contradiction. Let  $J$  be a  $\tau$ -forbidden itemset for  $\tau < 1$  and assume for the sake of contradiction that there exists a  $I \subsetneq J$  with  $\text{supp}(I, \mathcal{D}) < \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\sigma_I^{max} \times \tau}$ . Then observe the following:

$$\begin{aligned} \tau &\geq \text{lift}(J, \mathcal{D}) \geq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\text{supp}(I, \mathcal{D}) \times \text{supp}(J \setminus I, \mathcal{D})} \\ &> \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\sigma_I^{max} \times \tau} \times \text{supp}(J \setminus I, \mathcal{D})} \\ &> \frac{\sigma_I^{max} \times \tau \times |\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}) \times \text{supp}(J \setminus I, \mathcal{D})} \\ &> \frac{\tau \times \sigma_I^{max}}{\text{supp}(J \setminus I, \mathcal{D})} \\ &> \tau, \end{aligned}$$

where the last inequality follows from the fact that  $\sigma_I^{max} \geq \text{supp}(J \setminus I, \mathcal{D})$ . We thus have that  $\tau > \tau$  which is clearly impossible. Hence, every subset  $I$  of a  $\tau$ -forbidden itemset  $J$  must satisfy  $\text{supp}(I, \mathcal{D}) < \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\sigma_I^{max} \times \tau}$ .

PROPOSITION 2: For any three itemsets  $I, J$  and  $K$  such that  $I \subset J \subseteq K$  holds, if  $K$  is a  $\tau$ -forbidden itemset then  $\text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D}) \geq \frac{1}{\tau} - \frac{\sigma_I^{max}}{|\mathcal{D}|} > 0$

We show this by contradiction. Let  $K$  be a  $\tau$ -forbidden itemset for  $\tau < 1$  and assume for the sake of contradiction that there exist subsets  $I \subset J \subseteq K$  with  $\text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D}) < \frac{1}{\tau} - \frac{\sigma_I^{max}}{|\mathcal{D}|}$ . Let  $k = \text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D})$  and rewrite  $\text{supp}(K, \mathcal{D})$  as

$$\text{supp}(K, \mathcal{D}) = \text{supp}(K, \mathcal{D}) + \text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D}) - k.$$

Let  $L$  denote the subset  $K \setminus (J \setminus I)$ , it then holds that  $\text{supp}(L, \mathcal{D}) = \text{supp}(K, \mathcal{D}) + \text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D})$ . Hence,  $\text{supp}(K, \mathcal{D}) = \text{supp}(L, \mathcal{D}) - k$ . Furthermore, since  $\text{lift}(K, \mathcal{D}) \leq \tau$ , Prop. 1 requires that  $\text{supp}(L, \mathcal{D}) \geq 1/\tau$ . We show that this leads to a contradiction. Indeed,

$$\begin{aligned} \tau &\geq \text{lift}(K, \mathcal{D}) \geq \frac{|\mathcal{D}| \times \text{supp}(K, \mathcal{D})}{\text{supp}(K \setminus L, \mathcal{D}) \times \text{supp}(L, \mathcal{D})} \\ \tau &\geq \frac{|\mathcal{D}|}{\text{supp}(K \setminus L, \mathcal{D})} \times \frac{\text{supp}(L, \mathcal{D}) - k}{\text{supp}(L, \mathcal{D})} \\ \tau &\geq \frac{|\mathcal{D}|}{\sigma_I^{max}} \times \left(1 - \frac{k}{1/\tau}\right). \end{aligned}$$

From this, we can infer that  $\tau \times \sigma_I^{max} \geq |\mathcal{D}| - (\tau \times k \times |\mathcal{D}|)$ , or that  $\tau \times k \times |\mathcal{D}| \geq \frac{|\mathcal{D}|}{\tau} - \tau \times \sigma_I^{max}$ . Rearranging terms leads to  $k \geq \frac{|\mathcal{D}|}{\tau \times |\mathcal{D}|} - \frac{\tau \times \sigma_I^{max}}{\tau \times |\mathcal{D}|}$  and hence

$$k = \text{supp}(I, \mathcal{D}) - \text{supp}(J, \mathcal{D}) \geq \frac{1}{\tau} - \frac{\sigma_I^{max}}{|\mathcal{D}|},$$

contradicting our assumption.

