

# Mining Approximate Acyclic Schemes from Relations

## ABSTRACT

Acylic schemes have numerous applications in databases and in machine learning, such as improved design, more efficient storage, and increased performance for queries and machine learning algorithms. Multivalued dependencies (MVDs) are the building blocks of acyclic schemes. The discovery from data of both MVDs and acyclic schemes is more challenging than other forms of data dependencies, such as Functional Dependencies, because these dependencies do not hold on subsets of data, and because they are very sensitive to noise in the data; for example a single wrong or missing tuple may invalidate the schema. In this paper we present *Maimon*, a system for discovering *approximate* acyclic schemes and MVDs from data. We give a principled definition of approximation, by using notions from information theory, then describe the two components of *Maimon*: mining for approximate MVDs, then reconstructing acyclic schemes from approximate MVDs. We conduct an experimental evaluation of *Maimon* on 20 real-world datasets, and show that it can scale up to 1M rows, and up to 30 columns.

## 1 INTRODUCTION

Acylic schemes have numerous applications in databases and in machine learning. Originally introduced by Beeri [5], they have lead to Yannakakis celebrated linear time query evaluation algorithms [40], and are used widely today in database design [13, 31], to speed up query evaluation with multiple aggregates [25], and to speed up machine learning applications such as ridge linear regression, classification trees, and regression trees [24, 37, 38]. When considering which types of schemes to fit the data, acyclic schemes are the natural choice due to their many desirable properties [6]. In this paper we study the following discovery problem: given a database consisting of a single relation, generate a set of acyclic schemes that fit the data to a large extent. For a simple illustration, consider the database shown on the left of Figure 1. It can be decomposed into an acyclic schema with four relations, shown on the right.

The building blocks of an acyclic schema are Multivalued Dependencies, MVDs. Every acyclic schema can be fully specified by the set of MVDs that it implies, which we call its *support*. Therefore, when mining acyclic schemes, the first step is to mine the MVDs satisfied by the data. MVDs were first introduced by Fagin [13], which used them to introduce the 4th normal form, a generalization of the Boyce-Codd normal form (BCNF) [10]. They were studied extensively in the database literature [2, 4, 14, 28], have been proven

to be equivalent to Saturated Conditional Independence in graphical models [17], and the methods used to synthesize an acyclic schema from a set of MVDs are well known [3, 7, 13, 32]. However, despite their importance, there is little research on the *discovery* of MVDs from data [1].

Work most closely related to the discovery of MVDs has been on discovering Functional Dependencies (FDs) and Unique Column Combinations (UCCs) [8, 21, 26, 27, 33, 35, 39]. These are special cases of MVDs, but MVDs are more general. Discovering all FDs and all UCCs is insufficient for discovering acyclic schemes. The only work that addressed the discovery problem for MVDs is by Savnik and Flach [36] and a master thesis by Draeger [12], and none of them address the more challenging task of discovering acyclic schemes.

There are two major challenges that make the discovery of MVDs and acyclic schemes, much harder than that of FDs and UCCs. First, they don't hold on subsets of the data. If a relation satisfies an FD, or a UCC, then every subset also satisfies the FD, or UCC, and this is exploited by many discovery algorithms, e.g FastFD [39] mines FDs in all subsets of size 2, while HyFD [35] mines FDs in a small subset extracted from the data. This property fails for MVDs, preventing us from considering subsets of the data. Second, MVDs and acyclic schemes are much more sensitive to data errors than FDs and UCCs. Even a single missing tuple may invalidate an MVD or schema. Real-world data often has important dependencies that do not hold exactly, but, if discovered, are very useful for a variety of applications. For that reason, in this paper we study the problem of discovering *approximate* MVDs and consequently, *approximate* acyclic schemes.

We present *Maimon*<sup>1</sup>, the first system for discovering approximate MVDs and acyclic schemes in the data. We introduce a principled notion of approximation, based on information theory, and develop the necessary theory for reasoning about approximate MVDs and schemes. We then describe algorithms for mining MVDs and schemes, and evaluate their scalability on real-world datasets of up to 1M rows, and 30 attributes. By allowing approximations, *Maimon* finds more interesting schemes without incurring too high a loss (i.e., spurious tuples). We make several contributions.

Our first contribution is to introduce a principled definition of approximation, and study its properties. Kivinen and Mannila [26] give three definitions of approximate functional dependencies, and Kruse and Naumann use one of them in their approximate FDs and UCCs discovery algorithm [27].

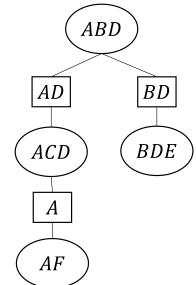
<sup>1</sup> *Maimon* stands for Multivalued Approximate Inference Mining and Normalization.

R:					
A	B	C	D	E	F
$a_1$	$b_1$	$c_1$	$d_1$	$e_1$	$f_1$
$a_2$	$b_2$	$c_1$	$d_1$	$e_2$	$f_2$
$a_2$	$b_2$	$c_2$	$d_2$	$e_3$	$f_2$
$a_1$	$b_2$	$c_1$	$d_2$	$e_3$	$f_1$
$a_1$	$b_2$	$c_1$	$d_2$	$e_2$	$f_1$

A	B	D
$a_1$	$b_1$	$d_1$
$a_2$	$b_2$	$d_1$
$a_2$	$b_2$	$d_2$
$a_1$	$b_2$	$d_2$
$a_1$	$b_2$	$d_2$

A	C	D
$a_1$	$c_1$	$d_1$
$a_2$	$b_2$	$d_1$
$a_2$	$b_2$	$d_2$
$a_1$	$c_1$	$d_2$
$a_1$	$c_1$	$d_2$

B	D	E
$b_1$	$d_1$	$e_1$
$b_2$	$d_1$	$e_2$
$b_2$	$d_2$	$e_3$
$b_2$	$d_2$	$e_2$
$b_2$	$d_2$	$e_2$



**Fig. 1** A relation R and its decomposition into an acyclic schema

We propose an alternative metric of approximation, based on information theory. Each MVD or acyclic schema is associated with an information theoretic expression, and its value represents the degree of approximation. Our definition builds on early work by Lee [30].

Second, we propose novel algorithms for mining approximate MVDs and approximate acyclic schemes. For mining MVDs, our theoretical results prove that we do not need to discover *all* approximate MVDs, but only the so-called *full* MVDs with *minimal* separators. Our algorithm builds on previous results by Gunopulos et al. [20] for discovering the most *specific* sentences in the data that meet a certain criterion (e.g., maximal sets of items whose frequency in the data is above a given threshold). Following the discovery of the MVDs that hold in the data, we turn to the task of enumerating the acyclic schemes that can be synthesized from the set of discovered MVDs. Our algorithm is based on an approach for efficiently enumerating the maximal independent sets of a graph [11, 22], which has also been applied to the problem of enumerating tree decompositions [9].

Third, we evaluate Maimon on 20 real-world datasets, and were able to scale it up to 1M rows, and up to 30 columns. We show that the loss entailed by the generated acyclic schemes (i.e., number of spurious tuples), monotonically depends on the information theoretic measure of approximation we develop herein. We also show that a larger degree of approximation enables the discovery of schemes that exhibit a larger degree of decomposition. These schemes generally have more relations, and the *width* of the schema (i.e., relation with the largest number of attributes), is smaller.

The most expensive operation of Maimon is the computation of the entropy  $H(X)$  of a set of attributes  $X$ . Each such computation requires a full scan over the data, and this is prohibitively expensive due to the exponential number of subsets of attributes. We describe a novel, efficient approach to computing entropy, which reduces the problem to a set of main-memory SQL queries. Our method is inspired by the *PLI cache* (Position List Indices) data structure used for mining both exact and approximate FDs [21, 27].

To sum up, the contributions of this work are as follows:

- (1) We define a principled notion of approximate data dependencies based on information theory, and study its properties; Sec. 4 and 5.
- (2) We describe a novel MVD enumeration algorithms and acyclic schema enumeration algorithm; Sec. 6 and 7.
- (3) We conduct an extensive experimental evaluation on 20 real datasets; Sec. 8.

## 2 RUNNING EXAMPLE

We will use the following running example in this paper. Consider the relation R over the signature  $\Omega = \{A, B, C, D, E, F\}$  in Figure 1. Ignore the probabilities, we will use them in Sec. 3. Also, ignore for now the last row (in red). The table with four rows can be decomposed into four tables, shown in the figure. More precisely, the following join dependency holds:  $R = R[ABD] \bowtie R[ACD] \bowtie R[BDE] \bowtie R[AF]$ . The schema of these four tables is *acyclic*, because it admits a join tree, shown in Fig. 2 (reviewed in Sec. 3). Our goal is to discover this acyclic schema from the data R. For that, we note that the acyclic schema can be entirely described by three Multivalued Dependencies:  $BD \rightarrow E|ACF$ ,  $AD \rightarrow CF|BE$ , and  $A \rightarrow F|BCDE$ . Each corresponds to one edge of the join tree: the left hand size of the MVD (that we call the *key*) is the label of that edge, while the two sets of attributes correspond to the subtrees connected by the edge. For example, the edge  $ACD \xrightarrow{AD} ABD$  in the join tree defines the MVD  $AD \rightarrow CF|BE$ . The key  $AD$  “separates” the attributes  $CF$  in one subtree from  $BE$  in the other subtree, and we will also call such a set a *separator*. Since MVDs are the building blocks of acyclic schemas, their discovery is a prerequisite for discovering acyclic schemas, and our first task is to discover MVDs from data, then use them to discover acyclic schemas.

Consider the 5'th row in R, shown in red. By adding it, we need to add a 4'th row to  $R[BDE]$ , also shown in red. However, now the join dependency no longer holds exactly, because  $R[ABD] \bowtie R[ACD] \bowtie R[BDE] \bowtie R[AF]$  contains a spurious tuple, namely  $(a_2, b_2, c_2, d_2, e_2, f_2)$ , which is not in R (it is not shown in the Figure); the first two MVDs no longer hold, only  $A \rightarrow F|BCDE$  still holds, and the acyclic schema is no longer a correct decomposition of R. Yet the schema

**Fig. 2** Join Tree

$\Omega$	set of variables (attributes)
$n =  \Omega $	number of variables (attributes)
$X, Y, A, B, \dots$	sets of variables $\subseteq \Omega$
$S$	a schema = $\{\Omega_1, \dots, \Omega_m\}$
$X \twoheadrightarrow Y Z$	a standard MVD
$X \twoheadrightarrow Y_1 Y_2 \dots Y_m$	an MVD [4]
$(\mathcal{T}, \chi)$	a join tree
$H(X)$	entropy of a set of variables $X$
$H(Y X), I(Y; Z X)$	entropic measures
$\mathcal{J}(\mathcal{T}, \chi)$	an entropic measure
$\mathcal{J}(S)$	$\mathcal{J}$ of any join tree for $S$
$\mathcal{J}(X \twoheadrightarrow Y_1 \dots Y_m)$	$\mathcal{J}$ of the schema $\{XY_1, \dots, XY_m\}$
$\mathcal{J}(X \twoheadrightarrow Y Z)$	$= I(Y; Z X)$
$R$	a relation
$N =  R $	number of tuples
$R \models AJD(S)$	$R$ satisfies an acyclic join dependency
$R \models_{\epsilon} AJD(S)$	$R$ $\epsilon$ -satisfies an acyclic join dependency

**Table 1: Notations**

can still be useful for many applications, even if it leads to a spurious tuple. Insisting on exact acyclic schemas would severely restrict their applications, and also make them very brittle since the addition of one single tuple would invalidate the schema. In this paper we compute *approximate* acyclic schemas, and *approximate* MVDs. By allowing approximations, the schema shown in the figure is still considered valid for the data, despite the spurious tuple.

### 3 BACKGROUND

Table 1 summarizes the notations in this paper. We denote by  $[n] = \{1, \dots, n\}$ . Let  $\Omega$  be a set of variables, also called attributes. If  $X, Y \subseteq \Omega$ , then  $XY$  denotes  $X \cup Y$ .

#### 3.1 Data Dependencies

Fix a relation instance  $R$  of size  $N = |R|$ , and schema  $\Omega$ . For  $Y \subseteq \Omega$  we let  $R[Y]$  denote the projection of  $R$  onto the attributes  $Y$ .

