

An Efficient Implementation of a Quasi-Polynomial Algorithm for Generating Hypergraph Transversals and its Application in Joint Generation*

E. Boros[†] K. Elbassioni[†] V. Gurvich[†] L. Khachiyan[‡]

Abstract

Given a finite set V , and a hypergraph $\mathcal{H} \subseteq 2^V$, the hypergraph transversal problem calls for enumerating all minimal hitting sets (transversals) for \mathcal{H} . This problem plays an important role in practical applications as many other problems were shown to be polynomially equivalent to it. Fredman and Khachiyan (1996) gave an incremental quasi-polynomial time algorithm for solving the hypergraph transversal problem. In this paper, we present an efficient implementation of this algorithm. While we show that our implementation achieves the same theoretical worst case bound, practical experience with this implementation shows that it can be substantially faster. We also show that a slight modification of the original algorithm can be used to give a stronger bound on the running time.

More generally, we consider a monotone property π over a bounded n -dimensional integral box. As an important application of the above hypergraph transversal problem, pioneered by Bioch and Ibaraki (1995), we consider the problems of incrementally generating simultaneously both families of all minimal subsets satisfying π and all maximal subsets not satisfying π , for properties given by a polynomial-time satisfiability oracle. Problems of this type arise in many practical applications. It is known that the above joint generation problem can be solved in incremental quasi-polynomial time performing a polynomial-time reduction to a generalization of the hypergraph transversal problem on integer boxes. In this paper we present an efficient implementation of this procedure, and present experimental results to evaluate our implementation for a number of interesting monotone properties π .

1 Introduction

Let V be a finite set of cardinality $|V| = n$. For a hypergraph $\mathcal{H} \subseteq 2^V$, let us denote by $\mathcal{I}(\mathcal{H})$ the family of its *maximal independent* sets, i.e. maximal subsets of V not containing any hyperedge of \mathcal{H} . The complement of a maximal independent subset is called a *minimal transversal* of \mathcal{H} (i.e. minimal subset of V intersecting all hyperedges of \mathcal{H}). The collection \mathcal{H}^d of minimal transversals is also called the *dual* or *transversal* hypergraph for \mathcal{H} . The *hypergraph transversal problem* is the problem of generating all transversals of a given hypergraph. This problem has important applications in combinatorics [24], artificial intelligence [13], game theory [17, 18], reliability theory [11], database theory [9, 13, 16], integer programming [5], learning theory [1], and data mining [2, 7, 9].

The *theoretically* best known algorithm for solving the hypergraph transversal problem is due to Fredman and Khachiyan [14] and works by performing $|\mathcal{H}^d| + 1$ calls to the following problem, known as *hypergraph dualization*:

*This research was supported by the National Science Foundation (Grant IIS-0118635). The third author is also grateful for the partial support by DIMACS, the National Science Foundation's Center for Discrete Mathematics and Theoretical Computer Science.

[†]RUTCOR, Rutgers University, 640 Bartholomew Road, Piscataway NJ 08854-8003; {boros,elbassio,gurvich}@rutcor.rutgers.edu

[‡]Department of Computer Science, Rutgers University, 110 Frelinghuysen Road, Piscataway NJ 08854-8003; leonid@cs.rutgers.edu

DUAL(\mathcal{H}, \mathcal{X}): Given a complete list of all hyperedges of \mathcal{H} , and a set of minimal transversals $\mathcal{X} \subseteq \mathcal{H}^d$, either prove that $\mathcal{X} = \mathcal{H}^d$, or find a new transversal $X \in \mathcal{H}^d \setminus \mathcal{X}$.

Two recursive algorithms were proposed in [14] to solve the hypergraph dualization problem. These algorithms have incremental quasi-polynomial time complexities of $\text{poly}(n) + m^{O(\log^2 m)}$ and $\text{poly}(n) + m^{o(\log m)}$ respectively, where $m = |\mathcal{H}| + |\mathcal{X}|$. Even though the second algorithm is theoretically more efficient, the first algorithm is much simpler in terms of its implementation overhead, making it more attractive for practical applications. In fact, as we have found out experimentally, in many cases the most critical parts of the dualization procedure, in terms of execution time, are operations performed in each recursive call, rather than the total number of recursive calls. With respect to this measure, the first algorithm is more efficient due to its simplicity. For that reason, we present in this paper an implementation of the first algorithm in [14], which is efficient with respect to the time per recursive call. We further show that this efficiency in implementation does not come at the cost of increasing the worst-case running time substantially.

Rather than considering the hypergraph dualization problem, we shall consider, in fact, the more general problem of dualization on boxes introduced in [5]. In this latter problem, we are given an integral box $\mathcal{C} = \mathcal{C}_1 \times \dots \times \mathcal{C}_n$, where \mathcal{C}_i is a finite set of consecutive integers, and a subset $\mathcal{A} \subseteq \mathcal{C}$. Denote by $\mathcal{A}^+ = \{x \in \mathcal{C} \mid x \geq a, \text{ for some } a \in \mathcal{A}\}$ and $\mathcal{A}^- = \{x \in \mathcal{C} \mid x \leq a, \text{ for some } a \in \mathcal{A}\}$, the ideal and filter generated by \mathcal{A} . Any element in $\mathcal{C} \setminus \mathcal{A}^+$ is called *independent of \mathcal{A}* , and we let $\mathcal{I}(\mathcal{A})$ denote the set of all maximal independent elements for \mathcal{A} . Call a family of vectors \mathcal{A} *Sperner* if \mathcal{A} is an antichain, i.e. if no two elements are comparable in \mathcal{A} . If \mathcal{C} is the Boolean cube $2^{[n]}$, we get the well-known definitions of a hypergraph \mathcal{A} and its family of maximal independent sets $\mathcal{I}(\mathcal{A})$. Given $\mathcal{A} \subseteq \mathcal{C}$ and a subset $\mathcal{B} \subseteq \mathcal{I}(\mathcal{A})$ of maximal independent elements of \mathcal{A} , problem **DUAL**($\mathcal{C}, \mathcal{A}, \mathcal{B}$) calls for generating a new element $x \in \mathcal{I}(\mathcal{A}) \setminus \mathcal{B}$, or proving that there is no such element:

DUAL($\mathcal{C}, \mathcal{A}, \mathcal{B}$): Given a family of vectors $\mathcal{A} \subseteq \mathcal{C}$, and a subset $\mathcal{B} \subseteq \mathcal{I}(\mathcal{A})$ of its maximal independent vectors, either find a new maximal independent vector $x \in \mathcal{I}(\mathcal{A}) \setminus \mathcal{B}$, or prove that no such vector exists, i.e., $\mathcal{B} = \mathcal{I}(\mathcal{A})$.

By performing $|\mathcal{I}(\mathcal{A})| + 1$ calls to problem **DUAL**($\mathcal{C}, \mathcal{A}, \mathcal{B}$), we can solve the following problem

GEN($\mathcal{C}, \mathcal{I}(\mathcal{A})$): Given an integral box \mathcal{C} , and a subset of vectors $\mathcal{A} \subseteq \mathcal{C}$, generate all maximal independent elements of \mathcal{A} .

Problem **GEN**($\mathcal{C}, \mathcal{I}(\mathcal{A})$) has several interesting applications in integer programming and data mining, see [5, 6, 7] and the references therein. In particular, it is known that the incremental generation of minimal infrequent elements of a given database of quantitative attributes, and the incremental generation of maximal boxes not containing any point from a given point set, both reduce in polynomial time to problem **GEN**(\mathcal{C}, \mathcal{A}), see [6, 7]. These problems are motivated, respectively, by the generation of association rules in databases with quantitative attributes [25], and by the discovery of missing associations or "holes" in data mining applications (see [21, 22]).

Extensions of the above mentioned hypergraph transversal algorithms to solve problem **DUAL**($\mathcal{C}, \mathcal{A}, \mathcal{B}$) were given in [5]. In this paper, we give an implementation of the first dualization algorithm in [5], which achieves efficiency in two directions:

- Re-use of the recursion tree: dualization-based techniques generate all maximal independent elements of a given subset $\mathcal{A} \subseteq \mathcal{C}$ by usually performing $|\mathcal{I}(\mathcal{A})| + 1$ calls to problem **DUAL**($\mathcal{C}, \mathcal{A}, \mathcal{B}$), thus building a new recursion tree for each call. However, as it will be illustrated, it is more efficient to use the same recursion tree to generate all the elements of $\mathcal{I}(\mathcal{A})$, since the recursion trees required to generate many elements may be nearly identical.
- Efficient implementation at each recursion tree node: Straight forward implementation of the algorithm in [5] requires $O(n|\mathcal{A}| + n|\mathcal{B}|)$ time per recursive call. However, this can be improved to $O(n|\mathcal{A}| + |\mathcal{B}| + n \log(|\mathcal{B}|))$ by maintaining a binary search tree on the elements of \mathcal{B} , and using randomization. Since in many natural examples $|\mathcal{B}|$ is usually much larger than $|\mathcal{A}|$, this gives a significant

improvement. Moreover, the expected memory required by the randomized version is only $O(nm)$ whereas the deterministic version may require upto $O(nm^2)$ space, where $m = |\mathcal{A}| + |\mathcal{B}|$.

Several heuristics are also used to improve the running time. For instance, we use random sampling to find the branching variable and its value, required to divide the problem at the current recursion node. We also estimate the numbers of elements of \mathcal{A} and \mathcal{B} that are active at the current node, and only actually compute these active elements when their numbers drop by a certain factor. As our experiments indicate, such heuristics can be very effective in practically improving the running time of the algorithm.

Let $\pi : \mathcal{C} \mapsto \{0, 1\}$ be a *monotone* property defined over the elements of an integral box $\mathcal{C} = \mathcal{C}_1 \times \dots \times \mathcal{C}_n$: if $x \in \mathcal{C}$ satisfies property π , i.e. $\pi(x) = 1$, then any $y \in \mathcal{C}$ such that $y \geq x$ also satisfies π . We assume that π is described by a polynomial satisfiability oracle \mathcal{O}_π , i.e. an algorithm that can decide whether a given vector $x \in \mathcal{C}$ satisfies π , in time polynomial in n and the size $|\pi|$ of the input description of π . Denote respectively by \mathcal{F}_π and \mathcal{G}_π the families of minimal elements satisfying property π , and maximal elements not satisfying property π . Then it is clear that $\mathcal{G}_\pi = \mathcal{I}(\mathcal{F}_\pi)$ for any monotone property π . Given a monotone property π , we consider the problem of jointly generating the families \mathcal{F}_π and \mathcal{G}_π :

GEN($\mathcal{C}, \mathcal{F}_\pi, \mathcal{G}_\pi, \mathcal{X}, \mathcal{Y}$): *Given a monotone property π , represented by a satisfiability oracle \mathcal{O}_π , and two explicitly listed vector families $\mathcal{X} \subseteq \mathcal{F}_\pi \subseteq \mathcal{C}$ and $\mathcal{Y} \subseteq \mathcal{G}_\pi \subseteq \mathcal{C}$, either find a new element in $(\mathcal{F}_\pi \setminus \mathcal{X}) \cup (\mathcal{G}_\pi \setminus \mathcal{Y})$, or prove that these families are complete: $(\mathcal{X}, \mathcal{Y}) = (\mathcal{F}_\pi, \mathcal{G}_\pi)$.*

It is clear that for a given monotone property π , described by a satisfiability oracle \mathcal{O}_π , we can generate both \mathcal{F}_π and \mathcal{G}_π simultaneously by starting with $\mathcal{X} = \mathcal{Y} = \emptyset$ and solving problem GEN($\mathcal{C}, \mathcal{F}_\pi, \mathcal{G}_\pi, \mathcal{X}, \mathcal{Y}$) for a total of $|\mathcal{F}_\pi| + |\mathcal{G}_\pi| + 1$ times, incrementing in each iteration either \mathcal{X} or \mathcal{Y} by the newly found vector $x \in (\mathcal{F}_\pi \setminus \mathcal{X}) \cup (\mathcal{G}_\pi \setminus \mathcal{Y})$, according to the answer of the oracle \mathcal{O}_π , until we have $(\mathcal{X}, \mathcal{Y}) = (\mathcal{F}_\pi, \mathcal{G}_\pi)$.

In most practical applications, the requirement is to generate either the family \mathcal{F}_π or the family \mathcal{G}_π , i.e. we consider the separate generation problems:

GEN($\mathcal{C}, \mathcal{F}_\pi, \mathcal{X}$) (**GEN**($\mathcal{C}, \mathcal{G}_\pi, \mathcal{Y}$)): *Given a monotone property π and a subfamily $\mathcal{X} \subseteq \mathcal{F}_\pi \subseteq \mathcal{C}$ (respectively, $\mathcal{Y} \subseteq \mathcal{G}_\pi \subseteq \mathcal{C}$), either find a new minimal satisfying vector $x \in \mathcal{F}_\pi \setminus \mathcal{X}$ (respectively, maximal non-satisfying vector $x \in \mathcal{G}_\pi \setminus \mathcal{Y}$), or prove that the given partial list is complete: $\mathcal{X} = \mathcal{F}_\pi$ (respectively, $\mathcal{Y} = \mathcal{G}_\pi$).*

Problems GEN($\mathcal{C}, \mathcal{F}_\pi, \mathcal{X}$) and GEN($\mathcal{C}, \mathcal{G}_\pi, \mathcal{Y}$) arise in many practical applications and in a variety of fields, see e.g. [9, 8, 5, 7, 13, 20]. It is easy to see that Problem GEN($\mathcal{C}, \mathcal{G}_\pi, \mathcal{X}$) includes as a special case the dualization problem DUAL($\mathcal{C}, \mathcal{A}, \mathcal{B}$). Indeed, for any $x \in \mathcal{C}$, if we define the monotone property $\pi(x)$ to be satisfied by x if and only if $x \geq a$ for some $a \in \mathcal{A}$, then the families \mathcal{F}_π and \mathcal{G}_π can be identified with the families \mathcal{A} and $\mathcal{I}(\mathcal{A})$ respectively.