PROPOSITION 3: For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , it holds that:

$$\text{lift}(J, \mathcal{D}) \geq \frac{\text{supp}(J, \mathcal{D}) \times |\mathcal{D}|}{\min_{S \subset I} \{ \text{supp}(S, \mathcal{D}) \times \text{supp}(I \setminus S, \mathcal{D}) \}}.$$

We show the following: For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , it holds that:

$$\begin{aligned} &\min_{S \subset I} \{ \text{supp}(S, \mathcal{D}) \times \text{supp}(I \setminus S, \mathcal{D}) \} \\ &\geq \min_{S \subset J} \{ \text{supp}(S, \mathcal{D}) \times \text{supp}(J \setminus S, \mathcal{D}) \}. \end{aligned}$$

Clearly, this implies the proposition. We next show that inequality holds. In fact, this follows from the fact any expression  $\text{supp}(S, \mathcal{D}) \times \text{supp}(I \setminus S, \mathcal{D})$  can be rewritten as  $\text{supp}(S, \mathcal{D}) \times \text{supp}((J \setminus S) \setminus X, \mathcal{D})$  where  $X = J \setminus I$ . Since  $\text{supp}((J \setminus S) \setminus X, \mathcal{D}) \geq \text{supp}(J \setminus S, \mathcal{D})$ , we obtain that the following holds for any  $S \subset I$ :  $\text{supp}(S, \mathcal{D}) \times \text{supp}(I \setminus S, \mathcal{D}) \geq \text{supp}(S, \mathcal{D}) \times \text{supp}(J \setminus S, \mathcal{D})$ .

PROPOSITION 4: If  $I$  is a  $\tau$ -forbidden itemset then its frequency is bounded by  $\text{freq}(I, \mathcal{D}) \leq \frac{2}{\tau} - 2\sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}} - 1$ . Furthermore, for small values of  $\tau$  this upper bound converges (from above) to  $\frac{\tau}{4}$ .

We first compute an upper bound on  $\text{supp}(I, \mathcal{D})$  for an itemset  $I$  which is  $\tau$ -forbidden for a given threshold  $\tau < 1$ . Denote by  $S$  a subset of  $I$  that minimizes  $\min_{\emptyset \subset S \subset I} \{ \text{supp}(S, \mathcal{D}) \times \text{supp}(I \setminus S, \mathcal{D}) \}$ . Observe that  $\text{supp}(S, \mathcal{D}) + \text{supp}(I \setminus S, \mathcal{D}) \leq |\mathcal{D}| + \text{supp}(I, \mathcal{D})$ . It can now be easily verified that the maximal value of the expression  $\text{supp}(S, \mathcal{D}) \times \text{supp}(I \setminus S, \mathcal{D})$  is obtained when  $\text{supp}(S, \mathcal{D}) = \frac{|\mathcal{D}| + \text{supp}(I, \mathcal{D})}{2}$ .

Since  $\text{lift}(I, \mathcal{D}) \leq \tau$  and the denominator in the definition of lift (Def. 4) is maximized for  $\text{supp}(S, \mathcal{D}) = \frac{|\mathcal{D}| + \text{supp}(I, \mathcal{D})}{2}$ , we have the following upper bound:

$$\begin{aligned} \text{supp}(I, \mathcal{D}) &\leq \frac{\tau}{4|\mathcal{D}|} \times (|\mathcal{D}|^2 + \\ &\quad 2|\mathcal{D}| \times \text{supp}(I, \mathcal{D}) + (\text{supp}(I, \mathcal{D}))^2) \end{aligned}$$