Let  $X, Y, Z \subseteq \Omega$ . A *schema* is a set  $S = \{\Omega_1, \dots, \Omega_m\}$  such that  $\bigcup_{i=1}^k \Omega_i = \Omega$  and  $\Omega_i \not\subseteq \Omega_j$  for  $i \neq j$ . We say that the relation instance  $R$  satisfies the *join dependency*  $JD(S)$ , and write  $R \models JD(S)$ , if  $R = \bowtie_{i=1}^k R[\Omega_i]$ . We say that  $R$  satisfies the *multivalued dependency* (MVD)  $\phi = X \twoheadrightarrow Y_1|Y_2|\dots|Y_m$  where  $m \geq 2$ , the  $Y_i$ s are pairwise disjoint, and  $XY_1 \dots Y_m = \Omega$ , if  $R = R[XY_1] \bowtie \dots \bowtie R[XY_m]$ . We call  $X$  the *key* of the MVD and  $\{Y_1, \dots, Y_m\}$  its *dependents*, denoted  $\text{key}(\phi) = X$  and  $\text{dep}(\phi) = \{Y_1, \dots, Y_m\}$ . Most of the literature considers only MVDs with  $m = 2$ , which we call here *standard MVDs*. Beeri et al. [4] noted that a generalized MVD can encode concisely multiple MVDs; for example  $X \twoheadrightarrow A|B|C$  holds iff

$X \twoheadrightarrow AB|C$ ,  $X \twoheadrightarrow A|BC$  and  $X \twoheadrightarrow AC|B$  hold. We review a *join tree* from [6]:

**Definition 3.1.** A *join tree* is a pair  $(\mathcal{T}, \chi)$  where  $\mathcal{T}$  is an undirected tree, and  $\chi$  is a function that maps each  $u \in \text{nodes}(\mathcal{T})$  to a set of variables  $\chi(u)$ , called a *bag*, such that the following *running intersection* property holds: for every variable  $X$ , the set  $\{u \in \text{nodes}(\mathcal{T}) \mid X \in \chi(u)\}$  is a connected component of  $\mathcal{T}$ . We denote by  $\chi(\mathcal{T}) \stackrel{\text{def}}{=} \bigcup_u \chi(u)$ , the set of variables of the join tree.

We often denote the join tree as  $\mathcal{T}$ , dropping  $\chi$  when it is clear from the context. The *schema* defined by  $\mathcal{T}$  is  $S = \{\Omega_1, \dots, \Omega_m\}$ , where  $\Omega_1, \dots, \Omega_m$  are the bags of  $\mathcal{T}$ . We call a schema  $S$  *acyclic* if there exists a join tree whose schema is  $S$ . Since we required  $\Omega_i \not\subseteq \Omega_j$  for  $i \neq j$ , one can prove that any acyclic schema with  $n$  attributes and  $m$  relations satisfies  $m \leq n$ . We say that a relation  $R$  satisfies the *acyclic join dependency*  $S$ , and denote  $R \models AJD(S)$ , if  $S$  is acyclic and  $R \models JD(S)$ . An MVD  $X \twoheadrightarrow Y_1|\dots|Y_m$  represents a simple acyclic schema, namely  $S = \{XY_1, XY_2, \dots, XY_m\}$ .

Let  $S = \{\Omega_1, \dots, \Omega_m\}$  be an acyclic schema with join tree  $(\mathcal{T}, \chi)$ . We associate to every  $(u, v) \in \text{edges}(\mathcal{T})$  an MVD  $\phi_{u,v}$  as follows. Let  $\mathcal{T}_u$  and  $\mathcal{T}_v$  be the two subtrees obtained by removing the edge  $(u, v)$ . Then, we denote by  $\phi_{u,v} \stackrel{\text{def}}{=} \chi(u) \cap \chi(v) \twoheadrightarrow \chi(\mathcal{T}_u)|\chi(\mathcal{T}_v)$ . We call the *support* of  $\mathcal{T}$  the set of  $m-1$  MVDs associated to its edges, in notation  $MVD(\mathcal{T}) = \{\phi_{u,v} \mid (u, v) \in \text{edges}(\mathcal{T})\}$ . If  $\mathcal{T}$  defines the acyclic schema  $S$ , then it satisfies  $R \models AJD(S)$  iff it satisfies all MVDs in its support:  $R \models \phi_{u,v}$  for all  $\phi_{u,v} \in MVD(\mathcal{T})$  [6, Thm. 8.8].

**Example 3.2.** We will illustrate with the running example from Sec. 2. The tree in Fig. 2 is a join tree. Its bags are the ovals labeled  $AF$ ,  $ACD$ ,  $ABD$ , and  $BDE$ , and it is custom to show the intersection of two bags on the connecting edge.  $MVD(\mathcal{T}) = \{BD \twoheadrightarrow E|ACF, AD \twoheadrightarrow CF|BE, A \twoheadrightarrow F|BCDE\}$ .

#### 3.2 Information Theory

Lee [29, 30] gave an equivalent formulation of data dependencies in terms of information measures; we review this briefly here, after a short background on information theory.

Let  $X$  be a random variable with a finite domain  $\mathcal{D}$  and probability mass  $p$  (thus,  $\sum_{x \in \mathcal{D}} p(x) = 1$ ). Its entropy is:

$$H(X) \stackrel{\text{def}}{=} \sum_{x \in \mathcal{D}} p(x) \log \frac{1}{p(x)} \quad (1)$$

If  $N = |\mathcal{D}|$  then  $H(X) \leq \log N$ , and equality holds iff  $p$  is uniform. For a set of jointly distributed random variables  $\Omega = \{X_1, \dots, X_n\}$  we define the function  $H : 2^\Omega \rightarrow \mathbb{R}$  as the entropy of the joint random variables in the set. For example,  $H(X_1X_2) = \sum_{x_1 \in \mathcal{D}_1, x_2 \in \mathcal{D}_2} p(x_1, x_2) \log \frac{1}{p(x_1, x_2)}$ . Let

$A, B, C \subseteq \Omega$ . The *mutual information*  $I(B; C|A)$  is defined as:

$$I(B; C|A) \stackrel{\text{def}}{=} H(AB) + H(AC) - H(ABC) - H(A) \quad (2)$$

It is known that the conditional independence  $p \models B \perp C | A$  (i.e.,  $B$  is independent of  $C$  given  $A$ ) holds iff  $I(B; C|A) = 0$ .

In this paper we use only the following two properties of the mutual information:

$$I(B; C|A) \geq 0 \quad (3)$$

$$I(B; CD|A) = I(B; C|A) + I(B; D|AC) \quad (4)$$

The first inequality follows from monotonicity and submodularity (it is in fact equivalent to them); the second equality is called the *chain rule*. All consequences of these two (in)equalities are called *Shannon inequalities*; for example, monotonicity  $H(AB) \geq H(A)$  is a Shannon inequality because it follows from (3) by setting  $B = C$ .

Let  $R$  be relation with attributes  $\Omega = \{X_1, \dots, X_n\}$  and  $N$  tuples. The *empirical distribution* is the uniform distribution over its tuples:  $\forall t \in R, p(t) = 1/N$ . Its entropy satisfies  $H(\Omega) = \log N$ . For  $\alpha \subseteq [n]$ , we denote by  $X_\alpha$  the set of variables  $X_i, i \in \alpha$ , and denote by  $R(X_\alpha = x_\alpha)$  the subset of tuples  $t \in R$  where  $t[X_\alpha] = x_\alpha$ , for fixed values  $x_\alpha$ . By uniformity, the marginal probability is  $p(X_\alpha = x_\alpha) = \frac{|R(X_\alpha = x_\alpha)|}{N}$ , and therefore:

$$H(X_\alpha) \stackrel{\text{def}}{=} \log N - \frac{1}{N} \sum_{x_\alpha \in \mathcal{D}_\alpha} |R(X_\alpha = x_\alpha)| \log |R(X_\alpha = x_\alpha)| \quad (5)$$

The sum above can be computed using a simple SQL query:  
Select  $X_\alpha$ , count(\*)  $\times$  log(count(\*)) From  $R$  Group By  $X_\alpha$ .

Lee [29, 30] formalized the following connection between database constraints, and entropic measures. Let  $(\mathcal{T}, \chi)$  be a join tree. We define the following expression:

$$\mathcal{J}(\mathcal{T}, \chi) \stackrel{\text{def}}{=} \sum_{v \in \text{nodes}(\mathcal{T})} H(\chi(v)) - \sum_{(v_1, v_2) \in \text{edges}(\mathcal{T})} H(\chi(v_1) \cap \chi(v_2)) - H(\chi(\mathcal{T})) \quad (6)$$

We abbreviate it with  $\mathcal{J}(\mathcal{T})$ , or  $\mathcal{J}$ , when  $\mathcal{T}, \chi$  are clear from the context; we will prove later (Th. 5.1) that  $\mathcal{J} \geq 0$  is a Shannon inequality. Lee proved that  $\mathcal{J}$  depends only on the schema  $S$  defined by the join tree, and not on the tree itself. To see this on a simple example, consider the MVD  $X \rightarrow U|V|W$  and its associated acyclic schema  $\{XU, XV, XW\}$ . If we consider the join tree  $XU - XV - XW$ , then  $\mathcal{J} = H(XU) + H(XV) + H(XW) - 2H(X) - H(XUVW)$ . Another join tree is  $XU - XW - XV$ , and  $\mathcal{J}$  is the same. Therefore, if  $S$  is acyclic, then we write  $\mathcal{J}(S)$  to denote  $\mathcal{J}(\mathcal{T})$  for any join tree of  $S$ . We denote by  $\mathcal{J}(X \rightarrow Y_1 | \dots | Y_m) \stackrel{\text{def}}{=} H(XY_1) + \dots + H(XY_m) - (m-1)H(X) - H(XY_1 \dots Y_m)$  for any sets of variables  $X, Y_1, \dots, Y_m$  where  $Y_1, \dots, Y_m$  are pairwise disjoint, even when  $XY_1 \dots Y_m$  is not necessarily  $\Omega$ . When  $m = 2$ , then  $J(X \rightarrow Y|Z) = I(Y; Z|X)$ . Lee proved the following:

**THEOREM 3.3.** ([30]) *Let  $H$  be the entropy of the empirical distribution on  $R$ , and let  $S$  be any acyclic schema. Then  $R \models AJD(S)$  iff  $\mathcal{J}(S) = 0$ .*

In the particular case of a standard MVD, Lee's result implies that  $R \models X \rightarrow Y|Z$  if and only if  $I(Y; Z|X) = 0$ .

**Example 3.4.** Continuing Example 3.2, the empirical distribution of the relation  $R$  in Fig 1 (without the red tuple) assigns probability  $1/4$  to each tuple. Thus,  $H(ABCDEF) = \log 4 = 2$ . The marginal probabilities need not be uniform, e.g. the marginals for  $BDE$  are  $1/4, 1/4, 1/2$ , and thus  $H(BDE) = 1/4 \log 4 + 1/4 \log 4 + 1/2 \log 2 = 3/2$ . The value of  $\mathcal{J}$  is:  $\mathcal{J}(\mathcal{T}) = H(AF) + H(ACD) + H(ABD) + H(BDE) - H(A) - H(AD) - H(BD) - H(ABCDEF)$ . For the empirical distribution in the figure, this quantity is 0.

## 4 PROBLEM STATEMENT

Our main goal is to discover an acyclic schema for a given relation instance  $R$ . Since exact schemas are very sensitive to data errors, Maimon discovers approximate schemas.

**Definition 4.1 (Approximate Acyclic Schema).** Fix a relation instance  $R$ , and  $\varepsilon \geq 0$ . We say that an acyclic schema  $S$  is an  $\varepsilon$ -schema for  $R$ , or simply *approximate schema*, if  $\mathcal{J}(S) \leq \varepsilon$ . In notation,  $R \models_\varepsilon AJD(S)$ .

Maimon takes as input  $\varepsilon \geq 0$  and discovers approximate acyclic schemas for  $R$ . By Lee's theorem, if we set  $\varepsilon = 0$ , then Maimon returns exact schemas. In practice, a relation  $R$  may not have any exact schemas, or may have very limited schemas; by allowing  $\varepsilon \geq 0$  we may find approximate schemas that are quite useful for many applications.

**PROBLEM 4.1 (SCHEMA ENUMERATION PROBLEM).** *Given a relational instance  $R$ , enumerate the approximate acyclic schemas of  $R$ .*

In practice, we are not interested in enumerating *all* approximate acyclic schemas of  $R$ . This would take a prohibitively long time, and some acyclic schemas are superior to others. For example, consider a relation over four attributes that satisfies the acyclic join dependency  $S = \{XA, XB, XC\}$ . The following acyclic join dependencies also hold in  $R$ :  $\{XAB, XC\}$ ,  $\{XAC, XB\}$ , and  $\{XA, XBC\}$ . The latter schemas are less useful than  $S = \{XA, XB, XC\}$  that leads to a larger degree of decomposition. Therefore, in this paper we address the problem of enumerating acyclic schemas that cannot be extended (i.e., with additional relational instances) while continuing to satisfy the accuracy threshold.

We derive the approximate schemas from the MVDs in their support. Since an MVD is, in particular, an acyclic schema, Def. 4.1 applies to them as well: a  $\varepsilon$ -MVD is one for which  $\mathcal{J}(X \rightarrow Y_1 | \dots | Y_m) \leq \varepsilon$ . Our second problem is:

**PROBLEM 4.2 (MVD ENUMERATION PROBLEM).** *Given a relational instance R, enumerate the approximate MVDs of R.*

Maimon works as follows. The user provides a parameter  $\varepsilon \geq 0$ . In the first phase, Maimon enumerates  $\varepsilon$ -MVDs, using the algorithm in Sec. 6. When it finishes, or after a timeout, it starts the second phase, where it enumerates approximate schemas with support from the set returned by the first phase, using the algorithm in Sec. 7. Since the support of a schema consists of  $m - 1$  MVDs, the algorithm reports schemas with  $\mathcal{J}(S) \leq (m-1)\varepsilon$ , where  $m$  is the number of relations in S but, since the enumeration algorithm is exhaustive, all schemas with  $\mathcal{J} \leq \varepsilon$  are reported eventually.