Even though the above two generation problems may be NP-hard in general (see e.g. [20]), it is known that the joint generation problem GEN($\mathcal{C}, \mathcal{F}_\pi, \mathcal{G}_\pi, \mathcal{X}, \mathcal{Y}$) can be solved in incremental $poly(n, \log \|\mathcal{C}\|_\infty) + m^{\text{poly} \log m}$ time (i.e., in incremental quasi-polynomial time) for any monotone property π described by a satisfiability oracle, where $\|\mathcal{C}\|_\infty = \sum_{i=1}^n |\mathcal{C}_i|$ and $m = |\mathcal{X}| + |\mathcal{Y}|$, see [4, 5, 14]. In particular, there is a polynomial-time reduction from the joint generation problem to the dualization problem on boxes, see [4, 5, 15]:

Proposition 1 *Problem GEN($\mathcal{C}, \mathcal{F}_\pi, \mathcal{G}_\pi, \mathcal{X}, \mathcal{Y}$) can be solved in time $\sum_{i=1}^n poly(\|\mathcal{C}\|_\infty, |\mathcal{A}|, |\mathcal{B}|) + T(\|\mathcal{O}_\pi\|) + T_{dual}$ for any monotone property π defined by a satisfiability oracle \mathcal{O}_π , where $T(\|\mathcal{O}_\pi\|)$ is the worst-case running time of the oracle on any $x \in \mathcal{C}$, and T_{dual} denotes the time required to solve problem DUAL($\mathcal{C}, \mathcal{A}, \mathcal{B}$).*

Proof. Let us first consider two subroutines which can be defined for any monotone property π . The first of these subroutines takes as input a vector $x \in \mathcal{F}_\pi^+$ and returns a minimal vector x^* in $\mathcal{F}_\pi^+ \cap \{x\}^-$. Such a vector $x^* = \min_{\mathcal{F}_\pi}(x)$ can, for instance, be computed by coordinate decent:

$$\begin{aligned}
x_1^* &\leftarrow \min\{y_1 \in \mathcal{C}_1 \mid (y_1, y_2, \dots, y_{n-1}, y_n) \in \mathcal{F}_\pi^+ \cap \{x\}^-\}, \\
x_2^* &\leftarrow \min\{y_2 \in \mathcal{C}_2 \mid (x_1^*, y_2, \dots, y_{n-1}, y_n) \in \mathcal{F}_\pi^+ \cap \{x\}^-\}, \\
&\vdots \\
x_n^* &\leftarrow \min\{y_n \in \mathcal{C}_n \mid (x_1^*, x_2^*, \dots, x_{n-1}^*, y_n) \in \mathcal{F}_\pi^+ \cap \{x\}^-\},
\end{aligned}$$

where $\min\{y_i \in \mathcal{P}_i \mid \dots\}$ returns an arbitrary minimal element with the specified property. Note that each of the n coordinate steps in the above procedures can be reduced via binary search to at most $\log(\|\mathcal{C}\|_\infty + 1)$ satisfiability oracle calls for the monotone property π .

The second subroutine is to compute, for a given vector $x \in \mathcal{I}(\mathcal{F}_\pi)^-$, a maximal vector $x^* \in \mathcal{I}(\mathcal{F}_\pi)^- \cap x^+$. Similarly, this problem can be done by coordinate decent. It is clear that these two subroutines can be executed in at most $O(\sum_{i=1}^n \log(\|\mathcal{C}\|_\infty + 1) \cdot T(\|\mathcal{O}_\pi\|))$ time.

Now the result stated in the proposition is obtained via the following algorithm.

Algorithm \mathcal{J}

Step 1. Check whether $\mathcal{B} \subseteq \mathcal{I}(\mathcal{A})$. Since $\mathcal{B} \subseteq \mathcal{I}(\mathcal{F}_\pi)$ and $\mathcal{A} \subseteq \mathcal{F}_\pi$, each vector $x \in \mathcal{B}$ is independent of \mathcal{A} and we only need to check the maximality of x for $\mathcal{I}(\mathcal{A})$. In other words, we have to check whether or not $y \geq \mathcal{A}^+$ for every immediate successor y of x . Since both \mathcal{A} and \mathcal{B} are explicitly given, this check can be done in $\text{poly}(n, |\mathcal{A}|, |\mathcal{B}|)$ comparisons. If there is an $x \in \mathcal{B} \setminus \mathcal{I}(\mathcal{A})$, then $x \notin \mathcal{F}_\pi^+$ because $x \in \mathcal{B} \subseteq \mathcal{I}(\mathcal{F}_\pi)$. This and the inclusion $\mathcal{A} \subseteq \mathcal{F}_\pi$ imply that $x \notin \mathcal{A}^+$. Since $x \notin \mathcal{I}(\mathcal{A})$, we can find an immediate successor $y \notin \mathcal{A}^+$ of x . By the maximality of x in $\mathcal{P} \setminus \mathcal{F}_\pi^+$, y belongs to \mathcal{F}_π^+ . Now letting $y^* = \min_{\mathcal{F}_\pi}(y)$ we conclude that $y^* \in \mathcal{F}_\pi \setminus \mathcal{A}$, i.e., y^* is a new minimal vector in \mathcal{F}_π .

Step 2 is similar to the previous step: we check whether $\mathcal{A} \subseteq \mathcal{I}^{-1}(\mathcal{B})$. If \mathcal{A} contains an element that is not minimal in $\mathcal{P} \setminus \mathcal{B}^-$, we can find a new vector in $\mathcal{I}(\mathcal{F}_\pi) \setminus \mathcal{B}$ and halt.

Step 3. Suppose that $\mathcal{B} \subseteq \mathcal{I}(\mathcal{A})$ and $\mathcal{A} \subseteq \mathcal{I}^{-1}(\mathcal{B})$. Then

$$(\mathcal{A}, \mathcal{B}) = (\mathcal{F}_\pi, \mathcal{I}(\mathcal{F}_\pi)) \iff \mathcal{B} = \mathcal{I}(\mathcal{A}).$$

To see this, assume that $\mathcal{B} = \mathcal{I}(\mathcal{A})$, and suppose on the contrary that there is an $x \in \mathcal{F}_\pi \setminus \mathcal{A}$. Since $x \notin \mathcal{A} = \mathcal{I}^{-1}(\mathcal{B})$ and $x \notin \mathcal{B}^- \subseteq \mathcal{I}(\mathcal{F}_\pi)^-$, there must exist a $y \in \mathcal{A} \subseteq \mathcal{F}_\pi$ such that $y \preceq x$. Hence we get two distinct elements $x, y \in \mathcal{F}_\pi$ such that $y \preceq x$, which contradicts the definition of \mathcal{F}_π . The existence of an $x \in \mathcal{I}(\mathcal{F}_\pi) \setminus \mathcal{B}$ leads to a similar contradiction.

To check the stopping criterion $\mathcal{B} = \mathcal{I}(\mathcal{A})$, we solve problem $\text{DUAL}(\mathcal{P}, \mathcal{A}, \mathcal{B})$. If $\mathcal{B} \neq \mathcal{I}(\mathcal{A})$, we obtain a new point $x \in \mathcal{I}(\mathcal{A}) \setminus \mathcal{B}$. By $\mathcal{F}_\pi^+ \cap \mathcal{I}(\mathcal{F}_\pi)^- = \emptyset$, either $x \in \mathcal{F}_\pi^+$, or $x \in \mathcal{I}(\mathcal{F}_\pi)^-$ and we can decide which of these two cases holds by calling the satisfiability oracle \mathcal{O}_π on x . In the first case, we conclude that $x^* = \min_{\mathcal{F}_\pi}(x)$ is a new vector in $\mathcal{F}_\pi \setminus \mathcal{A}$. In the second case, we can extend $\mathcal{I}(\mathcal{F}_\pi) \setminus \mathcal{B}$ by using the second subroutine to compute a maximal vector in $x^+ \cap \mathcal{I}(\mathcal{F}_\pi)^-$. \square

As mentioned above, the currently best known algorithm for dualization runs in time $\text{poly}(n) + m^{o(\log m)}$, see [5, 14]. Unfortunately, this joint generation may not be an efficient algorithm for solving either of $\text{GEN}(\mathcal{C}, \mathcal{F}_\pi, \mathcal{X})$ or $\text{GEN}(\mathcal{C}, \mathcal{G}_\pi, \mathcal{Y})$ separately for the simple reason that we do not control which of the families $\mathcal{F}_\pi \setminus \mathcal{X}$ and $\mathcal{G}_\pi \setminus \mathcal{Y}$ contains each new vector produced by the algorithm. Suppose we want to generate \mathcal{F}_π and the family \mathcal{G}_π is exponentially larger than \mathcal{F}_π . Then, if we are unlucky, we may get elements of \mathcal{F}_π with exponential delay, while getting large subfamilies of \mathcal{G}_π (which are not needed at all) in between. However, there are two reasons that we are interested in solving the joint generation problem efficiently. First, it is easy to see [16] that no satisfiability oracle based algorithm can generate \mathcal{F}_π in fewer than $|\mathcal{F}_\pi| + |\mathcal{G}_\pi|$ steps, in general:

Proposition 2 *Consider an arbitrary algorithm A , which generates the family \mathcal{F}_π by using only a satisfiability oracle \mathcal{O}_π for the monotone property π . Then, for the algorithm to generate the whole family \mathcal{F}_π , it must call the oracle at least $|\mathcal{F}_\pi| + |\mathcal{G}_\pi|$ times.*

Proof. Clearly, any monotone property π can be defined by its value on the boundary $\mathcal{F}_\pi \cup \mathcal{G}_\pi$. For any $y \in \mathcal{F}_\pi \cup \mathcal{G}_\pi$, let us thus define the boundary of the monotone property π_y as follows:

$$\pi_y(x) = \begin{cases} \pi(x) & \text{if } x \neq y \\ \bar{\pi}(x) & \text{if } x = y. \end{cases}$$

Then π_y is a monotone property different from π , and algorithm A must be able to distinguish the Sperner families described by π and π_y for every $y \in \mathcal{F}_\pi \cup \mathcal{G}_\pi$. \square

For a wide class of Sperner families (or equivalently, monotone properties), the so-called *uniformly dual-bounded* families, it was realized that the size of the dual family $\mathcal{I}(\mathcal{F}_\pi)$ is uniformly bounded by a (quasi-)polynomial in the size of \mathcal{F}_π and the oracle description:

$$|\mathcal{I}(\mathcal{X}) \cap \mathcal{I}(\mathcal{F}_\pi)| \leq (\text{quasi-})\text{poly}(|\mathcal{X}|, |\pi|) \quad (1)$$

for any non-empty subfamily $\mathcal{X} \subseteq \mathcal{F}_\pi$. An inequality of the form (1) implies that joint generation is an incrementally efficient way for generating the family \mathcal{F}_π (see Section 3.1 below and also [6] for several interesting examples):

Corollary 1 *Suppose \mathcal{F}_π is uniformly dual-bounded as in (1) and defined by a (quasi-) polynomial time membership oracle for \mathcal{F}_π^+ . Then problem $GEN(\mathcal{C}, \mathcal{F}_\pi, \mathcal{X})$ is (quasi-) polynomial-time reducible to at most (quasi-)poly($|\mathcal{X}|, \log(\|\mathcal{C}\|_\infty + 1), |\mathcal{O}_\pi|$) + 1 instances of problem $DUAL(\mathcal{C}, \mathcal{A}, \mathcal{B})$.*

Proof. Given a subset \mathcal{X} of \mathcal{F}_π , we repeatedly run Algorithm \mathcal{J} , starting with $\mathcal{A} = \mathcal{X}$ and $\mathcal{B} = \emptyset$, until it either produces a new element in $\mathcal{F}_\pi \setminus \mathcal{X}$ or proves that $\mathcal{X} = \mathcal{F}_\pi$ by generating the entire family $\mathcal{I}(\mathcal{F}_\pi)$. By Step 1, either $\mathcal{B} \subseteq \mathcal{I}(\mathcal{X})$ is maintained during the execution of the algorithm, or a new element $x \in \mathcal{F}_\pi \setminus \mathcal{X}$ can be found. Thus, as long as Algorithm \mathcal{J} outputs elements of $\mathcal{I}(\mathcal{F}_\pi)$, these elements also belong to $\mathcal{I}(\mathcal{X})$, and hence the total number of such elements does not exceed (quasi-)poly($|\mathcal{X}|, \log(\|\mathcal{C}\|_\infty + 1), \|\mathcal{O}_\pi\|$). \square

In Section 2.2, we present an efficient implementation of a quasi-polynomial dualization algorithm on boxes. Direct application of this implementation to solve the joint generation problem, as suggested by Proposition 1, may not be very efficient since each call to the dualization code requires the construction of a new recursion tree, wasting therefore information that might have been collected from previous calls. A much more efficient approach, which we implement in this paper, is to use the same recursion tree for generating all elements of the families \mathcal{F}_π and \mathcal{G}_π . The details of this method will be described in Section 3.2. In section 2.1 we introduce some basic terminology used throughout the paper, and briefly outline the quasi-polynomial algorithm [14] (or more precisely, its generalization to boxes [5]). Section 2.2 describes the data structure used in our implementation of this algorithm, and presents the algorithm. In Section 2.3, we show that the new version of the algorithm has, on the expected, the same quasi-polynomial bound on the running time as that of [5], and we also show how to get a slightly stronger bound on the running time. In Section 3.1, we give three examples of monotone properties, that will be used in our experimental study of the joint generation algorithm. Finally, Section 4 presents our experimental findings, and Section 5 provides some conclusions.