It is now a routine exercise to find the maximal value of  $\text{supp}(I, \mathcal{D})$ . Indeed, one simple needs to solve the equation

$$\begin{aligned} \frac{\tau}{4|\mathcal{D}|} \times (|\mathcal{D}|^2 + 2|\mathcal{D}| \times \text{supp}(I, \mathcal{D}) + \\ (\text{supp}(I, \mathcal{D}))^2) - \text{supp}(I, \mathcal{D}) = 0 \end{aligned}$$

We obtain the following two solutions:

$$\text{supp}(I, \mathcal{D})^\pm = \frac{2|\mathcal{D}|}{\tau} - |\mathcal{D}| \pm 2|\mathcal{D}| \sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}}$$

Since  $\frac{2|\mathcal{D}|}{\tau} - |\mathcal{D}| > |\mathcal{D}|$  for any  $\tau < 1$ , we see that  $\text{supp}(I, \mathcal{D})^+ = \frac{2|\mathcal{D}|}{\tau} - |\mathcal{D}| + 2|\mathcal{D}| \sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}} > |\mathcal{D}|$ . This is impossible since no itemset can have a support strictly greater than  $|\mathcal{D}|$ . Hence, we are left with the other root  $\text{supp}(I, \mathcal{D})^- = \frac{2|\mathcal{D}|}{\tau} - |\mathcal{D}| - 2|\mathcal{D}| \sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}}$ . To obtain the

upper bound on the frequency we divide by  $|\mathcal{D}|$ . This results in that the maximal frequency of a  $\tau$ -forbidden itemset is  $\frac{2}{\tau} - 2\sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}} - 1$ .

To show the second statement in the proposition, we consider the Taylor expansion of  $\frac{2}{\tau} - 2\sqrt{\frac{1}{\tau^2} - \frac{1}{\tau}} - 1$  given by

$$\frac{\tau}{4} + \frac{\tau^2}{8} + \frac{5\tau^3}{64} + \frac{7\tau^4}{128} + O(\tau^5).$$

Hence, for decreasing values of  $\tau$ , the maximal frequency of a  $\tau$ -forbidden itemset converges from above to  $\tau/4$ .  $\square$

### Proofs of Sect. VI-A

**PROPOSITION 5:** Let  $I$  and  $J$  be itemsets. Suppose that  $\text{lift}(J, \mathcal{D}) = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\text{supp}(I, \mathcal{D}) \times \text{supp}(J \setminus I, \mathcal{D})}$  for some  $I \subset J$ , with  $\text{supp}(I, \mathcal{D}) \leq \text{supp}(J \setminus I, \mathcal{D})$ . If  $J \in \text{FBI}(\mathcal{D}', \tau)$  for some  $\mathcal{D}'$  obtained from  $\mathcal{D}$  by at most  $k$  modifications, then  $\text{mpl}(J, I, k) \leq \tau$ .

As a first step, we compute the maximal drop in lift after performing a single modification. Consider an itemset  $J$  and consider  $\emptyset \subset I \subset J$  such that

$$\text{lift}(J, \mathcal{D}) = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\text{supp}(I, \mathcal{D}) \times \text{supp}(J \setminus I, \mathcal{D})}.$$

Without loss of generality, we assume that  $\text{supp}(I, \mathcal{D}) \leq \text{supp}(J \setminus I, \mathcal{D})$ . In general, when the supports of  $I, J$  and  $J \setminus I$  are all greater than 1, we distinguish between 10 possible changes for  $\text{lift}(J, \mathcal{D})$  after performing a single modification. Denote by  $\mathcal{D}'$  the dataset obtained from  $\mathcal{D}$  by this modification. Then, either