## 5 THREE MAIN TECHNIQUES

We describe here three main techniques that allow us to design efficient schema- and MVD-discovery algorithms. The first reduces the approximate schema discovery to approximate MVD discovery, the next two reduce the number of MVD's that need to be discovered.

### 5.1 From MVDs to Acyclic Schemas

Beeri et al. [6] showed that, for exact constraints, an acyclic schema over  $m$  relations is equivalent to the set of  $m-1$  MVDs in its support. We give here a non-trivial generalization to approximate schemas and MVDs. We start with two simple inequalities which we need throughout the paper:

**PROPOSITION 5.1.** *Let  $Y_1, Z_1, \dots, Y_m, Z_m$  be pairwise disjoint sets of variables, and let X be any set of variables. Then the following are Shannon inequalities:*

$$\mathcal{J}(X \rightarrow Y_1 | \dots | Y_m) \leq \mathcal{J}(X \rightarrow Y_1 Z_1 | \dots | Y_m Z_m) \quad (7)$$

$$\mathcal{J}(X Z_1 \dots Z_m \rightarrow Y_1 | \dots | Y_m) \leq \mathcal{J}(X \rightarrow Y_1 Z_1 | \dots | Y_m Z_m) \quad (8)$$

**PROOF.** The first inequality follows from this chain of inequalities:  $\mathcal{J}(X \rightarrow Y_1 | \dots | Y_m) \leq \mathcal{J}(X \rightarrow Y_1 Z_1 | Y_2 | \dots | Y_m) \leq \mathcal{J}(X \rightarrow Y_1 Z_1 | Y_2 Z_2 | \dots | Y_m) \leq \dots$ ; to prove it, we show only the first step (the others are similar), which follows by observing  $\mathcal{J}(X \rightarrow Y_1 | \dots | Y_m) + I(Z_1; Y_2 \dots Y_m | X Y_1) = \mathcal{J}(X \rightarrow Y_1 Z_1 | \dots | Y_m)$  then using inequality (3). The second inequality follows from a similar chain, where the first step follows from  $\mathcal{J}(X Z_1 \rightarrow Y_1 | \dots | Y_m) + \sum_{i=2}^m I(Y_i; Z_1 | X) = \mathcal{J}(X \rightarrow Y_1 Z_1 | Y_2 | \dots | Y_m)$  and the inequality follows from (3).  $\square$

Let  $(\mathcal{T}, \chi)$  be a join tree, defining an acyclic schema S over the variables  $\chi(\mathcal{T}) = \Omega$ . Choose an arbitrary root, orient the tree accordingly, and let  $u_1, \dots, u_m$  be a depth-first enumeration of nodes( $\mathcal{T}$ ). Thus,  $u_1$  is the root, and for every  $i > 1$ ,  $\text{parent}(u_i)$  is some node  $u_j$  with  $j < i$ . For every  $i$ , we define  $\Omega_i \stackrel{\text{def}}{=} \chi(u_i)$ ,  $\Omega_{i:j} \stackrel{\text{def}}{=} \bigcup_{\ell=i,j} \Omega_\ell$ , and  $\Delta_i \stackrel{\text{def}}{=} \chi(\text{parent}(u_i)) \cap \chi(u_i)$  (by the running intersection property this is equal to  $\Omega_{1:(i-1)} \cap \Omega_i$ ). We prove:

**THEOREM 5.1.** *The following hold:*

$$\mathcal{J}(\mathcal{T}) = \sum_{i=2}^m I(\Omega_{1:(i-1)}; \Omega_i | \Delta_i) \quad (9)$$

$$\max_{i=2,m} I(\Omega_{1:(i-1)}; \Omega_{i:m} | \Delta_i) \leq \mathcal{J}(\mathcal{T}) \leq \sum_{i=2}^m I(\Omega_{1:(i-1)}; \Omega_{i:m} | \Delta_i) \quad (10)$$

*The first is an identity, and the second is a Shannon inequality.*

The identity (9) captures precisely the intuition that the information measure associated with a join tree  $\mathcal{T}$  is equivalent to  $m - 1$  mutual information. This identity implies that  $\mathcal{J}(\mathcal{T}) \geq 0$ , because  $I(\dots) \geq 0$ . But the expressions  $I(\dots)$  in (9) do not correspond to MVDs, because they do not include all variables  $\Omega$ . The Shannon inequality (10) rectifies this, by showing that  $\mathcal{J}(\mathcal{T})$  lies between the max and the sum of  $m - 1$  MVDs. Notice that the MVDs  $\Delta_i \rightarrow \Omega_{1:(i-1)} | \Omega_{i:m}$ ,  $i = 2, m$  are precisely the support,  $\text{MVD}(\mathcal{T})$ , thus (10) generalizes Beeri's observation to approximate schemas. An immediate consequence of (10) is the following relationship between an acyclic schema S and its support.

**COROLLARY 5.2.** *Let S be an acyclic schema with join tree  $(\mathcal{T}, \chi)$ . Then: (1) if  $R \models_\varepsilon \text{AJD}(S)$  then  $R \models_\varepsilon \text{MVD}(\mathcal{T})$ . (2) If  $R \models_\varepsilon \text{MVD}(\mathcal{T})$  then  $R \models_{(m-1)\varepsilon} \text{AJD}(S)$ . In particular, (1) and (2) are equivalent if  $\varepsilon = 0$ . Here  $R \models_\varepsilon \text{MVD}(\mathcal{T})$  means  $R \models_\varepsilon \phi$ , for all  $\phi \in \text{MVD}(\mathcal{T})$ .*

**PROOF.** (of Theorem 5.1) Let  $\mathcal{T}_i$  denote the subtree consisting of the nodes  $u_1, \dots, u_i$ . We prove (9) by induction on  $m$ . Assume the identity holds for  $m - 1$ . Compared to  $\mathcal{T}_{m-1}$ , the tree  $\mathcal{T}_m$  has one extra node  $u_m$  and one extra edge  $(\text{parent}(u_m), u_m)$ , hence by the definition of  $\mathcal{J}$  in (6):

$$\begin{aligned} \mathcal{J}(\mathcal{T}_m) &= \mathcal{J}(\mathcal{T}_{m-1}) + H(\chi(u_m)) - H(\chi(u_m) \cap \chi(\text{parent}(u_m))) \\ &\quad + H(\chi(\mathcal{T}_{m-1})) - H(\chi(\mathcal{T}_m)) \\ &= \mathcal{J}(\mathcal{T}_{m-1}) + H(\Omega_m) - H(\Delta_m) + H(\Omega_{1:(m-1)}) - H(\Omega_{1:m}) \\ &= \mathcal{J}(\mathcal{T}_{m-1}) + I(\Omega_{1:(m-1)}; \Omega_m | \Delta_m) \end{aligned}$$

The claim follows from the induction hypothesis on  $\mathcal{J}(\mathcal{T}_{m-1})$ .

We prove (10). The right inequality follows from the fact that  $I(\Omega_{1:(i-1)}; \Omega_i | \Delta_i) \leq I(\Omega_{1:(i-1)}; \Omega_{i:m} | \Delta_i)$  (which holds by Eq. (8)). For the left inequality, we make the following observation. If  $\mathcal{T}$  is any join tree and  $\mathcal{T}'$  is obtained by merging two adjacent nodes  $(u, v) \in \text{edges}(\mathcal{T})$ , then  $\mathcal{J}(\mathcal{T}) \geq \mathcal{J}(\mathcal{T}')$ . This is because  $\mathcal{J}(\mathcal{T}) = \mathcal{J}(\mathcal{T}') + H(\chi(u)) + H(\chi(v)) - H(\chi(u) \cap \chi(v)) - H(\chi(u) \cup \chi(v)) = \mathcal{J}(\mathcal{T}') + I(\chi(u); \chi(v) | \chi(u) \cap \chi(v))$ . To prove (10), we fix one edge  $(\text{parent}(u_i), u_i)$  and repeatedly merge all other edges, until we end with a tree  $\mathcal{T}'$  with two bags,  $\Omega_{1:(i-1)}$  and  $\Omega_{i:m}$  respectively. Then  $\mathcal{J}(\mathcal{T}) \geq \mathcal{J}(\mathcal{T}') = I(\Omega_{1:(i-1)}; \Omega_{i:m} | \Delta_i)$ . The claim follows from the fact that this holds for any  $i = 2, m$ .  $\square$

**Example 5.3.** We illustrate the first part of the theorem on the running example in Fig. 2 and Example 3.4. Enumerating

the nodes depth-first ( $ABD, ACD, AF, BDE$ ), Eq. (9) and (10) become:

$$\begin{aligned}\mathcal{J}(\mathcal{T}) &= I(C; B|AD) + I(F; BCD|A) + I(ACF; E|BD) \\ \max(\dots) &\leq \mathcal{J}(\mathcal{T}) \leq I(CF; BE|AD) + I(F; BCDE|A) + I(ACF; E|BD)\end{aligned}$$

## 5.2 Full MVDs

The number of candidate MVD's is very large: there are<sup>2</sup>  $(3^n + 1)/2 - 2^n = O(3^n)$  standard MVD's  $X \rightarrow Y|Z$ , which is too large to consider for practical datasets. Here, and in the next section, we describe two techniques that allow us to restrict the search space. Consider a fixed key  $X$ . Beeri et al. [4] noted that, in the exact case, if any MVD  $X \rightarrow \dots$  holds on the data, then there exist a “best” one. For example if both  $X \rightarrow AB|C$  and  $X \rightarrow A|BC$  hold exactly, then so does  $X \rightarrow A|B|C$ , and it suffices to discover only the latter. Unfortunately, this fails for approximate MVDs, as we explain here.

We say that  $\phi = X \rightarrow A_1| \dots |A_m$  refines  $\psi = X \rightarrow B_1| \dots |B_k$ , denoted by  $\phi \geq \psi$  if they both have the same key (i.e.,  $\text{key}(\phi) = \text{key}(\psi) = X$ ) and for every  $A_i \in \text{dep}(\phi)$  there exists  $B_j \in \text{dep}(\psi)$  such that  $A_i \subseteq B_j$ . For example,  $X \rightarrow A|B|C$  refines  $X \rightarrow AB|C$ .

**PROPOSITION 5.2.** *If  $\phi \geq \psi$  then  $\mathcal{J}(\phi) \geq \mathcal{J}(\psi)$ .*

**PROOF.** It suffices to consider the case when two dependents in  $\phi$  are replaced by their union in  $\psi$ , e.g.  $\phi = X \rightarrow A|B| \dots$  and  $\psi = X \rightarrow AB| \dots$ , since any refinement is a sequence of such steps. In that case, by inspecting Eq.(6) we observe  $\mathcal{J}(\phi) = \mathcal{J}(\psi) + H(XA) + H(XB) - H(XAB) - H(X) = \mathcal{J}(\psi) + I(A; B|X) \geq \mathcal{J}(\psi)$  proving the claim.  $\square$

We say that an MVD  $\psi$  is  $\varepsilon$ -full, or simply full, if  $R \models_{\varepsilon} \psi$  and, for all strict refinements  $\phi > \psi$ ,  $R \not\models_{\varepsilon} \phi$ . We denote by  $\text{FULLMVD}_{\varepsilon}(R, X)$  the set of all full  $\varepsilon$ -MVDs with key  $X$ . Thus, we only need to discover the sets  $\text{FULLMVD}_{\varepsilon}(R, X)$ , for all  $X \subseteq \Omega$ , because all other MVDs can be derived using Shannon inequalities.

Beeri proved that, in the exact case,  $\text{FULLMVD}_0(R, X)$  has at most one element. We show here what happens in the approximate case, and derive Beeri’s result. Given two MVDs  $\phi = X \rightarrow A_1| \dots |A_m$  and  $\psi = X \rightarrow B_1| \dots |B_k$ , define their join as  $\phi \vee \psi = X \rightarrow C_{11}|C_{12}| \dots |C_{mk}$ , where  $C_{ij} = A_i \cap B_j$ . Clearly,  $\phi \vee \psi$  refines both  $\phi$  and  $\psi$ , i.e.  $\mathcal{J}(\phi \vee \psi) \geq \max(\mathcal{J}(\phi), \mathcal{J}(\psi))$ . We prove a weak form of converse:

**LEMMA 5.4.** *The following are Shannon inequalities:  $\mathcal{J}(\phi \vee \psi) \leq \mathcal{J}(\phi) + m\mathcal{J}(\psi)$  and  $\mathcal{J}(\phi \vee \psi) \leq k\mathcal{J}(\phi) + \mathcal{J}(\psi)$ .*

<sup>2</sup>There are  $3^n$  ways to partition  $\Omega$  into three sets  $X, Y, Z$ . We rule out the  $2^n$  partitions that have  $Y = \emptyset$  and the  $2^n$  partitions that have  $Z = \emptyset$ , and add back the 1 partition that has  $Y = Z = \emptyset$ , for a total of  $3^n - 2^{n+1} + 1$ . Finally, we divide by 2 since  $X \rightarrow Y|Z$  and  $X \rightarrow Z|Y$  are the same MVD.