2 Hypergraph transversal

2.1 Terminology and outline of the algorithm

Throughout this section, we assume that we are given an integer box $\mathcal{C}^* = \mathcal{C}_1^* \times \dots \times \mathcal{C}_n^*$, where $\mathcal{C}_i^* = [l_i^* : u_i^*]$, and $l_i^* \leq u_i^*$, are integers, and a subset $\mathcal{A}^* \subseteq \mathcal{C}^*$ of vectors for which it is required to generate all maximal

independent elements. The algorithm of [5], considered in this paper, solves problem $\text{DUAL}(\mathcal{C}, \mathcal{A}, \mathcal{B})$, by decomposing it into a number of smaller subproblems and solving each of them recursively. The input to each such subproblem is a sub-box \mathcal{C} of the original box \mathcal{C}^* and two subsets $\mathcal{A} \subseteq \mathcal{A}^*$ and $\mathcal{B} \subseteq \mathcal{B}^*$ of integral vectors, where $\mathcal{B}^* \subseteq \mathcal{I}(\mathcal{A}^*)$ denotes the subfamily of maximal independent elements that the algorithm has generated so far. Note that, by definition, the following condition holds for the original problem and all subsequent subproblems:

$$a \not\leq b, \text{ for all } a \in \mathcal{A}, b \in \mathcal{B}. \quad (2)$$

Given an element $a \in \mathcal{A}$ ($b \in \mathcal{B}$), we say that a coordinate $i \in [n] \stackrel{\text{def}}{=} \{1, \dots, n\}$ is *essential* for a (respectively, b), in the box $\mathcal{C} = [l_1 : u_1] \times \dots \times [l_n : u_n]$, if $a_i > l_i$ (respectively, if $b_i < u_i$). Let us denote by $\text{Ess}(x)$ the set of essential coordinates of an element $x \in \mathcal{A} \cup \mathcal{B}$. Finally, given a sub-box $\mathcal{C} \subseteq \mathcal{C}^*$, and two subsets $\mathcal{A} \subseteq \mathcal{A}^*$ and $\mathcal{B} \subseteq \mathcal{B}^*$, we shall say that \mathcal{B} is *dual to* \mathcal{A} in \mathcal{C} if $\mathcal{A}^+ \cup \mathcal{B}^- \supseteq \mathcal{C}$.

A key lemma, on which the algorithm in [5] is based, is that either (i) there is an element $x \in \mathcal{A} \cup \mathcal{B}$ with at most $1/\epsilon$ essential coordinates, where $\epsilon \stackrel{\text{def}}{=} 1/(1 + \log m)$ and $m \stackrel{\text{def}}{=} |\mathcal{A}| + |\mathcal{B}|$, or (ii) one can easily find a new maximal independent element $z \in \mathcal{C}$, by picking each element z_i independently at random from $\{l_i, u_i\}$ for $i = 1, \dots, n$; see subroutine $\text{Random solution}(\cdot, \cdot, \cdot)$ in the next section. In case (i), one can decompose the problem into two strictly smaller subproblems as follows. Assume, without loss of generality, that $x \in \mathcal{A}$ has at most $1/\epsilon$ essential coordinates. Then, by (2), there is an $i \in [n]$ such that $|\{b \in \mathcal{B} : b_i < x_i\}| \geq \epsilon |\mathcal{B}|$. This allows us to decompose the original problem into two subproblems $\text{DUAL}(\mathcal{C}', \mathcal{A}, \mathcal{B}')$ and $\text{DUAL}(\mathcal{C}'', \mathcal{A}'', \mathcal{B})$, where $\mathcal{C}' = \mathcal{C}_1 \times \dots \times \mathcal{C}_{i-1} \times [x_i : u_i] \times \mathcal{C}_{i+1} \times \dots \times \mathcal{C}_n$, $\mathcal{B}' = \mathcal{B} \cap \mathcal{C}^+$, $\mathcal{C}'' = \mathcal{C}_1 \times \dots \times \mathcal{C}_{i-1} \times [l_i : x_i - 1] \times \mathcal{C}_{i+1} \times \dots \times \mathcal{C}_n$, and $\mathcal{A}'' = \mathcal{A} \cap \mathcal{C}^-$. This way, the algorithm is guaranteed to reduce the cardinality of one of the sets \mathcal{A} or \mathcal{B} by a factor of at least $1 - \epsilon$ at each recursive step. For efficiency reasons, we do two modifications to this basic approach. First, we use sampling to estimate the sizes of the sets $\mathcal{B}', \mathcal{A}''$ (see subroutine $\text{Est}(\cdot, \cdot)$ below). Second, once we have determined the new sub-boxes $\mathcal{C}', \mathcal{C}''$ above, we do not compute the *active* families \mathcal{B}' and \mathcal{A}'' at each recursion step (this is called the Cleanup step in the next section). Instead, we perform the cleanup step only when the number of vectors reduces by a certain factor f , say $1/2$, for two reasons: First, this improves the running time since the elimination of vectors is done less frequently. Second, the expected total memory required by all the nodes of the path from the root of the recursion tree to a leaf is at most $O(nm + m/(1 - f))$, which is linear in m for constant f .

2.2 The hypergraph transversal algorithm

We used the following data structures in our implementation:

- Two arrays of vectors, A and B containing the elements of \mathcal{A}^* and \mathcal{B}^* respectively.
- Two (dynamic) arrays of indices, $\text{index}(\mathcal{A})$ and $\text{index}(\mathcal{B})$, containing the indices of vectors from \mathcal{A}^* and \mathcal{B}^* (i.e. containing pointers to elements of the arrays A and B), that appear in the current subproblem. These arrays are used to enable sampling from the sets \mathcal{A} and \mathcal{B} , and also to keep track of which vectors are currently active, i.e. intersect the current box.
- A balanced binary search tree $\mathbf{T}(\mathcal{B}^*)$, built on the elements of \mathcal{B}^* using lexicographic ordering. Each node of the tree contains an index of an element in the array B . This way, checking whether a given vector $x \in \mathcal{C}$ belongs to \mathcal{B}^* or not, takes only $O(n \log |\mathcal{B}^*|)$ time.

In the sequel, we let $m = |\mathcal{A}| + |\mathcal{B}|$ and $\epsilon = 1/(1 + \log m)$. We assume further that operations of the form $\mathcal{A}'' \leftarrow \mathcal{A}$ and $\mathcal{B}' \leftarrow \mathcal{B}$ are actually performed on the index arrays $\text{index}(A)$, $\text{index}(B)$, so that they only take $O(m)$ rather than $O(nm)$ time. We used the following subroutines in our implementation:

Procedure GEN-DUAL($\mathcal{C}, \mathcal{A}, \mathcal{B}$):
Input: A box $\mathcal{C} = \mathcal{C}_1 \times \cdots \times \mathcal{C}_n$ and subsets $\mathcal{A} \subseteq \mathcal{A}^* \subseteq \mathcal{C}$, and $\mathcal{B} \subseteq \mathcal{I}(\mathcal{A}^*)$.
Output: A subset $\mathcal{N} \subseteq \mathcal{I}(\mathcal{A}^*) \setminus \mathcal{B}$.

1. $\mathcal{N} \leftarrow \emptyset$.
2. While $|\mathcal{A}||\mathcal{B}| \leq 1$
 - 2.1. $z \leftarrow \text{Exhaustive duality}(\mathcal{C}, \mathcal{A}, \mathcal{B})$.
 - 2.2. If $z = \text{FALSE}$ then return(\mathcal{N}).
 - 2.3. $\mathcal{B} \leftarrow \mathcal{B} \cup \{z\}$, $\mathcal{N} \leftarrow \mathcal{N} \cup \{z\}$.
end while
3. $z \leftarrow \text{Random Solution}(\mathcal{C}, \mathcal{A}^*, \mathcal{B})$.
4. While ($z \neq \text{FALSE}$) do
 - 4.1. $\mathcal{B} \leftarrow \mathcal{B} \cup \{z\}$, $\mathcal{N} \leftarrow \mathcal{N} \cup \{z\}$.
 - 4.2. $z \leftarrow \text{Random Solution}(\mathcal{C}, \mathcal{A}^*, \mathcal{B})$.
end while
5. $x^* \leftarrow \text{argmin}\{|\text{Ess}(y)| : y \in (\mathcal{A} \cap \mathcal{C}^-) \cup (\mathcal{B} \cap \mathcal{C}^+)\}$.
6. If $x^* \in \mathcal{A}$ then
 - 6.1. $i \leftarrow \text{argmax}\{\text{Est}(\{b \in \mathcal{B} : b_j < x_j^*\}, \mathcal{C}) : j \in \text{Ess}(x^*)\}$.
 - 6.2. $\mathcal{C}' = \mathcal{C}_1 \times \cdots \times \mathcal{C}_{i-1} \times [x_i^* : u_i] \times \mathcal{C}_{i+1} \times \cdots \times \mathcal{C}_n$.
 - 6.3. If $\text{Est}(\mathcal{B}, \mathcal{C}') \leq f * |\mathcal{B}|$ then
 - 6.3.1. $\mathcal{B}' \leftarrow \text{Cleanup}(\mathcal{B}, \mathcal{C}')$.
 - 6.4. else
 - 6.4.1. $\mathcal{B}' \leftarrow \mathcal{B}$.
 - 6.5. $\mathcal{N}' \leftarrow \text{GEN-DUAL}(\mathcal{C}', \mathcal{A}, \mathcal{B}')$.
 - 6.6. $\mathcal{N} \leftarrow \mathcal{N} \cup \mathcal{N}'$, $\mathcal{B} \leftarrow \mathcal{B} \cup \mathcal{B}'$.
 - 6.7. $\mathcal{C}'' = \mathcal{C}_1 \times \cdots \times \mathcal{C}_{i-1} \times [l_i : x_i^* - 1] \times \mathcal{C}_{i+1} \times \cdots \times \mathcal{C}_n$.
 - 6.8. If $\text{Est}(\mathcal{A}, \mathcal{C}'') \leq f * |\mathcal{A}|$ then
 - 6.8.1. $\mathcal{A}'' \leftarrow \text{Cleanup}(\mathcal{A}, \mathcal{C}'')$.
 - 6.9. else
 - 6.9.1. $\mathcal{A}'' \leftarrow \mathcal{A}$.
 - 6.10. $\mathcal{N}'' \leftarrow \text{GEN-DUAL}(\mathcal{C}'', \mathcal{A}'', \mathcal{B})$.
 - 6.11. $\mathcal{N} \leftarrow \mathcal{N} \cup \mathcal{N}''$, $\mathcal{B} \leftarrow \mathcal{B} \cup \mathcal{B}''$.
7. else
 - 7.1-7.11. Symmetric versions for Steps 6.1-6.11 above (details omitted).
end if
8. Return (\mathcal{N}).

Maximization $\max_{\mathcal{A}}(z)$. It takes as input a vector $z \notin \mathcal{A}^+$ and returns a maximal vector z^* in $(\mathcal{C}^* \cap \{z\}^+) \setminus \mathcal{A}^+$. This can be done in $O(n|\mathcal{A}|)$ by initializing $c(a) = |\{i \in [n] : a_i > z_i\}|$ for all $a \in \mathcal{A}$, and repeating, for $i = 1, \dots, n$, the following two steps: (i) $z_i^* \leftarrow \min(u_i^*, \min\{a_i - 1 : a \in \mathcal{A}, c(a) = 1 \text{ and } a_i > z_i\})$ (where we assume $\min(\emptyset) = \infty$); (ii) $c(a) \leftarrow c(a) - 1$ for each $a \in \mathcal{A}$ such that $z_i < a_i \leq z_i^*$.

Exhaustive duality($\mathcal{C}, \mathcal{A}, \mathcal{B}$). Assuming $|\mathcal{A}||\mathcal{B}| \leq 1$, check duality in $O(n(|\mathcal{A}^*| + \log |\mathcal{B}|))$ as follows: First, if $|\mathcal{A}| = |\mathcal{B}| = 1$ then find an $i \in [n]$ such that $a_i > b_i$, where $\mathcal{A} = \{a\}$ and $\mathcal{B} = \{b\}$. (Such a coordinate is guaranteed to exist by (2).) If there is a $j \neq i$ such that $b_j < u_j$ then return $\max_{\mathcal{A}^*}(u_1, \dots, u_{i-1}, b_i, u_{i+1}, \dots, u_n)$. If there is a $j \neq i$ such that $a_j > l_j$ then return $(u_1, \dots, u_{j-1}, a_j - 1, u_{j+1}, \dots, u_n)$. If $b_i < a_i - 1$ then return $(u_1, \dots, u_{i-1}, a_i - 1, u_{i+1}, \dots, u_n)$. Otherwise return *FALSE* (meaning that \mathcal{A} and \mathcal{B} are dual in \mathcal{C}). Second, if $|\mathcal{A}| = 0$ then let $z = \max_{\mathcal{A}^*}(u)$, and return either *FALSE* or z depending on whether $z \in \mathcal{B}^*$ or not (this check can be done in $O(n \log |\mathcal{B}^*|)$ using the search tree $\mathbf{T}(\mathcal{B}^*)$). Finally, if $|\mathcal{B}| = 0$ then return either *FALSE* or $z = \max_{\mathcal{A}^*}(l)$ depending on whether $l \in \mathcal{A}^+$ or not (this check requires $O(n|\mathcal{A}|)$ time).

Random solution($\mathcal{C}, \mathcal{A}^*, \mathcal{B}$). Repeat the following for $k = 1, \dots, t_1$ times, where t_1 is a constant (say 10): Find a random point $z^k \in \mathcal{C}$, by picking each coordinate z_i^k randomly from $\{l_i, u_i\}$, $i = 1, \dots, n$. Let

$(z^k)^* \leftarrow \max_{\mathcal{A}^*} (z^k)$. If $(z^k)^* \notin \mathcal{B}^*$ then return $(z^k)^*$. If $\{(z^1)^*, \dots, (z^{t_1})^*\} \subseteq \mathcal{B}^*$ then return *FALSE*. This step takes $O(n(|\mathcal{A}^*| + \log |\mathcal{B}^*|))$ time, and is used to check whether $\mathcal{A}^+ \cup \mathcal{B}^-$ covers a large portion of \mathcal{C} .