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times (\text{supp}(J, \mathcal{D}) - 1)}{(\text{supp}(I, \mathcal{D}) - 1) \times \text{supp}(J \setminus I, \mathcal{D})} \quad (1)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times (\text{supp}(J, \mathcal{D}) - 1)}{\text{supp}(I, \mathcal{D}) \times (\text{supp}(J \setminus I, \mathcal{D}) - 1)} \quad (2)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times (\text{supp}(J, \mathcal{D}) - 1)}{(\text{supp}(I, \mathcal{D}) - 1) \times (\text{supp}(J \setminus I, \mathcal{D}) - 1)} \quad (3)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times (\text{supp}(J, \mathcal{D}) + 1)}{(\text{supp}(I, \mathcal{D}) + 1) \times \text{supp}(J \setminus I, \mathcal{D})} \quad (4)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times (\text{supp}(J, \mathcal{D}) + 1)}{\text{supp}(I, \mathcal{D}) \times (\text{supp}(J \setminus I, \mathcal{D}) + 1)} \quad (5)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times (\text{supp}(J, \mathcal{D}) + 1)}{(\text{supp}(I, \mathcal{D}) + 1) \times (\text{supp}(J \setminus I, \mathcal{D}) + 1)} \quad (6)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{(\text{supp}(I, \mathcal{D}) + 1) \times \text{supp}(J \setminus I, \mathcal{D})} \quad (7)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\text{supp}(I, \mathcal{D}) \times (\text{supp}(J \setminus I, \mathcal{D}) + 1)} \quad (8)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{(\text{supp}(I, \mathcal{D}) - 1) \times \text{supp}(J \setminus I, \mathcal{D})} \quad (9)$$

$$\text{lift}(J, \mathcal{D}') = \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D})}{\text{supp}(I, \mathcal{D}) \times (\text{supp}(J \setminus I, \mathcal{D}) - 1)} \quad (10)$$

Since  $\text{supp}(I, \mathcal{D}) \leq \text{supp}(J \setminus I, \mathcal{D})$ , it is clear that cases (1) and (3) result in a higher lift than (2). Similarly, case (8) gives

a higher lift than case (7) while case (9) gives a higher lift than case (10). Case (10) is in turn dominated by case (2). Cases (4) - (6) create a new co-occurrence which effectively increases the lift value if the original lift was  $< 1$ . This leaves cases (2) and (7) as valid options for minimizing the lift after an edit, which case is optimal depends on the specific support values of  $I$  and  $J \setminus I$ .

If we denote by  $\sigma_I, \sigma_J$  and  $\sigma_{J \setminus I}$  the supports of  $I, J$  and  $J \setminus I$  in  $\mathcal{D}$ , then it is readily verified that  $\text{MINPOSSLIFT}(\sigma_I, \sigma_{J \setminus I}, \sigma_J, 1) \leq \text{lift}(J, \mathcal{D}')$ .

Since  $\text{MINPOSSLIFT}(\sigma_I, \sigma_{J \setminus I}, \sigma_J, 1)$  is independent of the exact modification used to go from  $\mathcal{D}$  to  $\mathcal{D}'$ ,  $\text{MINPOSSLIFT}(\sigma_I, \sigma_{J \setminus I}, \sigma_J, 1) \leq \text{lift}(J, \mathcal{D}')$  holds for any such  $\mathcal{D}'$ .

A lower bound on the minimum lift of an itemset  $J$  after  $k$  modifications are made is computed iteratively, as there is no closed formula solution available. We compute the maximal drop in lift after performing a single edit, and repeat this procedure  $k$  times. Using a similar argument as above, it is guaranteed that  $\text{MINPOSSLIFT}(\sigma_I, \sigma_{J \setminus I}, \sigma_J, k) \leq \text{lift}(J, \mathcal{D}')$  holds for any  $\mathcal{D}'$  obtained from  $\mathcal{D}$  by at most  $k$  modifications.

This suffices to show Proposition 5. Indeed, consider  $J \in \text{FBI}(\mathcal{D}', \tau)$  where  $\mathcal{D}'$  is obtained from  $\mathcal{D}$  by at most  $k$  modifications. In other words,  $\text{lift}(J, \mathcal{D}') \leq \tau$  and since  $\text{MINPOSSLIFT}(\sigma_I, \sigma_{J \setminus I}, \sigma_J, k) \leq \text{lift}(J, \mathcal{D}')$ , also

$$\text{MINPOSSLIFT}(\sigma_I, \sigma_{J \setminus I}, \sigma_J, k) \leq \tau$$

holds, as desired.  $\square$

### PROOFS OF SECT. VI-B

**PROPOSITION 6:** For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , if  $J$  is a  $\tau$ -forbidden itemset in  $\mathcal{D}'$  then we have that  $\text{supp}(I, \mathcal{D}_{\text{clean}}) \geq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{\text{max}} \times \tau} - k$ .  $\square$