In the exact case,  $\mathcal{J}(\phi) = \mathcal{J}(\psi) = 0$  implies  $\mathcal{J}(\phi \vee \psi) = 0$ , which proves Beeri’s theorem that  $\text{FULLMVD}_{\varepsilon}(R, X)$  has at most one element, because if  $\phi_1, \phi_2, \dots$  are all MVD’s with key  $X$  that hold exactly on  $R$ , then  $\phi_1 \vee \phi_2 \vee \dots$  refines all of them and holds too. This property was also used by Draeger [12] in his MVD discovery algorithm. When  $\varepsilon > 0$  however, then this fails. For a very simple example, consider

a relation with two tuples,

X	A	B	C
0	0	0	0
0	1	1	1

and fix  $\varepsilon = 1$ .

Then  $R \models_{\varepsilon} X \rightarrow AB|C, X \rightarrow AC|B, X \rightarrow BC|A$ , but  $\not\models_{\varepsilon} X \rightarrow A|B|C$ ; indeed,  $H(\emptyset) = H(X) = 0$  and  $H(W) = 1$  for all other sets  $W$ , and the reader can check  $\mathcal{J}(X \rightarrow AB|C) = \mathcal{J}(X \rightarrow AC|B) = \mathcal{J}(X \rightarrow BC|A) = 1$  but  $\mathcal{J}(X \rightarrow A|B|C) = 2$ .

In summary, our algorithm discovers  $\text{FULLMVD}_{\varepsilon}(R, X)$ , for every  $X$ . Unlike the exact case,  $\text{FULLMVD}_{\varepsilon}(R, X)$  may contain more than one element.

## 5.3 Minimal Separators

We now show that it is not necessary to discover the sets  $\text{FULLMVD}_{\varepsilon}(R, X)$  for all subset of attributes  $X \subset \Omega$ , but only those where  $X$  is a *minimal separator*.

**Definition 5.5.** Fix a relation  $R$  and  $\varepsilon \geq 0$ . We say that a set  $X$  separates two variables  $A, B \notin X$  if there exists an  $\varepsilon$ -MVD  $X \rightarrow Y_1| \dots |Y_m$  that separates  $A, B$ , i.e.  $A, B$  occur in different sets  $Y_i, Y_j$ . We say  $X$  is a *minimal A, B-separator* if there is no  $X_0 \subsetneq X$  that separates  $A, B$ .

For a pair  $A, B \in \Omega$ , we denote by  $\text{MINSEP}_{\varepsilon}(R, A, B)$  the set of minimal  $A, B$  separators in  $R$ , and for a minimal  $AB$  separator  $X$  we denote by  $\text{FULLMVD}_{\varepsilon}(R, X, A, B)$  the set of full MVDs that separate  $A, B$ . Notice that:

$$\text{FULLMVD}_{\varepsilon}(R, X) = \bigcup_{A, B \in \Omega \setminus X} \text{FULLMVD}_{\varepsilon}(R, X, A, B).$$

**Example 5.6.** Let  $R$  be a relation over  $\Omega = \{A, \dots, E\}$ . Suppose  $R \models_{\varepsilon} CD \rightarrow A|BE$ . By (8) we also have  $R \models_{\varepsilon} CDE \rightarrow A|B$ , which means that  $CDE$  cannot be a minimal separator for  $A, B$ . To check that  $CD$  is a minimal  $A, B$ -separator, we need to check that neither  $C$  nor  $D$  separates  $A, B$ .

The main result in this section is that we only need to compute the full MVDs with minimal separators, denoted as:

$$\mathcal{M}_{\varepsilon} \stackrel{\text{def}}{=} \bigcup_{A, B \in \Omega} \bigcup_{\substack{X \in \\ \text{MINSEP}_{\varepsilon}(R, A, B)}} \text{FULLMVD}_{\varepsilon}(R, X, A, B) \quad (11)$$

because, as we show, every  $\varepsilon$ -MVD can be derived from the set  $\mathcal{M}_{\varepsilon}$  by a Shannon inequality.

**THEOREM 5.7.** *Let  $X \rightarrow Y|Z$  be an  $\varepsilon$ -MVD for  $R$ . Then there exist  $\phi_1, \dots, \phi_m \in \mathcal{M}_{\varepsilon}$ , where  $m = |Y| \cdot |Z|$ , such that the following is a Shannon inequality:  $I(Y; Z|X) \leq \sum_i \mathcal{J}(\phi_i)$ .*

---

**Algorithm** MVDMiner( $R, \Omega, \varepsilon$ )

---

```

1:  $\mathcal{M}_\varepsilon \leftarrow \emptyset$ 
2: for all pairs  $A, B \in \Omega$  do
3:    $\text{MINSEP}_\varepsilon(R, A, B) \leftarrow \text{MineMinSeps}(R, \Omega, \varepsilon, (A, B))$ 
4:   for all  $X \in \text{MINSEP}_\varepsilon(R, A, B)$  do
5:      $\mathcal{M}_\varepsilon \leftarrow \mathcal{M}_\varepsilon \cup \text{getFullMVDs}(X, \varepsilon, (A, B), \infty)$ 
6: return  $\mathcal{M}_\varepsilon$ 

```

---

**Fig. 3** Discover the set  $\mathcal{M}_\varepsilon = \cup_{S \in \text{MINSEP}_\varepsilon} \text{FULLMVD}_\varepsilon(S)$ .

In summary, our algorithm will iterate over pairs of attributes  $A, B$ , will compute  $\text{MINSEP}_\varepsilon(R, A, B)$ , then, for each  $X$  in this set will compute  $\text{FULLMVD}_\varepsilon(R, X, A, B)$ , and return their union,  $\mathcal{M}_\varepsilon$ ; we describe it in the next section. We end this section with the proof of Theorem 5.7.

**PROOF.** Let  $Y = A_1 \dots A_m$ , and  $Z = B_1 \dots B_k$ . By the chain rule (4) it holds that:

$$I(Y; Z|X) = \sum_{i=1}^m \sum_{j=1}^k I(A_i; B_j | X A_1 \dots A_{i-1} B_1 \dots B_{j-1})$$

It suffices to prove that, for each  $i, j$ , there exists an MVD  $\phi \in \mathcal{M}_\varepsilon$  such that the following is a Shannon inequality:

$$I(A_i; B_j | X A_1 \dots A_{i-1} B_1 \dots B_{j-1}) \leq \mathcal{J}(\phi)$$

Since  $X \rightarrow Y|Z$  is a  $\varepsilon$ -MVD for the relation  $R$ , then  $X$  is an  $A_i, B_j$  separator. Let  $S \subseteq X$  be any minimal  $A_i, B_j$  separator, thus  $S \in \text{MINSEP}_\varepsilon(R, A_i, B_j)$ , and let  $\phi = S \rightarrow U_1 | \dots | U_p$  be a full MVD in  $\text{FULLMVD}_\varepsilon(R, S, A_i, B_j) \subseteq \mathcal{M}_\varepsilon$  that separates  $A_i, B_j$ . Assume w.l.o.g.  $A_i \in U_1, B_j \in U_2$ , and let  $\psi \stackrel{\text{def}}{=} S \rightarrow W|V$ , where  $W = U_1, V = U_2 U_3 \dots U_p$ . Thus,  $\phi \geq \psi$ , and therefore by Prop. 5.2 the following Shannon inequality holds:  $\mathcal{J}(\phi) \geq \mathcal{J}(\psi)$ . Write  $\psi$  as  $\psi = S \rightarrow W_0 W_1 | V_0 V_1$ , where  $W_0 = W \cap (X A_1 \dots A_i B_1 \dots B_j), W_1 = W - W_0$ , and similarly  $V_0 = V \cap (X A_1 \dots A_i B_1 \dots B_j), V_1 = V - V_0$ . By Prop. 5.1 (7) we have the following Shannon inequality  $\mathcal{J}(\psi) = \mathcal{J}(S \rightarrow W_0 W_1 | V_0 V_1) \geq \mathcal{J}(S \rightarrow W_0 | V_0)$ . Finally, we notice that the set  $S W_0 V_0$  is the same as  $X A_1 \dots A_i B_1 \dots B_j$  and that  $A_i \in W_0, B_j \in V_0$ , therefore by Prop. 5.1, (8),  $\mathcal{J}(S \rightarrow W_0 | V_0) \geq \mathcal{J}(X A_1 \dots A_{i-1} B_1 \dots B_{j-1} \rightarrow A_i | B_j)$ , proving the claim.  $\square$

## 6 DISCOVERING $\varepsilon$ -MVDs

In this section we present the first phase of Maimon: the algorithm for the discovery of  $\varepsilon$ -MVDs in a relation  $R$ , called MVDMiner, and shown in Figure 3. As explained, the algorithm returns the set  $\mathcal{M}_\varepsilon$ , defined in Eq.(11); this set is used in the second phase of Maimon to compute  $\varepsilon$ -schemes.

MVDMiner iterates over all pairs of attributes  $A, B \in \Omega$ . It first computes the set  $\text{MINSEP}_\varepsilon(R, A, B)$  of minimal  $A, B$ -separators (line 3): we describe this step in Sec. 6.1. Then, for each  $X \in \text{MINSEP}_\varepsilon(R, A, B)$ , it computes  $\text{FULLMVD}_\varepsilon(R, X, A, B)$  (line 5): we describe this step in Sec. 6.2. Finally, the algorithm returns their union,  $\mathcal{M}_\varepsilon$ . Both steps require access to an oracle  $\text{getEntropy}_R(X)$  for computing the entropy  $H(X)$ , according to Eq. (5), where  $H$  is the entropy associated with the empirical distribution over  $R$ . We describe the implementation and optimization of  $\text{getEntropy}_R(X)$  in Section 6.3.

### 6.1 Discovering the Minimal Separators

We describe here how we compute all minimal  $A, B$ -separators,  $\text{MINSEP}_\varepsilon(R, A, B)$  (line 3 of MVDMiner). One possible way to do this could be to iterate over sets  $X$  top down, because it enables pruning: if  $X$  is not an  $A, B$ -separator, then neither is any subset of  $X$ , by (8) in Prop. 5.1. This suggests a top-down algorithm, which starts from the largest set  $X = \Omega \setminus \{A, B\}$ , and checks if it is an  $A, B$ -separator. If not, then none exists. Otherwise it exhaustively searches over subsets of  $X$ , from largest to smallest, returning the minimal (with regard to inclusion) sets that separate  $A, B$ . Such an exhaustive search will explore *all* separators, while we only want to find the *minimal* ones. Our approach takes advantage of the fact that we need to find only the *minimal* separators, and builds on a result by Gunopulos et al. [20].

Let  $\mathbf{C} = \{C_1, \dots, C_m\}$  be a set of distinct subsets of  $\Omega$ . A set  $D \subset \Omega$  is a *transversal* of  $\mathbf{C}$  if  $D \cap C_i \neq \emptyset$  for every  $C_i \in \mathbf{C}$ . For a set  $D \subseteq \Omega$ , we denote by  $\bar{D}$  the complement set  $\Omega \setminus D$ .

**THEOREM 6.1.** Let  $\mathbf{C} = \{C_1, \dots, C_n\}$  denote a set of minimal  $A, B$  separators in  $R$ . Then there exists a minimal  $A, B$ -separator  $X \notin \mathbf{C}$  iff there exists a minimal transversal  $D$  of  $\mathbf{C}$  such that  $\bar{D}$  is an  $A, B$ -separator.

**PROOF.** **only if.** Since  $D$  is a transversal of  $\mathbf{C}$  then:

$$\bigwedge_{i=1}^n (C_i \cap D \neq \emptyset) \iff \bigwedge_{i=1}^n (\bar{D} \not\supseteq C_i) \quad (12)$$

Since  $\bar{D}$  is an  $A, B$  separator, there exists some minimal separator  $X \subseteq \bar{D}$ . Assume, by contradiction, that  $X \supseteq C_i$  for some  $C_i \in \mathbf{C}$ . Then  $\bar{D} \supseteq X \supseteq C_i$ , contradicting (12).

**if.** Since  $X$  is a minimal  $A, B$  separator that is not in  $\mathbf{C}$ , then  $\bigwedge_{i=1}^n (X \not\supseteq C_i)$ , meaning that  $\bar{X}$  is a transversal of  $\mathbf{C}$ . Then any minimal transversal  $D \subseteq \bar{X}$  satisfies the claim.  $\square$

Algorithm MineMinSeps (Fig. 5) for discovering all minimal  $A, B$  separators,  $\text{MINSEP}_\varepsilon(R, A, B)$  is based on Theorem 6.1, and proceeds as follows:

- (1) Initialize  $\mathbf{C}$  with a single minimal  $A, B$ -separator (Line 3-5).
- (2) Iterate over all minimal transversals  $D$  of  $\mathbf{C}$  (Line 8):

- (3) If  $\bar{D}$  separates  $A, B$  (Line 11), then:
- Find any minimal  $A, B$  separator  $X \subseteq \bar{D}$  (Line 12).
  - $C \leftarrow C \cup \{X\}$ .