Count estimation. For a subset $\mathcal{X} \subseteq \mathcal{A}$ (or $\mathcal{X} \subseteq \mathcal{B}$), use sampling to estimate the number $\text{Est}(\mathcal{X}, \mathcal{C})$ of elements of $\mathcal{X} \subseteq \mathcal{A}$ (or $\mathcal{X} \subseteq \mathcal{B}$) that are active with respect to the current box \mathcal{C} , i.e. the elements of the set $\mathcal{X}' \stackrel{\text{def}}{=} \{a \in \mathcal{X} \mid a^+ \cap \mathcal{C} \neq \emptyset\}$ ($\mathcal{X}' \stackrel{\text{def}}{=} \{b \in \mathcal{X} \mid b^- \cap \mathcal{C} \neq \emptyset\}$). This can be done as follows. For $t_2 = O(\log(|\mathcal{A}| + |\mathcal{B}|)/\epsilon)$, pick elements $x^1, \dots, x^{t_2} \in \mathcal{A}$ at random, and let the random variable $Y = \frac{|\mathcal{A}|}{t_2} * |\{x^i \in \mathcal{X}' : i = 1, \dots, t_2\}|$. Repeat this step independently for a total of $t_3 = O(\log(|\mathcal{A}| + |\mathcal{B}|))$ times to obtain t_3 estimates Y^1, \dots, Y^{t_3} , and let $\text{Est}(\mathcal{X}, \mathcal{C}) = \min\{Y^1, \dots, Y^{t_3}\}$. This step requires $O(n \log^3 m)$ time.¹

Cleanup(\mathcal{A}, \mathcal{C}) (Cleanup(\mathcal{B}, \mathcal{C})). Set $\mathcal{A}' \leftarrow \{a \in \mathcal{A} \mid a^+ \cap \mathcal{C} \neq \emptyset\}$ (respectively, $\mathcal{B}' \leftarrow \{b \in \mathcal{B} \mid b^- \cap \mathcal{C} \neq \emptyset\}$), and return \mathcal{A}' (respectively, \mathcal{B}'). This step takes $O(n|\mathcal{A}|)$ (respectively, $O(n|\mathcal{B}|)$).

Now, we describe the implementation of procedure $\text{GEN-DUAL}(\mathcal{A}, \mathcal{B}, \mathcal{C})$ which is called initially using $\mathcal{C} \leftarrow \mathcal{C}^*$, $\mathcal{A} \leftarrow \mathcal{A}^*$ and $\mathcal{B} \leftarrow \emptyset$. At the return of this call, \mathcal{B} is extended by the elements in $\mathcal{I}(\mathcal{A}^*)$. Below we assume that $f \in (0, 1)$ is a constant.

2.3 Analysis of the expected running time

Let $C(v)$ be the expected number of recursive calls on a subproblem $\text{GEN-DUAL}(\mathcal{C}, \mathcal{A}, \mathcal{B})$ of volume $v \stackrel{\text{def}}{=} |\mathcal{A}||\mathcal{B}|$. Consider a particular recursive call of the algorithm and let \mathcal{A}, \mathcal{B} and \mathcal{C} be the current inputs to this call. Let x^* be the element with minimum number of essential coordinates found in Step 5, and assume without loss of generality that $x^* \in \mathcal{A}$. As mentioned before, we assume also that the factor f used in Steps 6.3 and 6.8 is $1/2$. For $i = 1, \dots, n$, let $\mathcal{B}_i \stackrel{\text{def}}{=} \{b \in \mathcal{B} : b_i < x_i^*\}$, and denote by $\overline{\mathcal{B}} = \mathcal{B} \cap \mathcal{C}^+$ and $\overline{\mathcal{B}}_i = \mathcal{B}_i \cap \mathcal{C}^+$ the subsets of \mathcal{B} and \mathcal{B}_i that are active with respect to the current box \mathcal{C} . In this section, we show that our implementation has, with high probability, almost the same quasi-polynomial bound on the running time as the algorithm of [5].

Lemma 1 *Suppose that $k \in \text{Ess}(x^*)$ satisfies $|\overline{\mathcal{B}}_k| \geq \epsilon |\overline{\mathcal{B}}|$. Let $i \in [n]$ be the coordinate obtained in Step 6.1 of the algorithm, and $v = |\mathcal{A}||\mathcal{B}|$. Then*

$$\Pr \left[|\overline{\mathcal{B}}_i| \geq \epsilon \frac{|\overline{\mathcal{B}}|}{4} \right] \geq 1 - \frac{1}{\epsilon v}. \quad (3)$$

Proof. For $j = 1, \dots, n$, let $Y_j \stackrel{\text{def}}{=} \text{Est}(\mathcal{B}_j, \mathcal{C})$. Then the random variable $X_j \stackrel{\text{def}}{=} t_2 Y_j / |\mathcal{B}|$ is Binomially distributed with parameters t_2 and $|\overline{\mathcal{B}}_j|/|\mathcal{B}|$, and thus by Chernoff Bound

$$\Pr[Y_j < \frac{\mathbb{E}[Y_j]}{2}] < e^{-\mathbb{E}[X_j]/8}, \quad \text{for } j = 1, \dots, n.$$

In particular, for $j = k$, we get $\Pr[Y_k < \epsilon \frac{|\overline{\mathcal{B}}|}{2}] < e^{-\mathbb{E}[X_k]/8}$ since $\mathbb{E}[Y_k] \geq \epsilon |\overline{\mathcal{B}}|$. Note that, since $\text{Est}(\mathcal{B}, \mathcal{C})$ is the minimum over t_3 independent trials, it follows by Markov Inequality that if $|\overline{\mathcal{B}}| < |\mathcal{B}|/4$, then $\Pr[\text{Est}(\mathcal{B}, \mathcal{C}) < |\mathcal{B}|/2] > 1 - 2^{-t_3}$, and the cleanup step will be performed with high probability. On the other hand, if $|\overline{\mathcal{B}}| \geq |\mathcal{B}|/4$ then $\mathbb{E}[X_k] = t_2 |\overline{\mathcal{B}}_k|/|\mathcal{B}| \geq \epsilon t_2/4$. Thus, it follows that $\Pr[Y_k < \epsilon \frac{|\overline{\mathcal{B}}|}{2}] < e^{-\epsilon t_2/32} + 2^{-t_3}$. Moreover, for any $j \in \text{Ess}(x^*)$ for which $|\overline{\mathcal{B}}_j|/|\mathcal{B}| < \epsilon/4$, we have $\Pr[Y_j \geq \epsilon \frac{|\overline{\mathcal{B}}|}{2}] < 2^{-t_3}$. Consequently,

$$\Pr \left[Y_k \geq \frac{\epsilon |\overline{\mathcal{B}}|}{2} \text{ and } Y_j < \frac{\epsilon |\overline{\mathcal{B}}|}{2} \text{ for all } j \in \text{Ess}(x^*) \text{ such that } \frac{|\overline{\mathcal{B}}_j|}{|\mathcal{B}|} < \frac{\epsilon}{4} \right] > 1 - 2^{-t_3} |\text{Ess}(x^*)| - e^{-\epsilon t_2/32} \geq 1 - \frac{1}{\epsilon v},$$

¹Note that these sample sizes were chosen to theoretically get a guarantee on the expected running time of the algorithm. However, as our experiments indicate, smaller (usually constant) sample sizes are enough to provide practically good performance.

where the last inequality follows by our selection of t_2 and t_3 . Since, in Step 6.1, we select the index $i \in [r]$ maximizing Y_i , we have $Y_i \geq Y_k$ and thus, with probability at least $1 - 1/(\epsilon v)$, we have $|\overline{\mathcal{B}}_i|/|\overline{\mathcal{B}}| \geq \epsilon/4$. \square

Lemma 2 *The expected number of recursive calls until a new maximal independent element is output, or procedure GEN-DUAL($\mathcal{C}, \mathcal{A}, \mathcal{B}$) terminates is $nm^{O(\log^2 m)}$.*

Proof. For a node N of the recursion tree, denote by $\overline{\mathcal{A}} = \overline{\mathcal{A}}(N)$, $\overline{\mathcal{B}} = \overline{\mathcal{B}}(N)$ the subsets of \mathcal{A} and \mathcal{B} intersecting the box specified by node N , and let $v(N) = |\overline{\mathcal{A}}(N)||\overline{\mathcal{B}}(N)|$. Now consider the node N at which the latest maximal independent element was generated. If $s \stackrel{\text{def}}{=} \sum_{a \in \overline{\mathcal{A}}(N)} (1/2)^{|\text{Ess}(a)|} + \sum_{b \in \overline{\mathcal{B}}(N)} (1/2)^{|\text{Ess}(b)|} < 1/2$, then the probability that the point $z \in \mathcal{C}$, picked randomly in Steps 3 or 4.2 of the procedure, belongs to $\overline{\mathcal{A}}(N)^+ \cup \overline{\mathcal{B}}(N)^-$ is at most $\sigma_1 \stackrel{\text{def}}{=} (1/2)^{t_1}$. Thus, in this case, with probability at least $1 - \sigma_1$, we find a new maximal independent element. Assume therefore that $s \geq 1/2$, let x^* be the element with $|\text{Ess}(x^*)| \leq 1/\epsilon$ found in Step 5, and assume without loss of generality that $x^* \in \mathcal{A}$. Then, by (2), there exists a coordinate $k \in \text{Ess}(x^*)$ such that $|\overline{\mathcal{B}}_k| \geq \epsilon|\overline{\mathcal{B}}|$. By Lemma 1, with probability at least $\sigma_2 \stackrel{\text{def}}{=} 1 - 1/(\epsilon v)$, we can reduce the volume of one of the subproblems, of the current problem, by a factor of at least $1 - \epsilon/4$. Thus for the expected number of recursive calls at node N , we get the following recurrence

$$C(v) \leq \frac{1}{\sigma_2} \left[1 + C(v-1) + C\left(\left(1 - \frac{\epsilon}{4}\right)v\right) \right], \quad (4)$$

where $v = v(N)$. This recurrence gives $C(v) \leq v^{O(\log^2 v)}$. Now consider the path $N_0 = N, N_1, \dots, N_r$ from node N to the root of the recursion tree N_r . Since a large number of new maximal independent elements may have been added at node N (and of course to all its ancestors in the tree), recurrence (4) may no longer hold at nodes N_1, \dots, N_r . However, since we count the number of recursive calls from the time of the last generation that happened at node N , each node N_i , that has N_{i-1} as a right child in the tree, does not contribute to this number. Furthermore, the number of recursive calls resulting from the right child of each node N_i , that has N_{i-1} as a left child, is at most $C(v(N_r))$. Since the number of such nodes does not exceed the depth of the tree, which is at most nm , the expected total number of recursive calls is at most $nmC(v(N_r))$ and the lemma follows. \square

We show further that, if $|\mathcal{B}| \gg |\mathcal{A}|$, i.e. if the output size is much bigger than the input size, then the number of recursive calls required for termination, after the last dual element is generated by GEN-DUAL($\mathcal{A}, \mathcal{B}, \mathcal{C}$), is $nm^{o(\log^2 m)}$.

Lemma 3 *Suppose that \mathcal{A} and \mathcal{B} are dual in \mathcal{C} , then the expected number of recursive calls until GEN-DUAL($\mathcal{C}, \mathcal{A}, \mathcal{B}$) terminates is $nm^{O(\delta \log m)}$, where $m = |\mathcal{A}| + |\mathcal{B}|$ and $\delta = \min\{\log \alpha, \frac{\log(\beta/\alpha)}{c(\alpha, \beta/\alpha)}, \frac{\log(\alpha/\beta)}{c(\beta, \alpha/\beta)}\} + 1$, $\alpha = |\mathcal{A}|$, $\beta = |\mathcal{B}|$, and $c = c(a, b)$ is the unique positive root of the equation*

$$2^c \left(a^{c/\log b} - 1 \right) = 1. \quad (5)$$

Proof. Let $r = \min\{|\text{Ess}(y)| : y \in \mathcal{A} \cup \mathcal{B}\}$, $p = \left(1 + \left(\frac{\beta}{\alpha} \right)^{\frac{1}{r-1}} \right)^{-1}$, and let $z \in \mathcal{C}$ be a random element obtained by picking each coordinate independently with $\Pr[z_i = l_i] = p$ and $\Pr[z_i = u_i] = 1 - p$. Then the probability that $z \in \mathcal{A}^+ \cup \mathcal{B}^-$ is at most $\sum_{a \in \mathcal{A}} (1-p)^{|\text{Ess}(a)|} + \sum_{b \in \mathcal{B}} p^{|\text{Ess}(b)|} \leq \alpha(1-p)^r + \beta p^r = \beta p^{r-1}$. Since \mathcal{A} and \mathcal{B} are dual in \mathcal{C} , it follows that $\beta p^{r-1} \geq 1$, and thus

$$r - 1 \leq \frac{\log \beta}{\log(1 + (\beta/\alpha)^{\frac{1}{r-1}})}. \quad (6)$$

The maximum value that r can achieve is when both sides of (6) are equal, i.e. r is bounded by the root r' of the equation $\beta^{1/(r'-1)} = 1 + (\beta/\alpha)^{1/(r'-1)}$. If $\alpha = \beta$, then $r' = \log \alpha + 1$. If $\beta > \alpha$, then letting

$(\beta/\alpha)^{1/(r'-1)} = 2^c$, we get

$$r' = 1 + \frac{\log(\beta/\alpha)}{c(\alpha, \beta/\alpha)}, \quad (7)$$

where $c(\cdot, \cdot)$ is as defined in (5). The case for $\alpha > \beta$ is similar and the lemma follows from (2) and Lemma 2. \square

Note that, if β is much larger than α , then the root r' in (7) is approximately

$$r' \sim 1 + \frac{\log(\beta/\alpha)}{\log(\log(\beta/\alpha)/\log \alpha)}$$

and thus the expected running time of procedure GEN-DUAL($\mathcal{A}, \mathcal{B}, \mathcal{C}$), from the time of last output till termination, is $nm^{o(\log^2 m)}$. In fact, one can use Lemma 2 together with the method of conditional expectations to obtain an incremental deterministic algorithm for solving problem GEN(\mathcal{C}, \mathcal{A}), whose delay between any two successive outputs is of the order given by Lemma 2.