We show this by contradiction. Let  $J$  be a  $\tau$ -forbidden itemset in  $\mathcal{D}'$  for  $\tau < 1$  and assume for the sake of contradiction that there exists a  $I \subsetneq J$  with  $\text{supp}(I, \mathcal{D}_{\text{clean}}) < \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{\text{max}} \times \tau} - k$ . Recall that for any itemset  $I$  and dataset  $\mathcal{D}'$  obtained from  $\mathcal{D}$  by modifying at most  $k$  dirty objects, it holds that:

$$\text{supp}(I, \mathcal{D}_{\text{clean}}) \leq \text{supp}(I, \mathcal{D}') \leq \text{supp}(I, \mathcal{D}_{\text{clean}}) + k$$

It follows that:

$$\text{supp}(I, \mathcal{D}') \leq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{\text{max}} \times \tau} - k + k$$

And hence  $\text{supp}(I, \mathcal{D}') \leq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{\text{max}} \times \tau}$ . According to Prop. 1, this implies that  $J$  cannot be a  $\tau$ -forbidden itemset in  $\mathcal{D}'$ , which contradicts our initial assumption. Hence, every subset  $I$  of a  $\tau$ -forbidden itemset  $J$  in  $\mathcal{D}'$  must satisfy  $\text{supp}(I, \mathcal{D}_{\text{clean}}) \geq \frac{|\mathcal{D}| \times \text{supp}(J, \mathcal{D}')}{\sigma_{I, \mathcal{D}'}^{\text{max}} \times \tau} - k$ .  $\square$

**PROPOSITION 7:** For any three itemsets  $I, J$  and  $K$  such that  $I \subset J \subset K$ , if  $K$  is a  $\tau$ -forbidden itemset in  $\mathcal{D}'$  then  $\text{supp}(I, \mathcal{D}_{\text{clean}}) - \text{supp}(J, \mathcal{D}_{\text{clean}}) \geq \frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}^{\text{max}}}{|\mathcal{D}|} - k$ .  $\square$

## ADDITIONAL PLOTS: ALMOST FORBIDDEN ITEMSETS

We show this by contradiction. Let  $K$  be a  $\tau$ -forbidden itemset in  $\mathcal{D}'$  for  $\tau < 1$  and assume for the sake of contradiction that there exist subsets  $I \subset J \subseteq K$  with  $\text{supp}(I, \mathcal{D}_{clean}) - \text{supp}(J, \mathcal{D}_{clean}) < \frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}}{|\mathcal{D}'|} - k$ . Recall that for any itemset  $I$  and dataset  $\mathcal{D}'$  obtained from  $\mathcal{D}$  by modifying at most  $k$  dirty objects, it holds that:

$$\text{supp}(I, \mathcal{D}_{clean}) \leq \text{supp}(I, \mathcal{D}') \leq \text{supp}(I, \mathcal{D}_{clean}) + k$$

Using these inequalities for  $I$ , we obtain that:

$$\text{supp}(I, \mathcal{D}') - \text{supp}(J, \mathcal{D}_{clean}) < \frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}}{|\mathcal{D}'|} - k + k$$

Since  $\text{supp}(J, \mathcal{D}_{clean}) \leq \text{supp}(J, \mathcal{D}')$ , it follows that:

$$\text{supp}(I, \mathcal{D}') - \text{supp}(J, \mathcal{D}') < \frac{1}{\tau} - \frac{\sigma_{I, \mathcal{D}'}}{|\mathcal{D}'|}$$

By Prop. 2, this means that the itemset  $J$  cannot be  $\tau$ -forbidden in  $\mathcal{D}'$ , contradicting our initial assumption.  $\square$