The function `ReduceMinSep` called in lines 4 and 12 takes a separator ( $\Omega \setminus \{A, B\}$  or  $\bar{D}$  respectively) and finds *any* subset that is a minimal separator; this is done greedily in `ReduceMinSep` (Fig. 4). The function `getFullMVDs` called in line 10 of `MineMinSeps`, and in line 4 of `ReduceMinSep`, takes as input an attribute set  $X$ , a pair of attributes  $A, B$ , and a threshold  $\varepsilon$ , and computes full  $\varepsilon$ -MVDs with key  $X$  that separate  $A, B$ ; a parameter  $K > 0$  is used to limit the number of full MVDs returned, and here we set  $K = 1$  because we only check if one exists; in line 5 of the main algorithm (Fig. 3) we set  $K = \infty$ .

---

**Algorithm** `ReduceMinSep`( $\varepsilon, X, (A, B)$ )

---

- 1: Let  $p = X_1, \dots, X_m$  be a predefined ordering of  $X$ .
  - 2:  $S \leftarrow X$
  - 3: **for all**  $i = 1$  to  $m$  **do**
  - 4:    $M_i \leftarrow \text{getFullMVDs}(S \setminus \{X_i\}, \varepsilon, (A, B), 1)$
  - 5:   **if**  $M_i \neq \emptyset$  **then**
  - 6:      $S \leftarrow S \setminus \{X_i\}$
  - 7: **return**  $S$
- 

**Fig. 4** Given a set  $X \subset \Omega$ , and a pair  $(A, B) \in \Omega \setminus X$ , find a subset  $S \subseteq X$  s.t.  $S$  is a minimal  $A, B$ -separator in  $R$ .

The only sets of attributes returned in `MineMinSeps` are minimal  $AB$ -separators returned by `ReduceMinSep` in lines 4 and 12. The proof of completeness (i.e., the algorithm returns all minimal  $AB$ -separators) follows techniques similar to those by Gunopulos et al. [20], and is given in the full version of the paper:

**THEOREM 6.2.** *Algorithm MineMinSeps in Figure 5 enumerates all minimal  $A, B$ -separators in  $R$ .*

We now analyze the runtime between consecutive discoveries of minimal  $A, B$ -separators in `MineMinSeps`. We let  $\Omega$  be a finite set of cardinality  $n$ , and let  $C \subseteq 2^\Omega$  be a finite set of sets. The problem of discovering all minimal transversals of  $C$  is called the *hypergraph transversal problem* [23]. The theoretically best known algorithm for solving the hypergraph transversal problem is due to Fredman and Khachiyan [16] and has a quasi incremental-polynomial delay of  $\text{poly}(n) + m^{O(\log^2 m)}$  where  $m = |C| + n$ . Note the dependence on the size of the discovered minimal separators  $|C|$ . We denote by  $T_{\minTrans}(n, C)$  the delay of the minimal transversal algorithm. However, not every minimal transversal  $D$  leads to the discovery of a minimal separator if  $\bar{D}$  does not separate  $A$  and  $B$  (i.e.,  $\phi = \emptyset$  in line 11 of `MineMinSeps`).

---

**Algorithm** `MineMinSeps`( $R, \Omega, \varepsilon, (A, B)$ )

---

- 1:  $C \leftarrow \emptyset$
  - 2:  $X \leftarrow \text{nil}$
  - 3: **if**  $I(A; B | \Omega \setminus \{A, B\}) \leq \varepsilon$  {by `getEntropyR`} **then**
  - 4:    $X \leftarrow \text{ReduceMinSep}(\varepsilon, \Omega \setminus \{A, B\}, (A, B))$
  - 5:    $C \leftarrow C \cup \{X\}$
  - 6: **else**
  - 7:   Return  $\emptyset$
  - 8: **while**  $(D \leftarrow \text{nextMinTransversal}(C)) \neq \text{nil}$  **do**
  - 9:    $\bar{D} \leftarrow \Omega \setminus D$
  - 10:    $\phi \leftarrow \text{getFullMVDs}(\bar{D}, \varepsilon, (A, B), 1)$
  - 11:   **if**  $\phi \neq \emptyset$  **then**
  - 12:      $X \leftarrow \text{ReduceMinSep}(\varepsilon, \bar{D}, (A, B))$
  - 13:      $C \leftarrow C \cup \{X\}$
  - 14: **return**  $C$
- 

**Fig. 5** Given a relation  $R$  with schema  $\Omega$ , two attributes  $A, B \in \Omega$ , and a threshold  $\varepsilon$  enumerate all minimal  $A, B$ -separators in  $R$ .

In the full version of this paper we show that the number of minimal transversals processed in lines 9-13 before a new minimal separator is discovered (e.g., in line 12), or before the loop exists, is bounded by  $n \cdot |C|$ . This allows us to formalize the delay between the discovery of minimal  $A, B$ -separators. We denote by  $T(\text{getFullMVDs})$  the runtime of `getFullMVDs`, which we analyze in the next section.

**COROLLARY 6.3.** *Algorithm MineAllMinseps enumerates the minimal  $A, B$ -separators in  $R$  with a delay of  $O(n \cdot |C| \cdot T_{\minTrans}(n, C) \cdot T(\text{getFullMVDs}))$ , where  $n = |\Omega|$ .*

## 6.2 Discovering the Full MVDs

Returning to our main algorithm, `MVDMiner`, we have shown how to compute  $\text{MINSEP}_\varepsilon(R, A, B)$ , the set of minimal  $A, B$  separators in  $R$ . Next, for each minimal  $A, B$  separator  $X \in \text{MINSEP}_\varepsilon(R, A, B)$ , we compute all full MVDs with key  $X$  that separate  $A$  and  $B$ , i.e. the set  $\text{FULLMVD}_\varepsilon(R, X, A, B)$ ; this is line 5 of `MVDMiner`. Recall that *full* means that the MVD cannot be further refined.

The algorithm `getFullMVDs` starts by checking the most refined MVD with key  $X$ , namely  $\varphi = X \rightarrow Y_1 | \dots | Y_n$  where  $Y_1, \dots, Y_n$  are all attributes not in  $X$  (including  $A, B$ ). If  $J(\varphi) \leq \varepsilon$  then we are done. Otherwise, the algorithm considers all possible ways to *merge* two dependents, while keeping  $A$  and  $B$  in different dependents; i.e. it tries  $X \rightarrow Y_1 Y_2 | \dots | Y_n, X \rightarrow Y_1 Y_3 | Y_2 | \dots | Y_n$ , etc. We denote the MVD that results from merging dependents  $Y_i$  and  $Y_j$  in  $\text{dep}(\varphi)$

**Algorithm** getFullMVDs( $S, \varepsilon, (A, B), K$ )

---

```

1:  $\mathcal{P} \leftarrow \emptyset$  {Output set}
2:  $Q \leftarrow \emptyset$  { $Q$  is a stack}
3:  $\phi_0 = S \rightarrow X_1 | \dots | X_n$  where  $X_i$  are singletons.
4:  $Q.push(\phi_0)$ 
5: while  $Q \neq \emptyset$  and  $|\mathcal{P}| < K$  do
6:    $\varphi \leftarrow Q.pop()$ 
7:   Compute  $\mathcal{J}(\varphi)$  {using getEntropyR}
8:   if  $\mathcal{J}(\varphi) \leq \varepsilon$  then
9:      $\mathcal{P} \leftarrow \mathcal{P} \cup \{\varphi\}$ 
10:    else
11:      for all  $\phi \in Nbr(\varphi)$  do
12:         $Q.push(\phi)$  {See (13)}
13: return  $\mathcal{P}$ 

```

---

**Fig. 6** Returns a set of at most  $K$  full MVDs with key  $S$  that approximately hold in  $R$  (w.r.t  $\varepsilon$ ) in which  $A$  and  $B$  are in distinct components.

by  $\text{merge}_{ij}(\varphi)$ . Since  $\varphi$  refines  $\text{merge}_{ij}(\varphi)$  then, by Proposition 5.2, it holds that  $\mathcal{J}(\text{merge}_{ij}(\varphi)) \leq \mathcal{J}(\varphi)$ . This procedure for searching for a full  $\varepsilon$ -MVD can be viewed as a graph traversal algorithm where every node  $\varphi$  is an  $\varepsilon$ -MVD candidate with key  $X$ , dependents  $Z_1, \dots, Z_k$ , and its neighbors  $Nbr(\varphi)$  are the  $\varepsilon$ -MVD candidates:

$$Nbr(\varphi) \stackrel{\text{def}}{=} \{\text{merge}_{ij}(\varphi) : Z_i, Z_j \in \text{dep}(\varphi), A, B \notin Z_i Z_j\} \quad (13)$$

Clearly, if  $A, B$  were separated in  $\varphi$ , then they remain separated in every MVD in  $Nbr(\varphi)$ . We present the algorithm as a depth-first traversal, which is how we implemented it. The pseudocode is presented in Figure 6.

**6.2.1 An Optimization to getFullMVDs.** In the worst case, if  $S$  is not an  $A, B$ -separator then Algorithm getFullMVDs will traverse the complete search space of size  $O(2^n)$ . While, in general, this is unavoidable, we implemented an optimization, described in the complete version of this paper, that leads to a significant reduction in the search space.

### 6.3 Computing Entropies Efficiently

We describe the procedure getEntropy<sub>R</sub> for calculating the joint entropy of a set of attributes. The efficiency of this procedure is crucial to the performance of MVDMiner, which needs to repeatedly compute mutual information values  $I(Y; Z|X)$ , and each such computation requires four entropic values  $H(XY)$ ,  $H(XZ)$ ,  $H(XYZ)$ , and  $H(X)$ . Repeatedly computing values of the form  $H(X_\alpha)$ , for  $\alpha \subseteq [n]$  requires multiple scans over the data that resides in external memory.

We build on ideas introduced in the PLI cache data structure [21, 27], and reduce the problem of computing  $H(X_\alpha)$  to a main memory join-group-by query. To describe the algorithm, we repeat here the entropy formula (5) for convenience:

$$H(X_\alpha) \stackrel{\text{def}}{=} \log N - \frac{1}{N} \sum_{x_\alpha \in \mathcal{D}_\alpha} |R(X_\alpha=x_\alpha) \log |R(X_\alpha=x_\alpha)| \quad (14)$$

The algorithm uses two ideas: (1) if  $x_\alpha$  is a singleton (i.e., its frequency  $|R(X_\alpha=x_\alpha)|=1$ ) then it can be ignored because its contribution to the total entropy in (14) is 0 (due to the logarithm), and (2) given two relations mapping the distinct values of attribute sets  $X_\alpha$ , and  $X_\beta$ , respectively, to the tuple ids in the relation  $R$  that contain them, then we can derive this mapping for  $X_\alpha \cup X_\beta$  by simply joining the two mappings on the tuple IDs. Ignoring singleton valuations makes these mappings highly compressed, enabling us to store them in main memory and perform the join using a main memory database system. We used the in-memory database  $H2$  [41]. We describe the details next. We let  $Hash$  denote a hash function. In our implementation we use the hash function provided by the database system.

Alg. getEntropy<sub>R</sub> maintains two sets of relations indexed by  $\alpha \subseteq [n]$ :  $CNT_\alpha(\text{val}, \text{cnt})$  and  $TID_\alpha(\text{val}, \text{tid})$  defined as:

$$CNT_\alpha = \{(Hash(x_\alpha), \text{cnt}) \mid \text{cnt} = |R(X_\alpha = x_\alpha)|, \text{cnt} > 1\}$$

$$TID_\alpha = \{(Hash(x_\alpha), t[\text{tid}]) \mid t \in R, t[X_\alpha] = x_\alpha, Hash(x_\alpha) \in \Pi_{\text{val}}(CNT_\alpha)\}$$

We compute  $H(X_\alpha)$  by scanning table  $CNT_\alpha$ . The algorithm starts by computing two sets of relations: (1)  $\{CNT_{\{i\}}\}$  and (2)  $\{TID_{\{i\}}\}$  for every  $i \in [n]$ . Assume that we have computed the relations  $CNT_\alpha$ ,  $CNT_\beta$  and  $TID_\alpha$ ,  $TID_\beta$  for some subsets  $\alpha, \beta \subset [n]$  such that  $\alpha \cap \beta = \emptyset$ . We first compute  $CNT_{\alpha \cup \beta}$  as:

```

Select Hash(A.val, B.val) as val, count(*) as cnt
From TID\alpha A, TID\beta B
Where A.tid = B.tid
Group By Hash(A.val, B.val) Having count(*) > 1

```

Next, we compute  $TID_{\alpha \cup \beta}$  as:

```

Select Hash(A.val, B.val) as val, A.tid as tid
From TID\alpha A, TID\beta B, CNT\alpha \cup \beta Z
Where A.tid = B.tid and Hash(A.val, B.val) = Z.val

```

Pruning the singleton values makes this technique very effective, because as we move up the lattice from smaller  $\alpha$ 's to larger  $\alpha$ 's, many more tuples  $x_\alpha$  are unique in the data, and the tables  $CNT_\alpha$  and  $TID_\alpha$  become smaller.