3 Joint generation

3.1 Examples of Monotone Properties

We consider the following three monotone properties in this paper:

Monotone systems of linear inequalities. Let $A \in \mathbb{R}^{r \times n}$ be a given non-negative real matrix, $b \in \mathbb{R}^r$ be a given r -vector, $c \in \mathbb{R}_+^n$ be a given non-negative n -vector, and consider the system of linear inequalities:

$$Ax \geq b, \quad x \in \mathcal{C} = \{x \in \mathbb{Z}^n \mid 0 \leq x \leq c\}. \quad (8)$$

For $x \in \mathcal{C}$, let $\pi_1(x)$ be the property that x satisfies (8). Then the families \mathcal{F}_{π_1} and \mathcal{G}_{π_1} correspond respectively to the minimal feasible and maximal infeasible vectors for (8). It is known [5] that the family \mathcal{F}_{π_1} is (uniformly) dual bounded:

$$|\mathcal{I}(\mathcal{F}_{\pi_1})| \leq rn|\mathcal{F}_{\pi_1}|. \quad (9)$$

Minimal infrequent and maximal frequent sets in binary databases. Let $\mathcal{D} : R \times V \mapsto \{0, 1\}$ be a given $r \times n$ binary matrix representing a set R of transactions over a set of attributes V . To each subset of columns $X \subseteq V$, let us associate the subset $S(X) = S_{\mathcal{D}}(X) \subseteq R$ of all those rows $i \in R$ for which $\mathcal{D}(i, j) = 1$ in every column $j \in X$. The cardinality of $S(X)$ is called the *support* of X . Given an integer t , a column set $X \subseteq V$ is called *t -frequent* if $|S(X)| \geq t$ and otherwise, is said to be *t -infrequent*. For each set $X \in \mathcal{C} \stackrel{\text{def}}{=} 2^V$, let $\pi_2(X)$ be the property that X is t -infrequent. Then π_2 is a monotone property and the families \mathcal{F}_{π_2} and \mathcal{G}_{π_2} correspond respectively to minimal infrequent and maximal frequent sets for \mathcal{D} . It is known [9] that

$$|\mathcal{I}(\mathcal{F}_{\pi_2})| \leq (r - t + 1)|\mathcal{F}_{\pi_2}|. \quad (10)$$

Problems GEN($\mathcal{C}, \mathcal{F}_{\pi_2}, \mathcal{Y}$) and GEN($\mathcal{C}, \mathcal{G}_{\pi_2}, \mathcal{Y}$) appear in data mining applications, see e.g. [16].

Sparse boxes for multi-dimensional data. Let \mathcal{S} be a set of points in \mathbb{R}^n , and $t \leq |\mathcal{S}|$ be a given integer. A *maximal t -box* is a closed n -dimensional interval which contains at most t points of \mathcal{S} in its interior, and which is maximal with respect to this property (i.e., cannot be extended in any direction without strictly enclosing more points of \mathcal{S}). Define $\mathcal{C}_i = \{p_i \mid p \in \mathcal{S}\}$ for $i = 1, \dots, n$ and consider the family of boxes $\mathcal{B} = \{[a, b] \subseteq \mathbb{R}^n \mid a, b \in \mathcal{C}_1 \times \dots \times \mathcal{C}_n, a \leq b\}$. For $i = 1, \dots, n$, let $u_i = \max \mathcal{C}_i$, and let $\mathcal{C}_{i+n} \stackrel{\text{def}}{=} \{u_i - p \mid p \in \mathcal{C}_i\}$ be the chain ordered in the direction opposite to \mathcal{C}_i . Consider the $2n$ -dimensional box $\mathcal{C} = \mathcal{C}_1 \times \dots \times \mathcal{C}_n \times \mathcal{C}_{n+1} \times \dots \times \mathcal{C}_{2n}$ and let us represent every n -dimensional interval $[a, b] \in \mathcal{B}$ as

the $2n$ -dimensional vector $(a, u - b) \in \mathcal{C}$, where $u = (u_1, \dots, u_n)$. This gives a monotone injective mapping $\mathcal{B} \mapsto \mathcal{C}$ (not all elements of \mathcal{C} define a box, since $a_i > b_i$ is possible for $(a, u - b) \in \mathcal{C}$). Let us now define the monotone property π_3 to be satisfied by an $x \in \mathcal{C}$ if and only if x does not define a box, or the box defined by x contains at most t points of S in its interior. Then the sets \mathcal{F}_{π_3} and \mathcal{G}_{π_3} can be identified respectively with the set of maximal t -boxes (plus a polynomial number of non-boxes), and the set of minimal boxes of $x \in \mathcal{B} \subseteq \mathcal{C}$ which contain at least $k + 1$ points of S in their interior. It is known [7] that the family of maximal t -boxes is (uniformly) dual-bounded:

$$|\mathcal{I}(\mathcal{F}_{\pi_3})| \leq |\mathcal{S}| |\mathcal{F}_{\pi_3}|. \quad (11)$$

The problem of generating all elements of \mathcal{F}_{π_3} has been studied in the machine learning and computational geometry literatures (see [10, 12, 23]), and is motivated by the discovery of missing associations or “holes” in data mining applications (see [3, 21, 22]). [12] gives an algorithm, for solving this problem, whose worst-case time complexity is exponential in the dimension n of the given point set.

3.2 The joint generation algorithm

As before, we assume that we are given an integer box $\mathcal{C}^* = \mathcal{C}_1^* \times \dots \times \mathcal{C}_n^*$, where $\mathcal{C}_i^* = [l_i^* : u_i^*]$, and $l_i^* \leq u_i^*$, are integers, and a monotone property π , described by a polynomial time satisfiability oracle, for which it is required to generate the families \mathcal{F}_π and \mathcal{G}_π . The generation algorithm, considered in this paper, is based on a reduction to a dualization algorithm of [5], for which an efficient implementation was given in the previous section. Again, the problem is solved by decomposing it into a number of smaller subproblems and solving each of them recursively. The input to each such subproblem is a sub-box \mathcal{C} of the original box \mathcal{C}^* and two subsets $\mathcal{F} \subseteq \mathcal{F}^*$ and $\mathcal{G} \subseteq \mathcal{G}^*$ of integral vectors, where $\mathcal{F}^* \subseteq \mathcal{F}_\pi$ and $\mathcal{G}^* \subseteq \mathcal{G}_\pi$ denote respectively the subfamilies of minimal satisfying and maximal non-satisfying vectors that have been generated so far. Note that, by definition, the following condition holds for the original problem and all subsequent subproblems:

$$a \not\leq b, \quad \text{for all } a \in \mathcal{F}, b \in \mathcal{G}. \quad (12)$$

Given an element $a \in \mathcal{F}$ ($b \in \mathcal{G}$), we say that a coordinate $i \in [n] \stackrel{\text{def}}{=} \{1, \dots, n\}$ is *essential* for a (respectively, b), in the box $\mathcal{C} = [l_1 : u_1] \times \dots \times [l_n : u_n]$, if $a_i > l_i$ (respectively, if $b_i < u_i$). Let us denote by $\text{Ess}(x)$ the set of essential coordinates of an element $x \in \mathcal{F} \cup \mathcal{G}$. Finally, given a sub-box $\mathcal{C} \subseteq \mathcal{C}^*$, and two subsets $\mathcal{F} \subseteq \mathcal{F}^*$ and $\mathcal{G} \subseteq \mathcal{G}^*$, we shall say that \mathcal{F} is *dual to* \mathcal{G} in \mathcal{C} if $\mathcal{F}^+ \cup \mathcal{G}^- \supseteq \mathcal{C}$.

The same decomposition lemma, stated in Section 2.1, holds also here. More precisely, either (i) there is an element $x \in \mathcal{F} \cup \mathcal{G}$ with at most $1/\epsilon$ essential coordinates, where $\epsilon \stackrel{\text{def}}{=} 1/(1 + \log m)$ and $m \stackrel{\text{def}}{=} |\mathcal{F}| + |\mathcal{G}|$, or (ii) one can easily find a new element $z \in \mathcal{C} \setminus (\mathcal{F}^+ \cup \mathcal{G}^-)$, by picking each element z_i independently at random from $\{l_i, u_i\}$ for $i = 1, \dots, n$. As before, in case (i), one can decompose the problem into two strictly smaller subproblems, and the algorithm is guaranteed to reduce the cardinality of one of the sets \mathcal{F} or \mathcal{G} by a factor of at least $1 - \epsilon$ at each recursive step.

We used the following data structures in our implementation:

- Two arrays of vectors, F and G containing the elements of \mathcal{F}^* and \mathcal{G}^* respectively.
- Two (dynamic) arrays of indices, $\text{index}(\mathcal{F})$ and $\text{index}(\mathcal{G})$, containing the indices of vectors from \mathcal{F}^* and \mathcal{G}^* (i.e. containing pointers to elements of the arrays F and G), that appear in the current subproblem. These arrays are used to enable sampling from the sets \mathcal{F} and \mathcal{G} , and also to keep track of which vectors are currently active, i.e. intersect the current box.
- Two balanced binary search trees $\mathbf{T}(\mathcal{F}^*)$ and $\mathbf{T}(\mathcal{G}^*)$, built on the elements of \mathcal{F}^* and \mathcal{G}^* respectively using lexicographic ordering. Each node of the tree $\mathbf{T}(\mathcal{F}^*)$ ($\mathbf{T}(\mathcal{G}^*)$) contains an index of an element in the array F (G). This way, checking whether a given vector $x \in \mathcal{C}$ belongs to \mathcal{F}^* (\mathcal{G}^*) or not, takes only $O(n \log |\mathcal{F}^*|)$ ($O(n \log |\mathcal{G}^*|)$) time.

In the sequel, we let $m = |\mathcal{F}| + |\mathcal{G}|$ and $\epsilon = 1/(1 + \log m)$. We use the following subroutines in our implementation:

Minimization $\min_{\mathcal{F}}(z)$. It takes as input a vector $z \in \mathcal{F}^+$ and returns a minimal vector z^* in $\mathcal{F}^+ \cap \{z\}^-$. As explained in the proof of Proposition 1, such a vector $z^* = \min_{\mathcal{F}}(z)$ can be computed by coordinate descent. More efficient procedures can be obtained if we specialize this routine to the specific monotone property under consideration. For instance, for property π_1 this operation can be performed in $O(nr)$ steps as follows. For $j = 1, \dots, r$, let $a_j x \geq b_j$ be the j th inequality of the system. We initialize $z^* \leftarrow z$ and $w_j = a_j z^* - b_j$ for $j = 1, \dots, r$. For the i th step of the coordinate descend operation, we let $z_i^* \leftarrow z_i^* - \lfloor \lambda \rfloor$ and $w_j \leftarrow w_j - \lfloor \lambda \rfloor a_{ji}$ for $j = 1, \dots, r$, where

$$\lambda = \min_{1 \leq j \leq r} \left\{ \frac{w_j}{a_{ji}} : a_{ji} > 0 \right\}.$$

Now consider property π_2 . Given a set $Z \in \mathcal{F}_{\pi_2}^+$, the operation $\min_{\mathcal{F}_{\pi_2}}(Z)$ can be done in $O(nr)$ by initializing $Z^* \leftarrow Z$, $s \leftarrow |S(Z)|$ and $c(Y) \leftarrow |Z \setminus Y|$ for all $Y \in \mathcal{D}$, and repeating, for $i \in Z$, the following two steps: (i) $\mathcal{Y} \leftarrow \{Y \in \mathcal{D} : c(Y) = 1, Y \not\ni i\}$; (ii) if $|\mathcal{Y}| + s \leq t - 1$ then 1. $Z^* \leftarrow Z^* \setminus \{i\}$, 2. $s \leftarrow s + |\mathcal{Y}|$, and 3. $c(Y) \leftarrow c(Y) - 1$ for each $Y \in \mathcal{D}$ such that $Y \not\ni i$.

For the monotone property π_3 and $z \in \mathcal{C}$, the operation $\min_{\mathcal{F}_{\pi_3}}(z)$ can be done in $O(n|S|)$ as follows. For each point $p \in S$, let $p' \in \mathbb{R}^{2n}$ be the point with components $p'_i = p_i$ for $i = 1, \dots, n$, and $p'_i = u_{i-n} - p_{i-n}$ for $i = n+1, \dots, 2n$. Initialize $s(z) \leftarrow |\{p \in S : p \text{ is in the interior of } z\}|$ and $c(p) \leftarrow |\{i \in [n] : p'_i \leq z_i\}|$ for all $p \in S$. Repeat, for $i = 1, \dots, 2n$, the following steps: (i) $z_i^* \leftarrow \min\{p'_i : p \in S, c(p) = 1, p'_i \leq z_i \text{ and } |\{q \in S : c(q) = 1, p'_i < q_i \leq z_i\}| \leq t - s(z)\}$; (ii) $c(p) \leftarrow c(p) - 1$ for each $p \in S$ such that $z_i^* < p'_i \leq z_i$. Note that (i) can be performed in $O(|S|)$ steps assuming that we know the sorted order for the points along each coordinate.

Maximization $\max_{\mathcal{G}}(z)$. It computes, for a given vector $z \in \mathcal{G}^-$, a maximal vector $z^* \in \mathcal{G}^- \cap z^+$. Similar to $\min_{\mathcal{F}}(z)$, this problem can be done, in general, by coordinate descent. For \mathcal{G}_{π_1} , \mathcal{G}_{π_2} and \mathcal{G}_{π_3} , this operation can be done in $O(nr)$, $O(nr)$, and $O(n|S|)$ respectively.

Below, we denote respectively by T_{min} and T_{max} the maximum time taken by the routines $\min_{\mathcal{F}_{\pi}}(z)$ and $\max_{\mathcal{G}_{\pi}}(z)$ on any point $z \in \mathcal{C}$.

Exhaustive duality $(\mathcal{C}, \pi, \mathcal{F}, \mathcal{G})$. Assuming $|\mathcal{F}||\mathcal{G}| \leq 1$, check if there are no other vectors in $\mathcal{C} \setminus (\mathcal{F}^+ \cup \mathcal{G}^-)$ as follows. First, if $|\mathcal{F}| = |\mathcal{G}| = 1$ then find an $i \in [n]$ such that $a_i > b_i$, where $\mathcal{F} = \{a\}$ and $\mathcal{G} = \{b\}$: (Such a coordinate is guaranteed to exist by (12))

1. If there is a $j \neq i$ such that $b_j < u_j$ then let $z = (u_1, \dots, u_{i-1}, b_i, u_{i+1}, \dots, u_n)$.
2. Otherwise, if there is a $j \neq i$ such that $a_j > l_j$ then let $z = (u_1, \dots, u_{j-1}, a_j - 1, u_{j+1}, \dots, u_n)$.
3. If $b_i < a_i - 1$ then let $z = (u_1, \dots, u_{i-1}, a_i - 1, u_{i+1}, \dots, u_n)$.

In cases 1, 2 and 3, return either $\min_{\mathcal{F}_{\pi}}(z)$ or $\max_{\mathcal{G}_{\pi}}(z)$ depending on whether $\pi(z) = 1$ or $\pi(z) = 0$, respectively.