PROPOSITION 8: For any two itemsets  $I$  and  $J$  such that  $I \subset J$ , it holds that  $\text{lift}(J, \mathcal{D}') \geq$

$$\frac{\text{supp}(J, \mathcal{D}') \times |\mathcal{D}'|}{\min_{S \subset I} \{(\text{supp}(S, \mathcal{D}_{clean}) + k) \times (\text{supp}(I \setminus S, \mathcal{D}_{clean}) + k)\}}.$$

$\square$

This is a straightforward extension of Prop. 3:

$$\text{lift}(J, \mathcal{D}') \geq \frac{\text{supp}(J, \mathcal{D}') \times |\mathcal{D}'|}{\min_{S \subset I} \{\text{supp}(S, \mathcal{D}') \times \text{supp}(I \setminus S, \mathcal{D}')\}}.$$

By applying the inequalities  $\text{supp}(S, \mathcal{D}') \leq \text{supp}(S, \mathcal{D}_{clean}) + k$  and  $\text{supp}(I \setminus S, \mathcal{D}') \leq \text{supp}(I \setminus S, \mathcal{D}_{clean}) + k$ , it follows that:

$$\frac{\text{supp}(J, \mathcal{D}') \times |\mathcal{D}'|}{\min_{S \subset I} \{(\text{supp}(S, \mathcal{D}_{clean}) + k) \times (\text{supp}(I \setminus S, \mathcal{D}_{clean}) + k)\}}.$$

$\square$

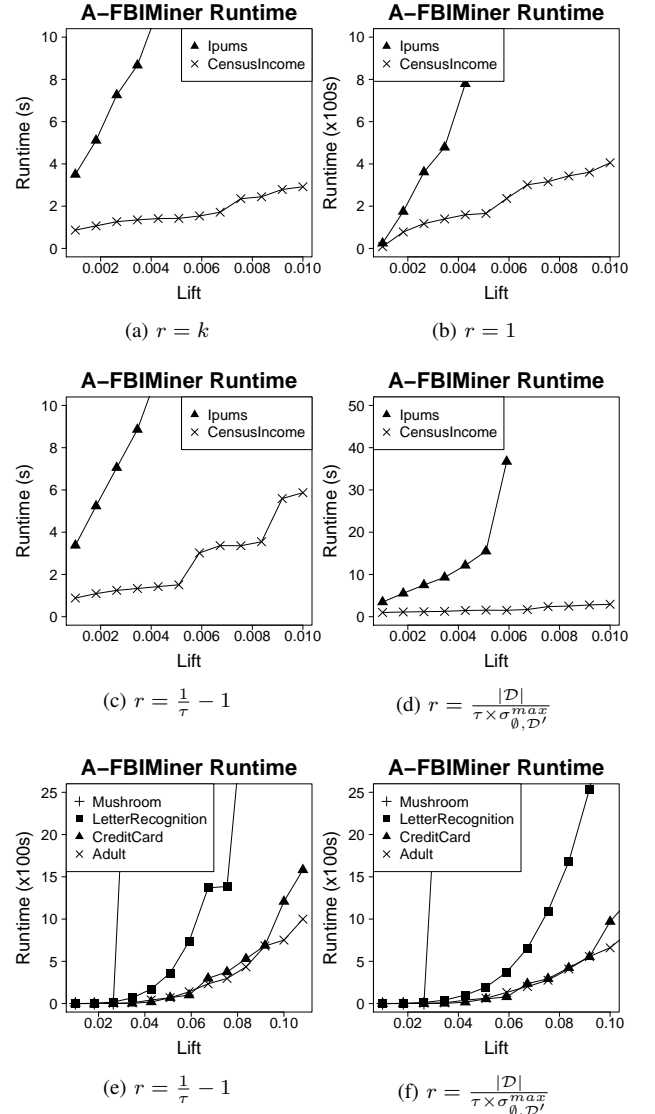


Figure 9. Runtime of A-FBIMINER in function of maximum lift threshold  $\tau$ , for various block sizes  $r$ .