*Example 6.4.* For a simple illustration, Fig. 7 shows the tables generated for a 3-attribute relation  $R$ . Both types of relations only contain values corresponding to non-singleton valuations in  $R$ .

However, even with our compression, generating and storing all  $2^n - 1$  tables  $\text{CNT}_\alpha$ , and  $\text{TID}_\alpha$  is intractable. Instead, we perform the following optimization. Fix a parameter  $L$  (in our implementation we chose  $L = 10$ ), and partition the set  $\Omega$  into  $\lceil \frac{n}{L} \rceil$  disjoint subsets  $\Omega_1, \Omega_2, \dots$  each of size at most  $L$ . For each  $i$ , compute the tables  $\text{TID}_\alpha$  and  $\text{CNT}_\alpha$  for all subsets  $\alpha \subseteq \Omega_i$ ; thus the total number of tables precomputed is  $2 \lceil \frac{n}{L} \rceil \cdot 2^L$ . In order to compute  $H(X_\alpha)$ , we express  $\alpha = (\alpha \cap \Omega_1) \cup (\alpha \cap \Omega_2) \cup \dots$ , where each union is treated as explained above for  $\alpha \cup \beta$ .

R	tid	A	B	C		CNT <sub>A</sub>	CNT <sub>B</sub>	CNT <sub>C</sub>			
						val	CNT	val	CNT	val	CNT
$t_1$	$a_1$	$b_2$	$c_3$								
$t_2$	$a_2$	$b_1$	$c_1$								
$t_3$	$a_2$	$b_2$	$c_2$			$a_2$	2	$b_2$	2	$c_3$	2
$t_4$	$a_3$	$b_3$	$c_3$			$a_3$	2	$b_3$	2		
$t_5$	$a_3$	$b_3$	$c_4$								
$\text{CNT}_{AB}$						TID <sub>A</sub>		TID <sub>B</sub>		TID <sub>C</sub>	
						val	tid	val	tid	val	tid
						$a_2$	$t_2$	$b_2$	$t_1$	$c_3$	$t_4$
$\text{Hash}(a_3, b_3)$				2		$a_2$	$t_3$	$b_2$	$t_3$	$c_3$	$t_4$
$\text{TID}_{AB}$						$a_3$	$t_4$	$b_3$	$t_4$		
						$a_3$	$t_5$	$b_3$	$t_5$		
$\text{Hash}(a_3, b_3)$				$t_4$							
$\text{Hash}(a_3, b_3)$				$t_5$							

Fig. 7 getEntropy<sub>R</sub> example.

---

#### Algorithm ASMiner( $\mathcal{M}_\epsilon$ )

---

```

1: schemes =  $\emptyset$ 
2: Construct the graph  $G = \{(\phi, \psi) \mid \phi, \psi \in \mathcal{M}_\epsilon, \phi \# \psi\}$ 
3: for all  $Q \in \text{MaxIndependentSet}(G)$  do
4:   schemes  $\leftarrow$  schemes  $\cup$  {BuildAcyclicSchema( $Q$ )}
5: return schemes

```

---

Fig. 8 Generate Acyclic Schemas from  $\mathcal{M}_\epsilon$ .

## 7 ENUMERATING ACYCLIC SCHEMAS

In this section we present the second phase of Maimon: given the set  $\mathcal{M}_\epsilon$  of full  $\epsilon$ -MVDs (Eq. (11)), generate acyclic  $\epsilon$ -schemas. The algorithm ASMiner is shown in Fig. 8. It searches for subsets of MVDs  $Q \subseteq \mathcal{M}_\epsilon$ , and reconstructs a schema from that set. The key to the algorithm's efficiency is our new definition of compatibility:

**Definition 7.1.** Let  $\phi_1 = X \twoheadrightarrow A_1 | \dots | A_m$  and  $\phi_2 = Y \twoheadrightarrow B_1 | \dots | B_k$  be two  $\epsilon$ -MVDs. We say that  $\phi_1$  and  $\phi_2$  are *compatible* if there exist an  $i \in \{1, \dots, m\}$ , and  $j \in \{1, \dots, k\}$  such that:

- (1)  $Y \subseteq XA_i$ , and  $X \subseteq YB_j$ . In this case we say that the two MVDs are *split-free* [6, 15, 19, 28].

---

#### Algorithm BuildAcyclicSchema( $Q$ )

---

```

1:  $S \leftarrow \{\Omega\}$ 
2: Sort  $Q$  by ascending order of key cardinality {e.g.,  $X \twoheadrightarrow A|B$  before  $XY \twoheadrightarrow C|D$ }
3: for all  $\phi \in Q$  do
4:   Let  $\phi = X \twoheadrightarrow C_1 | \dots | C_m$ 
5:   Let  $\Omega_i \in S$  s.t.  $X \subseteq \Omega_i$ 
6:    $D_\phi \leftarrow \{C_j X \cap \Omega_i \mid j \in [i, m]\} \setminus \{X\}$ 
7:   if  $|D_\phi| \geq 2$  then
8:     Replace  $\Omega_i \in S$  with  $D_\phi$  { $\phi$  is non-redundant}
9: return  $S$ 

```

---

Fig. 9 Gets a set  $Q$  of pairwise compatible MVDs, and returns an acyclic schema.

- (2) There exist two distinct indexes  $j_1, j_2 \in \{1, \dots, k\}$  such that  $XA_i \cap B_{j_1} \neq \emptyset$ , and  $XA_i \cap B_{j_2} \neq \emptyset$ . Likewise, there exist two distinct indexes  $i_1, i_2 \in \{1, \dots, m\}$  such that  $YB_j \cap A_{i_1} \neq \emptyset$ , and  $YB_j \cap A_{i_2} \neq \emptyset$ .

We write  $\phi_1 \# \phi_2$  to denote the fact that  $\phi_1, \phi_2$  are *incompatible*.

We say that a set  $Q$  of  $\epsilon$ -MVDs is *pairwise compatible* if every pair of  $\epsilon$ -MVDs in  $Q$  is compatible. Recall that every join tree  $\mathcal{T}$  with  $m$  nodes defines a set of  $m - 1$  MVDs called its *support* and denoted by  $\text{MVD}(\mathcal{T})$ .

**THEOREM 7.2.** Let  $S$  be an acyclic schema with join tree  $(\mathcal{T}, \chi)$ . Then the set  $\text{MVD}(\mathcal{T})$  is pairwise compatible.

Thus, it suffices to iterate over sets of pairwise compatible  $\epsilon$ -MVDs. Specifically, our algorithm enumerates the *maximal* sets of pairwise compatible  $\epsilon$ -MVDs, and for this task we use a graph algorithm from the literature. Define the graph  $G(\mathcal{M}_\epsilon, E)$  as follows:

$$E = \{(\phi_1, \phi_2) : \phi_1, \phi_2 \in \mathcal{M}_\epsilon \text{ and } \phi_1 \# \phi_2\} \quad (15)$$

By this definition every maximal independent set in  $G$  corresponds to a maximal set of pairwise compatible  $\epsilon$ -MVDs. We apply the following result.

**THEOREM 7.3.** ([11, 22]) Let  $G(V, E)$  be a graph. The maximal independent sets of  $G$  can be enumerated such that the delay between consecutive outputs is in  $O(|V|^3)$ .

In summary, algorithm ASMiner in Fig. 8 enumerates all maximal independent sets  $Q$ , then for each of them constructs one acyclic schema  $S$ , by calling BuildAcyclicSchema shown in Fig. 9, and described next.

Algorithm BuildAcyclicSchema starts with a schema that contains a single relation with all attributes (i.e.,  $S = \{\Omega\}$ ). It then builds the acyclic schema for  $R$  by repeatedly using an

$\varepsilon$ -MVD from  $Q$  to decompose one of the relations in  $S$ . The MVDs are processed in ascending order of the cardinality of their keys. Therefore, when an MVD  $S \rightarrow C_1 | \dots | C_m$  is processed, then we know that  $S$  is contained in exactly one of the relations in  $S$  (e.g., otherwise,  $S$  must be contained in a key of a previously processed  $\varepsilon$ -MVD). The algorithm then applies this  $\varepsilon$ -MVD to the single relation that contains it, and continues until all  $\varepsilon$ -MVDs in  $Q$  have been processed. An MVD is said to be *redundant* [18] if it does not split the single relation that contains it (i.e., condition of line 7 does not hold). Redundant MVDs are simply ignored in `BuildAcyclicSchema`.

**THEOREM 7.4.** *Algorithm `BuildAcyclicSchema` generates an acyclic schema  $S$  with join tree  $(\mathcal{T}, \chi)$  such that  $MVD(\mathcal{T}) \subseteq Q$ . If  $Q$  is a non-redundant set of  $\varepsilon$ -MVDs then  $MVD(\mathcal{T}) = Q$ . The algorithm runs in time  $O(n^3)$ .*

The novel insight of our algorithm is the characterization of (in)compatibility in Definition 7.1, which depends only on the pairwise relationship between the MVDs, and therefore enables the reduction to enumerating maximal independent sets in graphs. Previous characterizations [6, 15, 19, 28] are for entire sets of MVDs, and are not pairwise (more precisely, they have a different second condition called *intersection* which relies on the existence of a third MVD in the set). Goodman and Tay [18] present an algorithm for synthesizing an acyclic schema from a set  $Q$  of MVDs that satisfy the *subset property*. As in Theorem 7.4, they show that if the set  $Q$  is non-redundant then the synthesized acyclic schema has a join tree whose support is precisely the set  $Q$ . However, we are not aware of any characterization of non-redundant MVDs. While the subset property is pairwise, it is applicable only to binary MVDs, while our MVDs may have any number of dependents. Algorithms for constructing a (single) acyclic schema from data dependencies have been previously developed by Bernstein [7] where the input is a set of functional dependencies, and by Beeri et al. and Lien whose algorithms work by combining *conflict-free* MVDs [6, 32].

## 8 EVALUATION

In this section we conduct an experimental evaluation of `Maimon`. We start by evaluating the relationship between the threshold  $\varepsilon$ , and the number of spurious tuples generated by the acyclic schema in Section 8.2. Next, we evaluate the efficiency and scalability of the algorithms, starting with mining the minimal separators, which forms the most critical part of the algorithm, in Section 8.3. We then measure the cost of generating all full approximate MVDs in Section 8.4. Finally, in Section 8.5 we evaluate the efficiency of the enumeration algorithm for approximate acyclic schemes.

Table 2 presents the 20 datasets we used for the evaluation along with the time it took to discover the set of full MVDs when the threshold is set to  $\varepsilon = 0$ . We ran these experiments

Dataset	Dataset			Full MVDs threshold=0.0	
	Dataset	Cols.	Rows	Runtime [sec]	Full MVDs
Ditag Feature	13	3960124	TL	NA	
Four Square (Spots)	15	973516	17017	105	
Image	12	777676	3747	151	
FD_Reduced_30	30	250000	8024	21	
FD_Reduced_15	15	250000	1006	21	
Census	42	199524	TL	NA	
SG_Bioentry	7	184292	101	3	
Atom Sites	26	160000	TL	242	
Classification	12	70859	1327	27	
Adult	15	32561	1083	58	
Entity Source	33	26139	14155	153	
Reflns	27	24769	TL	543	
Letter	17	20000	605	44	
School Results	27	14384	7202	2394	
Voter State	45	10000	TL	262	
Abalone	9	4177	602	36	
Breast-Cancer	11	699	5	30	
Hepatitis	20	155	479	2953	
Echocardiogram	13	132	6	104	
Bridges	13	108	3.8	60	

**Table 2: Evaluation datasets sorted by number of rows. We counted only columns with non-constant attributes containing more than one value. We present runtimes (in seconds) for mining full MVDs with threshold 0.0. We report on the number of full MVDs generated within a time limit (TL) of 5 hours.**

with a time limit of 5 hours. For three datasets, Atom Sites, REFLNS, and Voter State, our algorithm did not complete within 5 hours, yet it discovered a large number of full MVDs. For the two largest datasets DITAG Feature, and Census, the algorithm did not find MVDs on the entire data set, yet we include them because we used subsets in our scalability experiments, Section 8.3.