4. Otherwise return *FALSE* (meaning that \mathcal{F} and \mathcal{G} are dual in \mathcal{C}).

Second, if $|\mathcal{F}| = 0$ then check satisfiability of u . If $\pi(u) = 1$ then return $\min_{\mathcal{F}_{\pi}}(u)$. Else, if $\pi(u) = 0$ then let $z = \max_{\mathcal{G}_{\pi}}(u)$, and return either *FALSE* or z depending on whether $z \in \mathcal{G}^*$ or not (this check can be done in $O(n \log |\mathcal{G}^*|)$ using the search tree $\mathbf{T}(\mathcal{G}^*)$). Finally, if $|\mathcal{G}| = 0$ then check satisfiability of l . If $\pi(l) = 0$ then return $\max_{\mathcal{G}_{\pi}}(l)$. Else, if $\pi(l) = 1$ then let $z = \min_{\mathcal{F}_{\pi}}(l)$, and return either *FALSE* or z depending on whether $z \in \mathcal{F}^*$ or not. This step takes $O(\max\{n \log |\mathcal{F}^*| + T_{min}, n \log |\mathcal{G}^*| + T_{max}\})$ time.

Procedure GEN($\mathcal{C}, \pi, \mathcal{F}, \mathcal{G}$):
Input: A box $\mathcal{C} = \mathcal{C}_1 \times \dots \times \mathcal{C}_n$, a monotone property π , and subsets $\mathcal{F} \subseteq \mathcal{F}_\pi$, and $\mathcal{G} \subseteq \mathcal{G}_\pi$.
Output: Subsets $\mathcal{F}^o \subseteq \mathcal{F}_\pi \setminus \mathcal{F}$ and $\mathcal{G}^o \subseteq \mathcal{G}_\pi \setminus \mathcal{G}$.

1. $\mathcal{F}^o \leftarrow \emptyset, \mathcal{G}^o \leftarrow \emptyset$.
2. While $|\mathcal{F}||\mathcal{G}| \leq 1$
3. $z \leftarrow$ Exhaustive duality($\mathcal{C}, \pi, \mathcal{F}, \mathcal{G}$).
4. If $z = FALSE$ then return($\mathcal{F}^o, \mathcal{G}^o$).
5. else if $\pi(z) = 1$ then $\mathcal{F} \leftarrow \mathcal{F} \cup \{z\}, \mathcal{F}^o \leftarrow \mathcal{F}^o \cup \{z\}, \mathcal{F}^* \leftarrow \mathcal{F}^* \cup \{z\}$.
6. else $\mathcal{G} \leftarrow \mathcal{G} \cup \{z\}, \mathcal{G}^o \leftarrow \mathcal{G}^o \cup \{z\}, \mathcal{G}^* \leftarrow \mathcal{G}^* \cup \{z\}$.
7. end while
8. $z \leftarrow$ Random Solution($\mathcal{C}, \pi, \mathcal{F}^*, \mathcal{G}^*$).
9. While ($z \neq FALSE$) do
10. if $\pi(z) = 1$ then $\mathcal{F} \leftarrow \mathcal{F} \cup \{z\}, \mathcal{F}^o \leftarrow \mathcal{F}^o \cup \{z\}, \mathcal{F}^* \leftarrow \mathcal{F}^* \cup \{z\}$.
11. else $\mathcal{G} \leftarrow \mathcal{G} \cup \{z\}, \mathcal{G}^o \leftarrow \mathcal{G}^o \cup \{z\}, \mathcal{G}^* \leftarrow \mathcal{G}^* \cup \{z\}$.
12. $z \leftarrow$ Random Solution($\mathcal{C}, \pi, \mathcal{F}^*, \mathcal{G}^*$).
13. end while
14. $x^* \leftarrow \operatorname{argmin}\{|\operatorname{Ess}(y)| : y \in (\mathcal{F} \cap \mathcal{C}^-) \cup (\mathcal{G} \cap \mathcal{C}^+)\}$.
15. If $x^* \in \mathcal{F}$ then
16. $i \leftarrow \operatorname{argmax}\{\operatorname{Est}(\{b \in \mathcal{G} : b_j < x_j^*\}, \mathcal{C}) : j \in \operatorname{Ess}(x^*)\}$.
17. $\mathcal{C}' = \mathcal{C}_1 \times \dots \times \mathcal{C}_{i-1} \times [x_i^* : u_i] \times \mathcal{C}_{i+1} \times \dots \times \mathcal{C}_n$.
18. If $\operatorname{Est}(\mathcal{G}, \mathcal{C}') \leq f * |\mathcal{G}|$ then
19. $\mathcal{G}' \leftarrow \operatorname{Cleanup}(\mathcal{G}, \mathcal{C}')$.
20. else
21. $\mathcal{G}' \leftarrow \mathcal{G}$.
22. $(\mathcal{F}_l, \mathcal{G}_l) \leftarrow \operatorname{GEN}(\mathcal{C}', \pi, \mathcal{F}, \mathcal{G}')$.
23. $\mathcal{F}^o \leftarrow \mathcal{F}^o \cup \mathcal{F}_l, \mathcal{F} \leftarrow \mathcal{F} \cup \mathcal{F}_l, \mathcal{F}^* \leftarrow \mathcal{F}^* \cup \mathcal{F}_l$.
24. $\mathcal{G}^o \leftarrow \mathcal{G}^o \cup \mathcal{G}_l, \mathcal{G} \leftarrow \mathcal{G} \cup \mathcal{G}_l, \mathcal{G}^* \leftarrow \mathcal{G}^* \cup \mathcal{G}_l$.
25. $\mathcal{C}'' = \mathcal{C}_1 \times \dots \times \mathcal{C}_{i-1} \times [l_i : x_i^* - 1] \times \mathcal{C}_{i+1} \times \dots \times \mathcal{C}_n$.
26. If $\operatorname{Est}(\mathcal{F}, \mathcal{C}'') \leq f * |\mathcal{F}|$ then
27. $\mathcal{F}'' \leftarrow \operatorname{Cleanup}(\mathcal{F}, \mathcal{C}'')$.
28. else
29. $\mathcal{F}'' \leftarrow \mathcal{F}$.
30. $(\mathcal{F}_r, \mathcal{G}_r) \leftarrow \operatorname{GEN}(\mathcal{C}'', \pi, \mathcal{F}'', \mathcal{G})$.
31. $\mathcal{F}^o \leftarrow \mathcal{F}^o \cup \mathcal{F}_r, \mathcal{F}^* \leftarrow \mathcal{F}^* \cup \mathcal{F}_r, \mathcal{G}^o \leftarrow \mathcal{G}^o \cup \mathcal{G}_r, \mathcal{G}^* \leftarrow \mathcal{G}^* \cup \mathcal{G}_r$.
32. else
- 33-48. Symmetric versions for Steps 16-31 above (details omitted).
49. end if
50. Return ($\mathcal{F}^o, \mathcal{G}^o$).

Random solution($\mathcal{C}, \pi, \mathcal{F}^*, \mathcal{G}^*$). Repeat the following for $k = 1, \dots, t_1$ times, where t_1 is a constant (say 10): Find a random point $z^k \in \mathcal{C}$, by picking each coordinate z_i^k randomly from $\{l_i, u_i\}$, $i = 1, \dots, n$. If $\pi(z^k) = 1$ then let $(z^k)^* \leftarrow \min_{\mathcal{F}_\pi}(z^k)$, and if $(z^k)^* \notin \mathcal{F}^*$ then return $(z^k)^* \in \mathcal{F}_\pi \setminus \mathcal{F}^*$. If $\pi(z^k) = 0$ then let $(z^k)^* \leftarrow \max_{\mathcal{G}_\pi}(z^k)$, and if $(z^k)^* \notin \mathcal{G}^*$ then return $(z^k)^* \in \mathcal{G}_\pi \setminus \mathcal{G}^*$. If $\{(z^1)^*, \dots, (z^{t_1})^*\} \subseteq \mathcal{F}^* \cup \mathcal{G}^*$ then return *FALSE*. This step takes $O(\max\{n \log |\mathcal{F}^*| + T_{min}, n \log |\mathcal{G}^*| + T_{max}\})$ time, and is used to check whether $\mathcal{F}^+ \cup \mathcal{G}^-$ covers a large portion of \mathcal{C} .

Count estimation. For a subset $\mathcal{X} \subseteq \mathcal{F}$ (or $\mathcal{X} \subseteq \mathcal{G}$), use sampling to estimate the number $\operatorname{Est}(\mathcal{X}, \mathcal{C})$ of elements of $\mathcal{X} \subseteq \mathcal{F}$ (or $\mathcal{X} \subseteq \mathcal{G}$) that are active with respect to the current box \mathcal{C} , i.e. the elements of the set $\mathcal{X}' \stackrel{\text{def}}{=} \{a \in \mathcal{X} \mid a^+ \cap \mathcal{C} \neq \emptyset\}$ ($\mathcal{X}'' \stackrel{\text{def}}{=} \{b \in \mathcal{X} \mid b^- \cap \mathcal{C} \neq \emptyset\}$). This can be done as follows. For $t_2 = O(\log(|\mathcal{F}| + |\mathcal{G}|)/\epsilon)$, pick elements $x^1, \dots, x^{t_2} \in \mathcal{F}$ at random, and let the random variable $Y =$

$\frac{|\mathcal{F}|}{t_2} * |\{x^i \in \mathcal{X}' : i = 1, \dots, t_2\}|$. Repeat this step independently for a total of $t_3 = O(\log(|\mathcal{F}| + |\mathcal{G}|))$ times to obtain t_3 estimates Y^1, \dots, Y^{t_3} , and let $\text{Est}(\mathcal{X}, \mathcal{C}) = \min\{Y^1, \dots, Y^{t_3}\}$. This step requires $O(n \log^3 m)$ time.

Cleanup(\mathcal{F}, \mathcal{C}) (Cleanup(\mathcal{G}, \mathcal{C})). Set $\mathcal{F}' \leftarrow \{a \in \mathcal{F} \mid a^+ \cap \mathcal{C} \neq \emptyset\}$ (respectively, $\mathcal{G}' \leftarrow \{b \in \mathcal{G} \mid b^- \cap \mathcal{C} \neq \emptyset\}$), and return \mathcal{F}' (respectively, \mathcal{G}'). This step takes $O(n|\mathcal{F}|)$ (respectively, $O(n|\mathcal{G}|)$).

Now, we describe the implementation of procedure $\text{GEN}(\mathcal{C}, \pi, \mathcal{F}, \mathcal{G})$ which is called initially using $\mathcal{C} \leftarrow \mathcal{C}^*$, $\mathcal{F} \leftarrow \emptyset$ and $\mathcal{G} \leftarrow \emptyset$. At the return of this call, the families \mathcal{F}^* and \mathcal{G}^* , which are initially empty, are extended respectively by the elements in \mathcal{F}_π and \mathcal{G}_π . Below we assume that $f \in (0, 1)$ is a constant, say $1/2$. The families \mathcal{F}^o and \mathcal{G}^o represent respectively the subfamilies of \mathcal{F}_π and \mathcal{G}_π that are generated at each recursion tree node.

The following result, regarding the expected running time of the algorithm, is inherent from Lemma 2.

Proposition 3 *The expected number of recursive calls until a new element in $(\mathcal{F}_\pi \setminus \mathcal{F}^*) \cup (\mathcal{G}_\pi \setminus \mathcal{G}^*)$ is output, or procedure $\text{GEN}(\mathcal{C}, \pi, \mathcal{F}, \mathcal{G})$ terminates is $nm^{O(\log^2 m)}$.*

However, as we shall see from the experiments, the algorithm seems to practically behave much more efficiently than indicated by Proposition 3. In fact, in most of the experiments we performed, we got an almost average linear delay (in m) for generating a new point in $(\mathcal{F}_\pi \setminus \mathcal{F}^*) \cup (\mathcal{G}_\pi \setminus \mathcal{G}^*)$.

4 Experimental results

4.1 Hypergraph transversal

We performed a number of experiments to evaluate our implementation of the hypergraph transversal algorithm. Five types of hypergraphs were used in the experiments:

- Random (denoted henceforth by $R(n, \alpha, d)$): this is a hypergraph with α hyperedges, each of which is picked randomly by first selecting its size k uniformly from $[2 : d]$ and then randomly selecting k elements of $[n]$ (in fact, in some experiments, we fix $k = d$ for all hyperedges).
- Matching ($M(n)$): this is a graph on n vertices (n is even) with $n/2$ edges forming an induced matching.
- Matching Dual ($MD(n)$): this is just $M(n)^d$, the transversal hypergraph of $M(n)$. In particular, it has $2^{n/2}$ hyperedges on n vertices.
- Threshold graph ($TH(n)$): this is a graph on n vertices numbered from 1 to n (where n is even), with edge set $\{\{i, j\} : 1 \leq i < j \leq n, j \text{ is even}\}$ (i.e., for $j = 2, 4, \dots, n$, there is an edge between i and j for all $i < j$). The reason we are interested in such kind of graphs is that they are known to have both a small number of edges (namely, $n^2/4$) and a small number of transversals (namely, $n/2 + 1$ for even n).
- Self-dualized threshold graph ($SDTH(n)$): this is a self-dual hypergraph \mathcal{H} on n vertices obtained from the threshold graph and its dual $TH(n-2)$, $TH(n-2)^d \subseteq 2^{[n-2]}$ as follows:

$$\mathcal{H} = \{\{n-1, n\}\} \cup \{\{n-1\} \cup H \mid H \in TH(n-2)\} \cup \{\{n\} \cup H \mid H \in TH(n-2)^d\}.$$

This gives a family of hypergraphs with polynomially bounded input and output sizes $|SDTH(n)| = |SDTH(n)^d| = (n-2)^2/4 + n/2 + 1$.