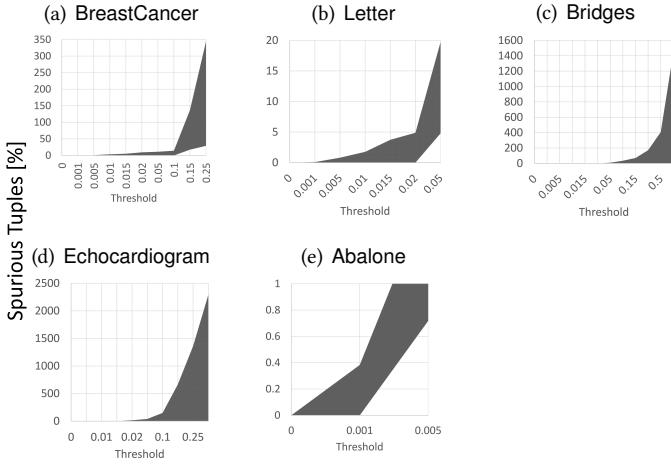
### 8.1 Experimental Setup

Our algorithm was implemented in Java 1.8. We perform all experiments on a 64bit Linux machine running Ubuntu 5.4.0. The machine has 120 CPUs and 1 TB of memory. Our algorithm is single-threaded, and thus runs on a single core. We evaluated our algorithm over real-world datasets<sup>3</sup> that are part of the Metanome data profiling project [34].

### 8.2 Schema Accuracy

We evaluate the accuracy of the enumerated acyclic schemes by examining the relationship between the threshold used

<sup>3</sup><https://hpi.de/naumann/projects/repeatability/data-profilng/fds.html#c168191>



**Fig. 10** The Y-axis represents the the number of spurious tuples as a percentage of the number of rows in the dataset. For every threshold we show the range of these percentages over all acyclic schemas generated for the dataset in that threshold (e.g., in *Letter* the percentage of spurious tuples for  $\varepsilon = 0.02$  is between 0% (minimum), and 5% (maximum)).

for generating the MVDs, and the number of *spurious tuples* generated by the decomposition defined by the acyclic schemes. Formally, for every generated acyclic schema  $S = \{\Omega_1, \dots, \Omega_m\}$  we calculate  $Err(S) = \frac{|R[\Omega_1] \bowtie \dots \bowtie R[\Omega_m]| - |R|}{|R|}$ . By Theorem 3.3 it holds that  $Err(S) = 0$  precisely when  $J(S) = 0$ . When  $\varepsilon = 0$  the acyclic schemas generated by the algorithm are perfect, and produce no spurious tuples. We ran the acyclic schema generation algorithm over several datasets, with different threshold (i.e.,  $\varepsilon$ ) values ranging between 0.0 and 0.5. For every dataset and every threshold we calculated the minimum and maximum number of spurious tuples in any of the schemes returned. The results are plotted in Figure 10. We see a clear and consistent relationship between the threshold and the number of spurious tuples, showing that the information theoretic measures precisely capture the error of the approximate acyclic schema.

### 8.3 Scalability Experiments

In this section we evaluate how the algorithm scales with the data, in particular with the number of rows, and the number of columns. We focus on the task of mining all minimal separators because this is the most expensive step of the algorithm.

In this experiment, we evaluate the scalability of the algorithm that mines all of the minimal separators from the relation for a given threshold  $\varepsilon$ . Notice that increasing the threshold does not necessarily increase nor decrease the number of minimal separators. To see why, consider a dataset

where no MVD holds exactly. Then, when  $\varepsilon = 0$  there are 0 minimal separators, when  $\varepsilon = \log N$  (for an  $N$ -tuple relation), then *every* MVD holds with this threshold, and  $\emptyset$  is the sole minimal separator, while when  $\varepsilon$  is in between, then we can get a arbitrary number of minimal separators.

**8.3.1 Row Scalability.** We evaluated the algorithm over three large datasets: *Image*, *foursquare*, and *Ditag Feature* (see Table 2). We executed the experiment with all columns of the datasets while taking between 10% to 100% of the tuples. The charts are presented in Figure 11. We note that for all thresholds we see that the runtime scales linearly with the dataset size. For both *Image* and *Ditag Feature* there are many more minimal separators generated with a threshold of 0.0, but this seems to have little affect on the runtime.

**8.3.2 Column Scalability.** We evaluated the algorithm over three datasets with a large number of columns. We executed the experiment with all rows of the datasets while taking between 10% to 100% of the columns. The results are presented in Figure 12. For *ncVoterState* the results for thresholds 0.0 and 0.01 are almost identical, so we display only the former. These datasets are particularly challenging for our algorithm, but even in the cases in which it timed out, the algorithm generated many minimal separators (e.g., 262 minimal separators for the full 45-column version of *Voter State*). In all three experiments, we see a high correlation between the running time and the number of minimal separators generated. A possible cause for this is the long time it takes the procedure *getFullMVDs* to arrive at a full MVD when there are many attributes.

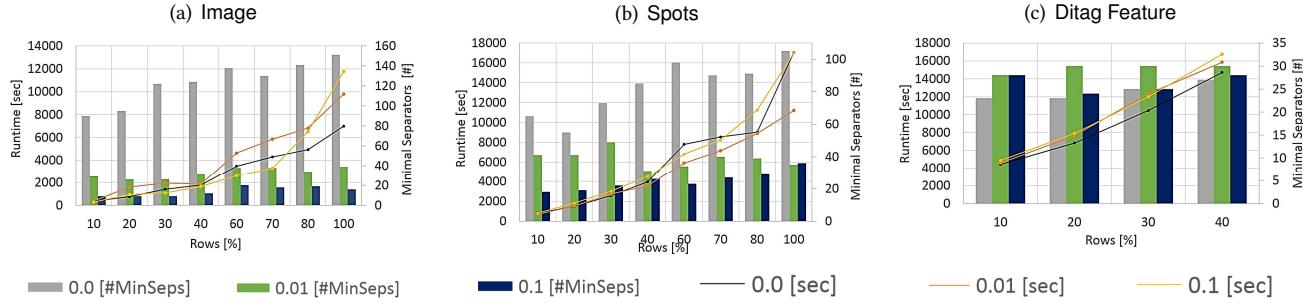
### 8.4 From minimal separators to full MVDs

We now experiment with the transition from minimal separators to full MVDs. We recall that an MVD  $\phi$  is full with regard to  $\varepsilon$  if  $R \models_{\varepsilon} \phi$  and for all MVDs  $\psi > \phi$  that strictly refine  $\phi$  then  $R \not\models_{\varepsilon} \psi$ .

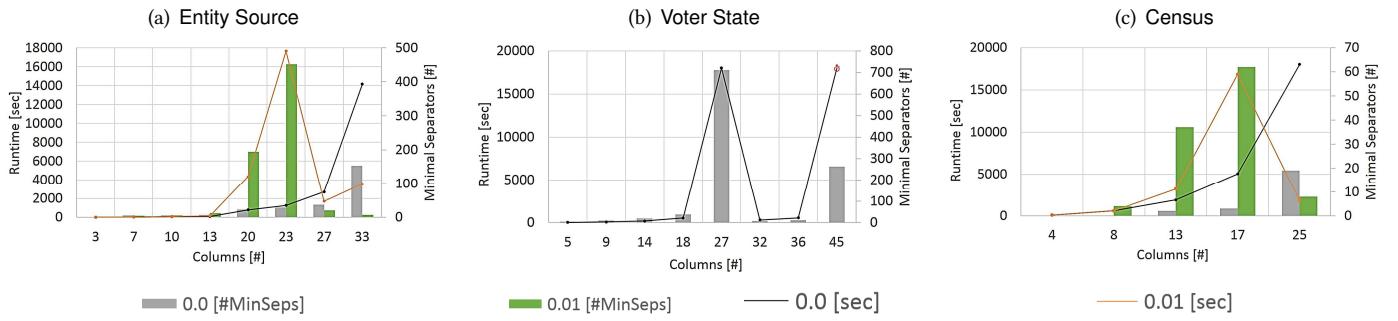
In this set of experiments we have, for every pair of attributes  $A, B \in \Omega$ , the set  $\text{MINSEP}_{\varepsilon, A, B}(R)$  of minimal  $AB$ -separators that hold in  $R$  w.r.t.  $\varepsilon$ , and we apply the algorithm for generating the set  $\text{FULLMVD}_{\varepsilon, A, B}$  by calling *getFullMVDs* (Fig. 6) with the pair  $(A, B)$ , and an unlimited number of MVDs to return (i.e.,  $K = \infty$ ).<sup>4</sup>

We conduct the experiment as follows. For every dataset we vary the threshold in the range  $[0, 0.5]$ , and for every threshold execute the procedure *getFullMVDsOpt* for a total of 30 minutes. The results are presented in Figure 13. When the threshold is  $\varepsilon = 0$  then the number of full MVDs is identical to the number of minimal separators as expected by

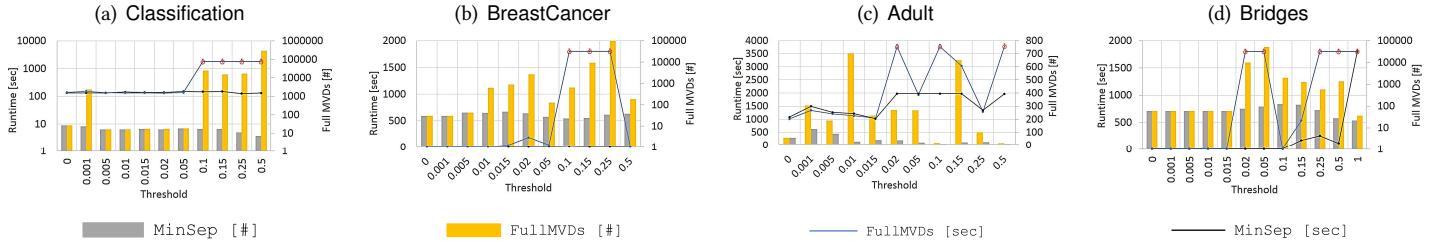
<sup>4</sup>We actually call the optimized version of this algorithm, *getFullMVDsOpt* described in the full version of this paper.



**Fig. 11** Row scalability experiments. Each chart contains the runtimes and number of minimal separators mined for three thresholds: 0.0, 0.01, and 0.1. All experiments presented completed within the timeout of 5 hours.



**Fig. 12** Column scalability experiments. Red stopwatch indicates that the algorithm timed out after five hours.



**Fig. 13** Full MVDs Experiments. Red stopwatch indicates that the algorithm stopped after 30 minutes.

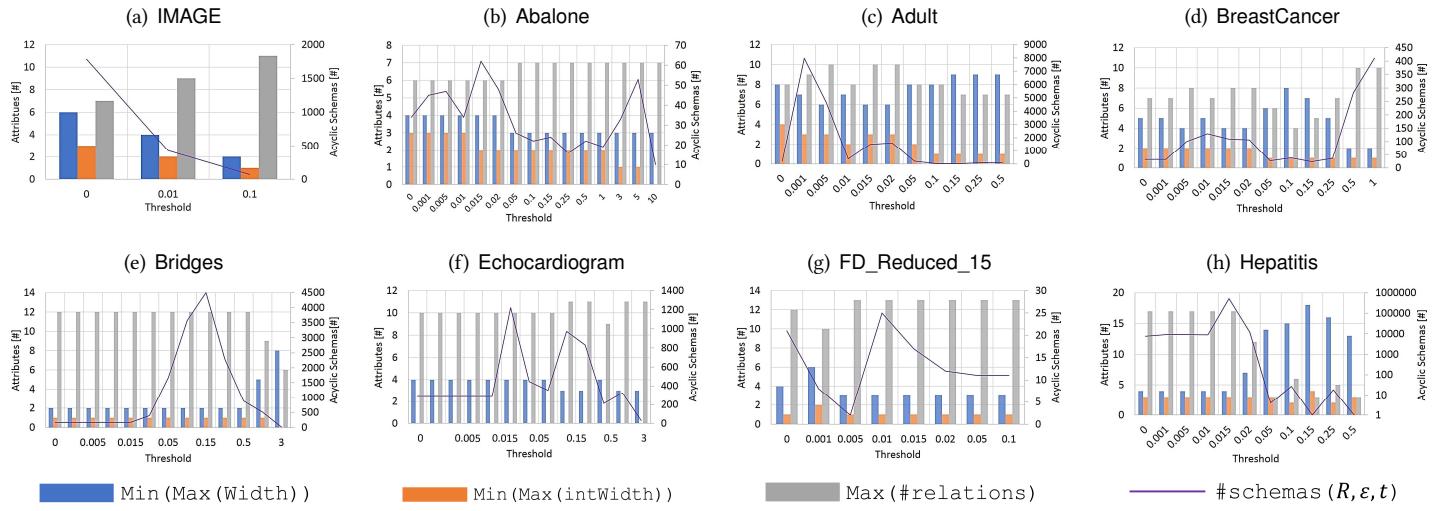
Lemma 5.4. In practice, when the threshold is 0, our algorithm for mining all minimal separators also discovers all full MVDs. As the threshold increases so does the difference between the number of minimal separators and the number of full MVDs. Overall, Algorithm getFullMVDsOpt for generating full MVDs is capable of reaching a rate of about 55 full MVDs per second (see Figures 13(a), 13(b), and 13(d)).