- Self-dualized Fano-plane product ($SDFP(n)$): this is constructed by starting with the hypergraph $\mathcal{H}_0 = \{\{1, 2, 3\}, \{1, 5, 6\}, \{1, 7, 4\}, \{2, 4, 5\}, \{2, 6, 7\}, \{3, 4, 6\}, \{3, 5, 7\}\}$ (which represents the set of lines in a Fano plane and is self-dual), taking $k = (n - 2)/7$ disjoint copies $\mathcal{H}_1, \dots, \mathcal{H}_k$ of \mathcal{H}_0 , and letting $\mathcal{H} = \mathcal{H}_1 \cup \dots \cup \mathcal{H}_k$. The dual hypergraph \mathcal{H}^d is just the hypergraph of all 7^k unions obtained by taking one hyperedge from each of the hypergraphs $\mathcal{H}_1, \dots, \mathcal{H}_k$. Finally, we define the hypergraph $SDFP(k)$ to be the hypergraph of $1 + 7k + 7^k$ hyperedges on n vertices, obtained by self-dualizing \mathcal{H} as we did for threshold graphs.

n	30			50			60		
$\alpha = \mathcal{H} $	275	213	114	507	441	342	731	594	520
Total Time (sec)	0.1	0.3	3.1	43.3	165.6	1746.8	322.2	2220.4	13329.5
Output size $ \mathcal{H}^d $	150	450	$5.7 * 10^3$	$1.7 * 10^4$	$6.4 * 10^4$	$4.7 * 10^5$	$7.5 * 10^4$	$4.7 * 10^5$	$1.7 * 10^6$
Time/Trans. (msec)	0.7	0.7	0.5	2.5	2.6	3.7	4.3	4.7	7.8

Table 1: Performance of the algorithm for the class $R(n, \alpha, d)$ of hypergraphs, where $2 \leq d \leq n - 1$.

n	20	22	24	26	28	30	32	34	36	38	40
$\alpha = \mathcal{H} $	10	11	12	13	14	15	16	17	18	19	20
Total Time (sec)	0.3	0.7	1.4	3.2	7.1	17.8	33.9	80.9	177.5	418.2	813.1
Output size $ \mathcal{H}^d $	2^{10}	2^{11}	2^{12}	2^{13}	2^{14}	2^{15}	2^{16}	2^{17}	2^{18}	2^{19}	2^{20}
Time/Trans. (msec)	0.3	0.3	0.3	0.4	0.4	0.5	0.5	0.6	0.7	0.8	0.8

Table 2: Performance of the algorithm for the induced matching $M(n)$.

n	20	22	24	26	28	30	32	34	36	38	40
$\alpha = \mathcal{H} $	2^{10}	2^{11}	2^{12}	2^{13}	2^{14}	2^{15}	2^{16}	2^{17}	2^{18}	2^{19}	2^{20}
Total Time (sec)	0.15	0.44	1.3	4.28	13.3	42.2	132.7	421.0	1330.3	4377.3	14010.5
Output size $ \mathcal{H}^d $	10	11	12	13	14	15	16	17	18	19	20
Time/Trans. (msec)	15	40	108	329	950	$2.8 * 10^3$	$8.2 * 10^3$	$2.5 * 10^4$	$7.4 * 10^4$	$2.3 * 10^5$	$7.0 * 10^5$

Table 3: Performance of the algorithm for the matching dual $MD(n)$.

The experiments were performed on a *Pentium 4* processor with 2.2 GHz of speed and 512M bytes of memory. Tables 1-6, summarizes our results for several instances of the different classes of hypergraphs listed above. In the tables, for the specified hypergraphs with the specified parameters, we show: (i) the total CPU time, in seconds, required to generate all transversals, (ii) the total size of the transversal hypergraph, and (iii) the average time per generated transversal. For random hypergraphs, the time and output size reported are the averages over 30 experiments. For random hypergraphs, we only show results for $n \leq 60$ vertices. For larger numbers of vertices, the number of transversals becomes very large (although the delay between successive transversals is still acceptable).

We also performed some experiments to compare different implementations of the algorithm and to study the effect of increasing the number of vertices and the number of hyperedges on the performance. In particular, Figure 1 shows the effect of rebuilding the tree each time a transversal is generated on the output rate. From this figure we see that the average time per transversal is almost constant if we do not rebuild the tree. In Figure 2, we show that the randomized implementation of the algorithm offers substantial improvement over the deterministic one. Figures 3 and 4, respectively, show how the average CPU time/transversal changes as the number of vertices n and the number of hyperedges α are increased. The plots show that the average CPU time/transversal does not increase more than linearly with increasing α or n .

4.2 Joint generation

We performed a number of experiments to evaluate our joint generation implementation on random instances of the three monotone properties described in Section 3.1. The experiments were also performed on a *Pentium*

n	20	40	60	80	100	120	140	160	180	200
$\alpha = \mathcal{H} $	100	400	900	1600	2500	3600	4900	6400	8100	10000
Total Time (sec)	0.02	0.4	1.9	6.0	18.4	40.2	78.2	142.2	232.5	365.0
Output size $ \mathcal{H}^d $	11	21	31	41	51	61	71	81	91	101
Time/Trans. (msec)	1.8	33.3	146.2	428.6	$1.2 * 10^3$	$2.5 * 10^3$	$4.6 * 10^3$	$7.9 * 10^3$	$12.2 * 10^3$	$18.3 * 10^3$

Table 4: Performance of the algorithm for the threshold graph $TH(n)$.

n	22	42	62	82	102	122	142	162	182	202
$\alpha = \mathcal{H} $	112	422	932	1642	2552	3662	4972	6482	8192	10102
Total Time (sec)	0.07	0.9	5.9	23.2	104.0	388.3	1164.2	2634.0	4820.6	8720.0
Output size $ \mathcal{H}^d $	112	422	932	1642	2552	3662	4972	6482	8192	10102
Time/Trans. (msec)	0.6	2.1	6.3	14.1	40.8	106	234	406	588	863

Table 5: Performance of the algorithm for the class $SDTH(n)$.

n	9	16	23	30	37
$\alpha = \mathcal{H} $	15	64	365	2430	16843
Total Time (sec)	0.01	0.1	4.8	198.1	11885.1
Output size $ \mathcal{H}^d $	15	64	365	2430	16843
Time/Trans. (msec)	0.7	1.6	13.2	81.5	706

Table 6: Performance of the algorithm for the class $SDFP(n)$.

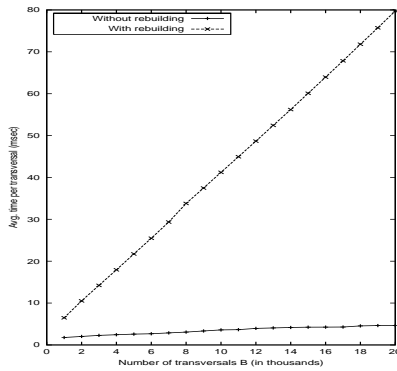


Figure 1: Effect of rebuilding the recursion tree. Each plot shows the average CPU time (in milli-seconds) per generated transversal versus the number of transversals, for hypergraphs of type $R(30, 100, 5)$.

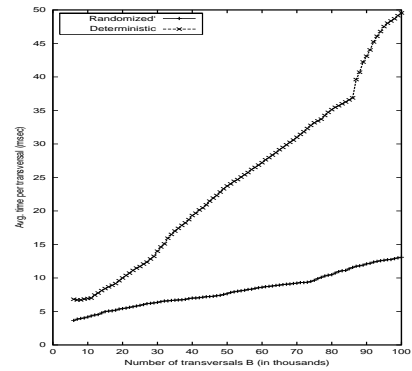


Figure 2: Comparing deterministic versus randomized implementations. Each plot shows the average CPU time/transversal versus the number of transversals, for hypergraphs of type $R(50, 100, 10)$.

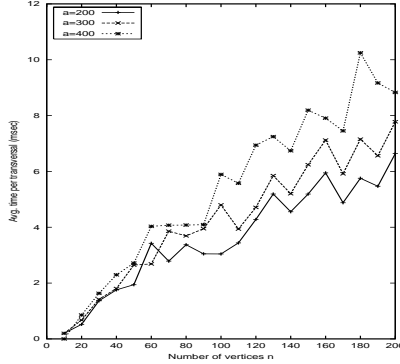


Figure 3: Average CPU time/transversal versus the number of vertices n for random hypergraphs $R(n, a, d)$, where $d = n/4$, and $a = 200, 300, 400$.

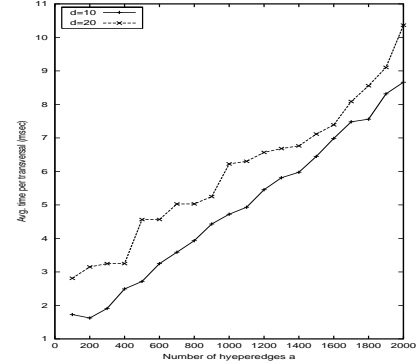


Figure 4: Average CPU time/transversal versus the number of hyperedges a for random hypergraphs $R(50, a, d)$, for $d = 10, 20$.

4 processor with 2.2 GHz of speed and 512M bytes of memory. For each monotone property π , we have limited the corresponding parameters defining the property to reasonable values such that the algorithm completes generation of the sets \mathcal{F}_π and \mathcal{G}_π in reasonable time. Using larger values of the parameters increases the output size, resulting in large total time, although the time per output remains almost constant. For each case, the experiments were performed 5 times, and the numbers shown in the tables below represent averages.

Tables 7 and 8 show our results for linear systems with n variables and r inequalities. Each element of the constraint matrix A and the right-hand side vector b is generated at random from 1 to 15. In the tables we show the output size, the total time taken to generate the output and the average time per each output vector. The parameter c denotes the maximum value that a variable can take. The last row of the table gives the ratio of the size of \mathcal{F}_{π_1} to the size of \mathcal{G}_{π_1} for comparison with the worst case bound of (9). Note that this ratio is relatively close to 1, making joint generation an efficient method for generating both families \mathcal{F}_{π_1} and \mathcal{G}_{π_1} .

n	10		20		30		40		50	
	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}
Output size (thousands)	0.31	0.19	9.9	5.7	49.6	20.0	127.3	59.5	195.3	74.7
Total Time (sec)	4.7	4.7	297	297	1627	1625	5844	5753	10703	10700
Time/output. (msec)	13	24	27	62	29	78	40	103	50	133
Ratio $ \mathcal{G}_{\pi_1} / \mathcal{F}_{\pi_1} $	0.60		0.57		0.40		0.47		0.38	

Table 7: Performance of the algorithm for property π_1 , where $r = 5$ and $c = 2$.

r	5		15		25		35		45	
	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}	\mathcal{F}_{π_1}	\mathcal{G}_{π_1}
Output size (thousands)	20.4	11.6	68.6	27.8	122.7	43.3	196.6	61.7	317.5	115.5
Total Time (sec)	408	408	2244	2242	6495	6482	15857	15856	30170	30156
Time/output. (msec)	20	50	32	90	50	158	76	258	75	260
Ratio $ \mathcal{G}_{\pi_1} / \mathcal{F}_{\pi_1} $	0.57		0.41		0.35		0.31		0.36	

Table 8: Performance of the algorithm for property π_1 , where $n = 30$ and $c = 2$.

Tables 9 and 10 show the results for minimal infrequent/maximal frequent sets. In the tables, n , r and t denote respectively the number of columns, the number of rows of the matrix, and the threshold. Each row of the matrix was generated uniformly at random. As seen from Table 9, for $t = 1, 2$, the bias between the numbers of maximal frequent sets and minimal infrequent sets, for the shown random examples, seem

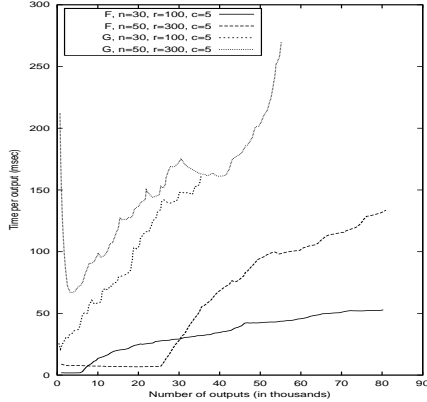


Figure 5: Average time per output, as a function of the number of outputs for minimal feasible/maximal infeasible solutions of linear systems, with $c = 5$ and $(n, r) = (30, 100), (50, 300)$.

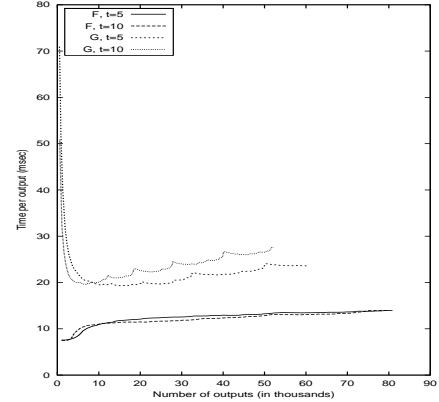


Figure 6: Average time per output, as a function of the number of outputs for minimal infrequent/maximal frequent sets, with $n = 30, r = 1000$ and $t = 5, 10$.

to be large. This makes joint generation an efficient method for generating minimal infrequent sets, but inefficient for generating maximal frequent sets for these examples. However, we observed that this bias in numbers decreases as the threshold t becomes larger. Table 10 illustrates this on a number of examples in which larger values of the threshold were used.

n, r, t	20,100,1		30,100,1		40,100,1		30,100,2		30,300,2		30,500,2	
	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}
Output size (thousands)	2.9	0.08	51.7	0.1	337.1	0.1	75.4	2.5	386.7	13.5	718.3	27.1
Total Time (sec)	60	55	1520	769	22820	3413	1962	1942	13269	13214	28824	28737
Time/output. (msec)	20	690	30	7742	68	34184	28	770	33	979	40	1062
Ratio $ \mathcal{G}_{\pi_2} / \mathcal{F}_{\pi_2} $	0.0280		0.0019		0.0002		0.0335		0.0350		0.0377	

Table 9: Performance of the algorithm for property π_2 for threshold $t = 1, 2$.

n, r, t	30,300,3		30,300,5		30,300,7		30,300,9		30,1000,20		30,1000,25	
	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}	\mathcal{F}_{π_2}	\mathcal{G}_{π_2}
Output size (thousands)	403.3	73.6	362.6	134.7	269.0	100.1	199.1	74.0	491.3	145.7	398.1	114.5
Total Time (sec)	7534	7523	6511	6508	4349	4346	3031	3029	13895	13890	9896	9890
Time/output. (msec)	19	102	18	48	17	43	15	41	28	95	25	86
Ratio $ \mathcal{G}_{\pi_2} / \mathcal{F}_{\pi_2} $	0.1826		0.3715		0.3719		0.3716		0.2965		0.2877	

Table 10: Performance of the algorithm for property π_2 for large threshold values.