## 8.5 Enumeration Experiments

We present results of enumeration experiments carried out on a variety of datasets. Let  $S = \{\Omega_1, \dots, \Omega_m\}$  be an acyclic schema generated by our algorithm. We aggregate the following measurements which, together, provide a good indication to the extent in which the relation is decomposed by  $S$ .

- (1) The number of relations in  $S$ , denoted  $\#relations(S)$ . A larger number indicates a larger degree of decomposition.
- (2) The largest number of attributes in any relation of  $S$ . Formally,  $width(S) \stackrel{\text{def}}{=} \max_{i \in [1, m]} |\Omega_i|$ . A smaller number is indicative of a larger degree of decomposition.
- (3) We denote by  $\text{intWidth}(S)$  the largest number of attributes in any intersection  $\Omega_i \cap \Omega_j$  of relational schemas in  $S$ . Formally,  $\text{intWidth}(S) \stackrel{\text{def}}{=} \max_{i, j \in [1, m]} |\Omega_i \cap \Omega_j|$ . The number of attributes in the intersection plays a significant role in the efficiency of any join operation, giving an advantage to a small  $\text{intWidth}(S)$ .

Let  $\text{schemes}(R, \varepsilon, t)$  denote the set of schemes generated for relation  $R$ , with error threshold  $\varepsilon$ , within  $t$  seconds. For every



**Fig. 14** Enumeration Experiments

schema  $S \in \text{schemes}(R, \varepsilon, t)$  we calculate the three measurements  $\#relations(S)$ ,  $width(S)$ , and  $intWidth(S)$ , and plot the best value of any schema in  $\text{schemes}(R, \varepsilon, t)$ . In our experiments we varied the threshold  $\varepsilon$  in the range  $[0.0, 3]$ , and executed the enumeration algorithm for a duration of  $t = 30$  minutes. We also plotted the number of schemes  $\#schemas(R, \varepsilon, t)$  generated during  $t = 30$  minutes. The results are summarized in Figure 14.

In some cases we see a clear advantage to a higher threshold. For example, in the Image dataset (Fig. 13(a)), the maximum width and intersection width go down from 6 to 2, and 3 to 1 respectively, while the number of relations goes up from 7 to 11. We see the same trend in the Abalone dataset (Fig. 13(b)) where both the width and intersection width decrease with threshold. When  $\varepsilon = 10$ , the algorithm generates a schema whose largest intersection (i.e.,  $intWidth$ ) is 0, which means that the decomposition is comprised of independent relations with no common attribute. In practice, this translates to a tradeoff between the extent of decomposition and accuracy, since we've seen in Section 8.2 that thresholds over 0.5 generally lead to a significant loss in accuracy.

In terms of the number of acyclic schemes generated, we see a large variety among the datasets. In the BreastCancer dataset (Fig. 13(d)) there are 32 acyclic schemas for a threshold of 0.0, while this number goes up to 411 when increasing the threshold to 1.0. This is caused by the sharp increase in the number of MVDs that hold in the dataset when increasing the threshold. In many cases, however, the number of schemes generated actually goes down with the threshold (e.g., Figures 14(a)-14(c), 14(e), 14(f), and 14(h)). The reason is that higher thresholds allow for smaller sets of attributes to hold as minimal separators, bringing down their cardinality.

## 9 CONCLUSIONS

We present Maimon, the first system for the discovery of approximate acyclic schemes and approximate MVDs from data. To define “approximate”, we used concepts from information theory, where each MVD or acyclic schema is defined by an expression over entropic terms; when the expression is 0, then the MVD or acyclic schema holds exactly. We then presented the two main algorithms in Maimon, mining all full  $\varepsilon$ -MVDs with minimal separators, and discovering acyclic schemes from a set of  $\varepsilon$ -MVDs. Both algorithms improve over prior work in the literature. We conducted an experimental evaluation of Maimon on over 20 real-world data sets.

Our approach of using information theory to define approximate data dependencies differs from the previous definitions that rely mostly on counting the number of offending tuples. On one hand, our definitions provide us with more powerful mathematical tools, on the other hand the connection to the actual data quality is less intuitive. We leave it up to future work to explore the connection between information theory and data quality.

Depending on the dataset, Maimon generates hundreds and even thousands of acyclic  $\varepsilon$ -schemas in as little as 30 minutes. As part of future work we intend to investigate acyclic schema generation in *ranked order*. The categories to rank on may be the extent of decomposition (e.g., width of the schema), or other measures indicative of how well the schema meets the requirements of the application.

## REFERENCES

- [1] Ziawasch Abedjan, Lukasz Golab, Felix Naumann, and Thorsten Papenbrock. Data profiling. *Synthesis Lectures on Data Management*, 10(4):1–154, 2018.
- [2] Catriel Beeri. On the membership problem for functional and multivalued dependencies in relational databases. *ACM Trans. Database Syst.*, 5(3):241–259, September 1980.
- [3] Catriel Beeri and Philip A. Bernstein. Computational problems related to the design of normal form relational schemas. *ACM Trans. Database Syst.*, 4(1):30–59, March 1979.
- [4] Catriel Beeri, Ronald Fagin, and John H. Howard. A complete axiomatization for functional and multivalued dependencies in database relations. In *Proceedings of the 1977 ACM SIGMOD International Conference on Management of Data, Toronto, Canada, August 3–5, 1977*, pages 47–61, 1977.
- [5] Catriel Beeri, Ronald Fagin, David Maier, Alberto O. Mendelzon, Jeffrey D. Ullman, and Mihalis Yannakakis. Properties of acyclic database schemes. In *Proceedings of the 13th Annual ACM Symposium on Theory of Computing, May 11–13, 1981, Milwaukee, Wisconsin, USA*, pages 355–362, 1981.
- [6] Catriel Beeri, Ronald Fagin, David Maier, and Mihalis Yannakakis. On the desirability of acyclic database schemes. *J. ACM*, 30(3):479–513, July 1983.
- [7] Philip A. Bernstein. Synthesizing third normal form relations from functional dependencies. *ACM Trans. Database Syst.*, 1(4):277–298, 1976.
- [8] Tobias BleiFuß, Susanne Bülow, Johannes Frohnhofer, Julian Risch, Georg Wiese, Sebastian Kruse, Thorsten Papenbrock, and Felix Naumann. Approximate discovery of functional dependencies for large datasets. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24–28, 2016*, pages 1803–1812, 2016.
- [9] Nofar Carmeli, Batya Kenig, and Benny Kimelfeld. Efficiently enumerating minimal triangulations. In *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2017, Chicago, IL, USA, May 14–19, 2017*, pages 273–287, 2017.
- [10] E. F. Codd. Further normalization of the data base relational model. *IBM Research Report, San Jose, California*, RJ909, 1971.
- [11] Sara Cohen, Benny Kimelfeld, and Yehoshua Sagiv. Generating all maximal induced subgraphs for hereditary and connected-hereditary graph properties. *J. Comput. Syst. Sci.*, 74(7):1147–1159, 2008.
- [12] Tim Draeger. Multivalued dependency discovery, 2016. Master’s Thesis, Hasso-Plattner-Institute, Potsdam.
- [13] Ronald Fagin. Multivalued dependencies and a new normal form for relational databases. *ACM Trans. Database Syst.*, 2(3):262–278, September 1977.
- [14] Ronald Fagin. Horn clauses and database dependencies. *J. ACM*, 29(4):952–985, 1982.
- [15] Ronald Fagin, Alberto O. Mendelzon, and Jeffrey D. Ullman. A simplified universal relation assumption and its properties. *ACM Trans. Database Syst.*, 7(3):343–360, 1982.
- [16] Michael L. Fredman and Leonid Khachiyan. On the complexity of dualization of monotone disjunctive normal forms. *Journal of Algorithms*, 21(3):618 – 628, 1996.
- [17] Dan Geiger and Judea Pearl. Logical and algorithmic properties of conditional independence and graphical models. *The Annals of Statistics*, 21(4):2001–2021, 1993.
- [18] Nathan Goodman and Y. C. Tay. A characterization of multivalued dependencies equivalent to a join dependency. *Inf. Process. Lett.*, 18(5):261–266, 1984.
- [19] Dirk Van Gucht. Interaction-free multivalued dependency sets. *Theor. Comput. Sci.*, 62(1-2):221–233, 1988.
- [20] Dimitrios Gunopulos, Roni Khardon, Heikki Mannila, Sanjeev Saluja, Hannu Toivonen, and Ram Sewak Sharma. Discovering all most specific sentences. *ACM Trans. Database Syst.*, 28(2):140–174, 2003.
- [21] Ykä Huhtala, Juha Kärkkäinen, Pasi Porkka, and Hannu Toivonen. TANE: an efficient algorithm for discovering functional and approximate dependencies. *Comput. J.*, 42(2):100–111, 1999.
- [22] David S. Johnson, Christos H. Papadimitriou, and Mihalis Yannakakis. On generating all maximal independent sets. *Inf. Process. Lett.*, 27(3):119–123, 1988.
- [23] Leonid Khachiyan, Endre Boros, Khaled Elbassioni, and Vladimir Gurvich. An efficient implementation of a quasi-polynomial algorithm for generating hypergraph transversals and its application in joint generation. *Discrete Applied Mathematics*, 154(16):2350 – 2372, 2006. Discrete Algorithms and Optimization, in Honor of Professor Toshihide Ibaraki at His Retirement from Kyoto University.
- [24] Mahmoud Abo Khamis, Hung Q. Ngo, XuanLong Nguyen, Dan Olteanu, and Maximilian Schleich. AC/DC: in-database learning thunderstruck. In *Proceedings of the Second Workshop on Data Management for End-To-End Machine Learning, DEEM@SIGMOD 2018, Houston, TX, USA, June 15, 2018*, pages 8:1–8:10, 2018.
- [25] Mahmoud Abo Khamis, Hung Q. Ngo, and Atri Rudra. FAQ: questions asked frequently. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2016, San Francisco, CA, USA, June 20 - July 01, 2016*, pages 13–28, 2016.
- [26] Jyrki Kivinen and Heikki Mannila. Approximate inference of functional dependencies from relations. *Theor. Comput. Sci.*, 149(1):129–149, 1995.
- [27] Sebastian Kruse and Felix Naumann. Efficient discovery of approximate dependencies. *PVLDB*, 11(7):759–772, 2018.
- [28] V. S. Lakshmanan. Split-freedom and mvd-intersection: A new characterization of multivalued dependencies having conflict-free covers. *Theor. Comput. Sci.*, 62(1-2):105–122, 1988.
- [29] Tony T. Lee. An information-theoretic analysis of relational databases - part I: data dependencies and information metric. *IEEE Trans. Software Eng.*, 13(10):1049–1061, 1987.
- [30] Tony T. Lee. An information-theoretic analysis of relational databases - part II: information structures of database schemas. *IEEE Trans. Software Eng.*, 13(10):1061–1072, 1987.
- [31] Mark Levene and George Loizou. Why is the snowflake schema a good data warehouse design? *Inf. Syst.*, 28(3):225–240, 2003.
- [32] Y. Edmund Lien. Hierarchical schemata for relational databases. *ACM Trans. Database Syst.*, 6(1):48–69, March 1981.
- [33] Jixue Liu, Jiuyong Li, Chengfei Liu, and Yongfeng Chen. Discover dependencies from data - A review. *IEEE Trans. Knowl. Data Eng.*, 24(2):251–264, 2012.
- [34] Thorsten Papenbrock, Tanja Bergmann, Moritz Finke, Jakob Zwiener, and Felix Naumann. Data profiling with metanome. *Proc. VLDB Endow.*, 8(12):1860–1863, August 2015.
- [35] Thorsten Papenbrock and Felix Naumann. A hybrid approach to functional dependency discovery. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 821–833, 2016.
- [36] Iztok Savnik and Peter A. Flach. Discovery of multivalued dependencies from relations. *Intell. Data Anal.*, 4(3-4):195–211, 2000.
- [37] Maximilian Schleich, Dan Olteanu, and Radu Ciucanu. Learning linear regression models over factorized joins. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 3–18, 2016.

- [38] Maximilian Schleich, Dan Olteanu, Mahmoud Abo Khamis, Hung Q. Ngo, and XuanLong Nguyen. A layered aggregate engine for analytics workloads. In *Proceedings of the 2019 International Conference on Management of Data, SIGMOD Conference 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019.*, pages 1642–1659, 2019.
- [39] Catharine M. Wyss, Chris Giannella, and Edward L. Robertson. Fastfds: A heuristic-driven, depth-first algorithm for mining functional dependencies from relation instances - extended abstract. In *Data Warehousing and Knowledge Discovery, Third International Conference, DaWaK 2001, Munich, Germany, September 5-7, 2001, Proceedings*, pages 101–110, 2001.
- [40] Mihalis Yannakakis. Algorithms for acyclic database schemes. In *Proceedings of the Seventh International Conference on Very Large Data Bases - Volume 7*, VLDB '81, pages 82–94. VLDB Endowment, 1981.
- [41] h2 main memory database. <https://www.h2database.com/html/main.html>.