Figures 5 and 6 show how the output rate changes for minimal feasible/maximal infeasible solutions of linear systems and for minimal infrequent/maximal frequent sets, respectively. For minimal feasible solutions, we can see that the output rate changes almost linearly as the number of outputs increases. This is not the case for the maximal infeasible solutions, where the algorithm efficiency decreases (the generation problem for maximal infeasible solutions is NP-hard). For minimal infrequent and maximal frequent sets, Figure 6 shows that the output rate increases very slowly. This illustrates somehow that the algorithm practically behaves much better than the quasi-polynomial bound stated in Proposition 3.

Table 11 shows the results for maximal sparse/minimal non-sparse boxes with dimension n , for a set of r random points, threshold t , and upper bound c on the coordinate of each point. As in the case of frequent sets, the bias between the numbers \mathcal{F}_{π_3} and \mathcal{G}_{π_3} is large for $t = 0$ but seems to decrease with larger values of the threshold. In fact, the table shows two examples in which the number of minimal non-sparse boxes is

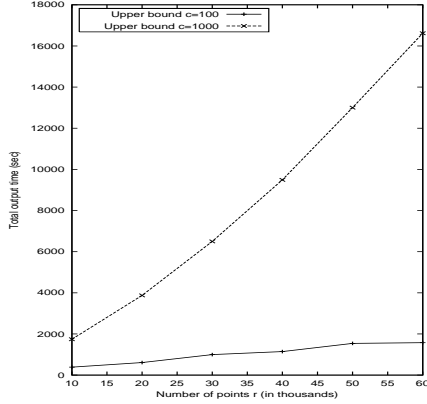


Figure 7: Total generation time as a function of the number of points r for maximal boxes with $n = 2$, $t = 0$, and $c = 100, 1000$.

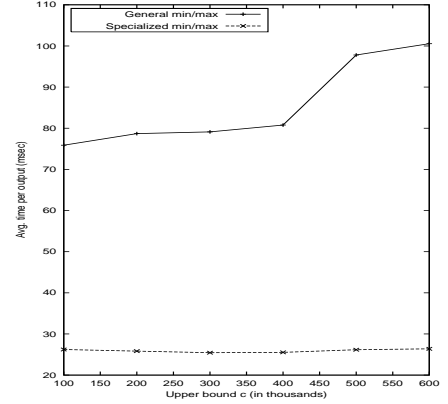


Figure 8: Comparing general versus specialized minimization for property π_3 . Each plot shows the average CPU time/maximal empty box generated versus the upper bound c , for $n = 5$ and $r = 500$.

larger than the number of maximal sparse boxes. We are not aware of any implementation of an algorithm for generating maximal sparse boxes except for [12] which presents some experiments for $n = 2$ and $t = 0$. Experiments in [12] indicated that the algorithm suggested there is almost linear in the the number of points r . Figure 7 illustrates a similar behaviour exhibited by our algorithm. In the figure, we show the total time required to generate all the 2-dimensional maximal empty boxes, as the number of points is increased from 10,000 to 60,000, for two different values of the upper bound c .

r, t	100,0		300,0		500,0		300,2		300,6		300,10	
	\mathcal{F}_{π_3}	\mathcal{G}_{π_3}	\mathcal{F}_{π_3}	\mathcal{G}_{π_3}	\mathcal{F}_{π_3}	\mathcal{G}_{π_3}	\mathcal{F}_{π_3}	\mathcal{G}_{π_3}	\mathcal{F}_{π_3}	\mathcal{G}_{π_3}	\mathcal{F}_{π_3}	\mathcal{G}_{π_3}
Output size (thousands)	16.1	0.1	49.1	0.3	72.9	0.5	228.5	88.7	373.4	466.3	330.4	456.5
Total Time (sec)	932	623	2658	1456	3924	2933	8731	8724	17408	17404	16156	16156
Time/output. (msec)	29	6237	27	4866	27	5889	19	98	23	37	24	35
Ratio $ \mathcal{G}_{\pi_3} / \mathcal{F}_{\pi_3} $	0.0062		0.0061		0.0068		0.3881		1.2488		1.3818	

Table 11: Performance of the algorithm for property π_3 with $n = 7$ and upper bound $c = 5$.

As mentioned in Section 3.2, it is possible in general to implement the procedures $\min_{\mathcal{F}}(z)$ and $\max_{\mathcal{G}}(z)$ using the coordinate decent method, but more efficient implementations can be obtained if we specialize these procedures to the monotone property under consideration. Figure 8 compares the two different implementations for the property π_3 . Clearly, the gain in performance increases as the upper bound c increases.

Let us finally point out that we have observed that the algorithm tends to run more efficiently when the sets \mathcal{F}_{π} and \mathcal{G}_{π} become closer in size. This observation is illustrated in Figure 9 which plots the average time per output (i.e. total time to output all the elements of $\mathcal{F}_{\pi} \cup \mathcal{G}_{\pi}$ divided by $|\mathcal{F}_{\pi} \cup \mathcal{G}_{\pi}|$) versus the ratio $|\mathcal{G}_{\pi}|/|\mathcal{F}_{\pi}|$. This indicates that, when the elements of the sets \mathcal{F}_{π} and \mathcal{G}_{π} are more uniformly distributed along the space, it becomes easier for the joint generation algorithm to find a new vector not in the already generated sets $\mathcal{F}^* \subseteq \mathcal{F}_{\pi}$ and $\mathcal{G}^* \subseteq \mathcal{G}_{\pi}$.

5 Conclusion

We have presented an efficient implementation of an algorithm for generating maximal independent elements for a family of vectors in an integer box. Experiments show that this implementation performs well in practice.

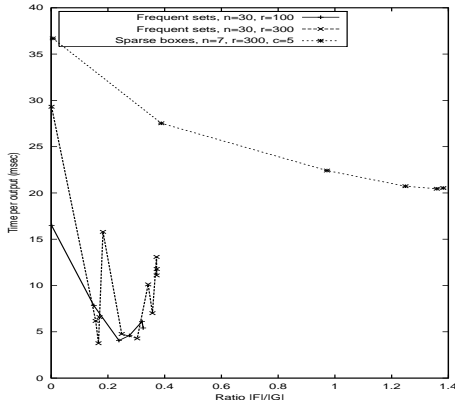


Figure 9: Average generation time as a function of the ratio $|\mathcal{G}_\pi|/|\mathcal{F}_\pi|$, for properties π_2 and π_3 .

We are not aware of any experimental evaluation of algorithms for generating hypergraph transversals except for [19] and [26] in which heuristics for solving this problem were described and experimentally evaluated. However, the results in [19] show the performance for relatively small instances which are easy cases for our implementation. On the other hand, the method described in this paper can handle much larger instances due to the fact that it scales nicely with the size of the problem. In particular, our code can produce, in a few hours, millions of transversals even for hypergraphs with hundreds of vertices and thousands of hyperedges. Furthermore, the experiments also indicate that the delay per transversal scales almost linearly with the number of vertices and number of hyperedges. As for the code in [26], we observed that, although this code seems to perform very well in the case of random hypergraphs, it increasingly slows down with increasing the number of vertices for non-random instances such as self-dualized threshold graphs (see Table 12 for comparison with our implementation for that class of hypergraphs).

n	62	82	102	122	142	162
Our Implementation	20.7	89.1	408.9	1022.5	2280.3	3140.1
[26]	17.9	105.6	514.5	1516.4	3792.0	8656.0

Table 12: Comparing the Performance of our Transversal algorithm with that of [26]. The table shows the total CPU time in seconds required to generate all transversals for hypergraphs in the class $SDTH(n)$.

We have also presented an efficient implementation for a quasi-polynomial algorithm for jointly generating the families \mathcal{F}_π and \mathcal{G}_π of minimal satisfying and maximal non-satisfying vectors for a given monotone property π . We provided experimental evaluation of the algorithm on three different monotone properties. Our experiments indicate that the algorithm behaves much more efficiently than its worst-case time complexity indicates. The algorithm seems to run faster on instances where the families \mathcal{F}_π and \mathcal{G}_π are not very biased in size. Finally, our experiments also indicate that such non-bias in size is not a rare situation (for random instances), despite the fact that inequalities of the form (9)-(11) may hold in general.

Acknowledgements

We thank Thomas Manoukian for providing us with the source code and the reference for [26].

References

- [1] M. Anthony and N. Biggs, *Computational Learning Theory*, Cambridge Univ. Press, 1992.
- [2] R. Agrawal, T. Imielinski and A. Swami, Mining associations between sets of items in massive databases, *Proc. 1993 ACM-SIGMOD Int. Conf.*, pp. 207-216.
- [3] R. Agrawal, H. Mannila, R. Srikant, H. Toivonen and A. I. Verkamo, Fast discovery of association rules, in *Advances in Knowledge Discovery and Data Mining* (U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth and R. Uthurusamy, eds.), pp. 307-328, AAAI Press, Menlo Park, California, 1996.
- [4] J. C. Bioch and T. Ibaraki (1995). Complexity of identification and dualization of positive Boolean functions. *Information and Computation*, **123**, pp. 50–63.
- [5] E. Boros, K. Elbassioni, V. Gurvich, L. Khachiyan and K. Makino, Dual-bounded generating problems: All minimal integer solutions for a monotone system of linear inequalities, *SIAM Journal on Computing*, **31** (5) (2002) pp. 1624-1643.
- [6] E. Boros, K. Elbassioni, V. Gurvich and L. Khachiyan, Generating Dual-Bounded Hypergraphs, *Optimization Methods and Software*, (OMS) **17** (5), Part I (2002), pp. 749–781.
- [7] E. Boros, K. Elbassioni, V. Gurvich, L. Khachiyan and K. Makino, An Intersection Inequality for Discrete Distributions and Related Generation Problems, in *Automata, Languages and Programming, 30-th International Colloquium, ICALP 2003, Lecture Notes in Computer Science (LNCS) 2719* (2003) pp. 543–555.
- [8] E. Boros, K. Elbassioni, V. Gurvich and L. Khachiyan, On enumerating minimal dicuts and strongly connected subgraphs, to appear in *the 10th Conference on Integer Programming and Combinatorial Optimization (IPCO X)*, DIMACS Technical Report 2003-35, Rutgers University, <http://dimacs.rutgers.edu/TechnicalReports/2003.html>.
- [9] E. Boros, V. Gurvich, L. Khachiyan and K. Makino, On the complexity of generating maximal frequent and minimal infrequent sets, in *19th Int. Symp. on Theoretical Aspects of Computer Science, (STACS)*, March 2002, LNCS 2285, pp. 133–141.
- [10] B. Chazelle, R. L. (Scot) Drysdale III and D. T. Lee, Computing the largest empty rectangle, *SIAM Journal on Computing*, 15(1) (1986) 550-555.
- [11] C. J. Colbourn, *The combinatorics of network reliability*, Oxford Univ. Press, 1987.
- [12] J. Edmonds, J. Gryz, D. Liang and R. J. Miller, Mining for empty rectangles in large data sets, in *Proc. 8th Int. Conf. on Database Theory (ICDT)*, Jan. 2001, *Lecture Notes in Computer Science* 1973, pp. 174–188.
- [13] T. Eiter and G. Gottlob, Identifying the minimal transversals of a hypergraph and related problems. *SIAM Journal on Computing*, 24 (1995) pp. 1278-1304.
- [14] M. L. Fredman and L. Khachiyan, On the complexity of dualization of monotone disjunctive normal forms, *Journal of Algorithms*, 21 (1996) pp. 618–628.
- [15] V. Gurvich and L. Khachiyan, On generating the irredundant conjunctive and disjunctive normal forms of monotone Boolean functions, *Discrete Applied Mathematics*, 96-97 (1999) pp. 363-373.
- [16] D. Gunopulos, R. Khardon, H. Mannila and H. Toivonen, Data mining, hypergraph transversals and machine learning, in *Proc. 16th ACM-PODS Conf.*, (1997) pp. 209–216.

- [17] V. Gurvich, To theory of multistep games, *USSR Comput. Math. and Math Phys.* **13-6** (1973), pp. 1485–1500.
- [18] V. Gurvich, Nash-solvability of games in pure strategies, *USSR Comput. Math. and Math. Phys.*, **15** (1975), pp. 357–371.
- [19] D. J. Kavvadias and E. C. Stavropoulos, Evaluation of an algorithm for the transversal hypergraph problem, in *Proc. 3rd Workshop on Algorithm Engineering (WAE'99)*, LNCS 1668, pp. 72–84, 1999.
- [20] E. Lawler, J. K. Lenstra and A. H. G. Rinnooy Kan, Generating all maximal independent sets: NP-hardness and polynomial-time algorithms, *SIAM Journal on Computing*, 9 (1980), pp. 558–565.
- [21] B. Liu, L.-P. Ku and W. Hsu, Discovering interesting holes in data, in *Proc. IJCAI*, pp. 930–935, Nagoya, Japan, 1997.
- [22] B. Liu, K. Wang, L.-F. Mun and X.-Z. Qi, Using decision tree induction for discovering holes in data, in *Proc. 5th Pacific Rim Int. Conf. on Artificial Intelligence*, pp. 182-193, 1998.
- [23] M. Orlowski, A new algorithm for the large empty rectangle problem, *Algorithmica* 5(1) (1990) 65-73.
- [24] R. C. Read, Every one a winner, or how to avoid isomorphism when cataloging combinatorial configurations, *Annals of Disc. Math.* 2 (1978) pp. 107-120.
- [25] R. Srikant and R. Agrawal, Mining quantitative association rules in large relational tables. In *Proc. of the ACM-SIGMOD 1996 Int. Conf. Management of Data*, pp. 1-12, Montreal, Canada, June 1996.
- [26] J. Bailey, T. Manoukian and K. Ramamohanarao, A fast algorithm for computing hypergraph transversals and its application in mining emerging patterns, to appear in *the Proceedings of the 3rd IEEE International Conference on Data Mining, (ICDM)*, 2003, November 19-22, Melbourne, Florida.