



# An average study of hypergraphs and their minimal transversals

Julien David, Loïck Lhote, Arnaud Mary, François Rioult

► **To cite this version:**

Julien David, Loïck Lhote, Arnaud Mary, François Rioult. An average study of hypergraphs and their minimal transversals. *Journal of Theoretical Computer Science (TCS)*, Elsevier, 2015, 596, pp.124-141. <10.1016/j.tcs.2015.06.052>. <hal-01086638>

**HAL Id: hal-01086638**

**<https://hal.archives-ouvertes.fr/hal-01086638>**

Submitted on 24 Nov 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# An Average Study of Hypergraphs and their minimal transversals

Julien David<sup>a</sup>, Loïck Lhote<sup>b</sup>, Arnaud Mary<sup>c</sup>, François Rioult<sup>b</sup>

<sup>a</sup>LIPN, Université Paris 13, and CNRS, UMR 7030. 99, av. J.-B. Clément, 93430 Villetaneuse, France

<sup>b</sup>GREYC, CNRS, ENSICAEN et Université de Caen, Caen, France

<sup>c</sup>Université Lyon 1, CNRS, UMR5558 LBBE, France

---

## Abstract

In this paper, we study some average properties of hypergraphs and the average complexity of algorithms applied to hypergraphs under different probabilistic models. Our approach is both theoretical and experimental since our goal is to obtain a random model that is able to capture the real-data complexity. Starting from a model that generalizes the Erdős-Renyi model [9, 10], we obtain asymptotic estimations on the average number of transversals, minimals and minimal transversals in a random hypergraph. We use those results to obtain an upper bound on the average complexity of algorithms to generate the minimal transversals of an hypergraph. Then we make our random model more complex in order bring it closer to real-data and identify cases where the average number of minimal transversals is at most polynomial, quasi-polynomial or exponential.

---

## 1. Introduction

A hypergraph is a pair  $\mathcal{H} = (V, \mathcal{E})$  where  $V = \{1, 2, \dots, n\}$  is the set of vertices and  $\mathcal{E} = (E_1, \dots, E_m)$  is the collection of hyperedges with  $E_i \subseteq V$  for all  $i$ .

A *transversal* is a set of vertices that intersects all the hyperedges. A set of vertices  $X$  is said to be *irredundant* if for all vertex  $i \in X$ , there exists a hyperedge  $H$  such that  $H \cap X = \{i\}$ .  $X$  is called a *minimal transversal* when it transversal and none of its subset is transversal. This is equivalent to being both irredundant and transversal.

Given a hypergraph  $\mathcal{H}$ , the set of all its minimal transversals forms a hypergraph called the *transversal hypergraph*.

The Transversal Hypergraph Generation problem (for short, THG-problem) consists in computing the transversal hypergraph of a given hypergraph. In the same way, the associated decision problem (in short, THD-problem) consists in deciding if a first hypergraph  $\mathcal{H}_1$  is the transversal hypergraph of a second one  $\mathcal{H}_2$ . This problem is known to be equivalent to the famous dualization of monotone boolean functions problem (see [7]). The Transversal Hypergraph Generation problem appears in very different domains: Artificial Intelligence and Logic [5, 6], Biology [2], Datamining and Machine Learning [12], mobile communications [20], *etc.* We refer to [13] for a more complete list of applications.

1  
2  
3  
4  
5  
6  
7  
8  
9 Since a hypergraph may have an exponential number of minimal transversals, the THG-problem does not belong to the class of polynomial problems. However, a long standing question is to decide whether there exists an algorithm to solve THG whose running time is a polynomial on the size of the hypergraph and on the number of minimal transversal. Such an algorithm is called an *output-polynomial* time algorithm.

10  
11  
12  
13  
14  
15 The complexity of the THG-problem is closely related to its associated decision problem THD. Precisely, if an output-polynomial algorithm solves THG, then THD can be solved in polynomial time. In addition, THD is clearly in the class of co-NP problems but there is no evidence of its co-NP-completeness. If THD is co-NP-complete, then no output polynomial algorithm is likely to exist for the generation problem THG (unless  $P=co-NP$ ) [5].

16  
17  
18  
19  
20  
21  
22 The best known algorithm to generate the transversal hypergraph is quasi-polynomial and is due to Fredman and Khachiyan in [11]. Its running time is of the form  $N^{o(\log N)}$  with  $N$  the size of the input plus the output. Nevertheless, this algorithm is not efficient for practical applications. Other algorithmic solutions were proposed and a list of them can be found in [8]. In this article, we focus on the MTMINER algorithm defined by Hébert, Bretto and Crémilleux [14]. MTMINER is closely related to the mining of the frequent patterns in data mining and is clearly output-exponential in the worst-case. We will study both average complexity and generic-case [15] output-sensitive complexity of the algorithm.

23  
24  
25  
26  
27  
28  
29  
30  
31  
32 In the previous quoted results, the complexity of the THG-problem and associated algorithms were mostly studied with the worst-case point of view. Indeed, very specific entries were exhibited in order to obtain worst-case lower or upper bounds on the behavior of the algorithms. But these entries do not generally occur in practice, and the existing worst-case analyses are then not sufficient to understand the practical complexity of THG. In this article, we adopt a probabilistic point of view. Though analytic combinatorics is often used to conduct an average study, the symbolic method does not seem to be relevant in our case, as it cannot be used to describe the patterns we are interested in.

33  
34  
35  
36  
37  
38  
39  
40  
41  
42 The study of random hypergraphs under various distributions is quite common and one of the most popular is the uniform distribution on  $k$ -uniform hypergraphs [1, 4, 17] (in which all hyperedges have the same cardinal  $k$ ). In [21], the authors prove that under the uniform distribution over all the simple hypergraphs with  $n$  vertices, the THG problem is output-polynomial with probability close to 1. In fact under this distribution, the size of the transversal hypergraph is with high probability exponential in  $n$  and even the naive algorithm that goes through the whole search space is almost surely output-polynomial. To the best of our knowledge, this is the only study on the average complexity of the THG problem.

43  
44  
45  
46  
47  
48  
49  
50  
51  
52 In this paper, we consider two random models in which the number  $n$  of vertices and the number  $m$  of hyperedges are given and suppose that  $m$  is a polynomial in  $n$ . The results we obtain are original and do not intersect with [21].

53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65 Section 2 is devoted to the average analysis of patterns in a random hy-

pergraph. In Section 2.1, we study a *single-parameter* model that generalizes the Erdős-Renyi model [9, 10] on graphs and obtain asymptotic estimations on the average number of transversals, irredundants and minimal transversals. In Section 2.2, we make our random model more complex so that the probability that each vertex appears in an hyperedge is given by a function. We identify functions for which the average number of minimal transversals is at most polynomial, quasi-polynomial or exponential. Section 3 is devoted to algorithms analysis. We study the average complexity of the MTMINER algorithm and the generic-case complexity of the THG-problem. The average complexity of MTMINER is closely related to the average number of irredundants: we obtain upper bounds on the average complexity for both models. Section 4 is devoted to experimental results. Using hypergraphs obtained from real datasets, we discuss the consistency our random models. Conclusion is devoted to perspectives and indications on a random model that might be interesting for a future work.

## 2. Pattern analysis

In this section, we study the average properties of hypergraphs under two probabilistic models. For both models we suppose that:

- The number of hyperedges  $m$  is at most polynomial in the number of vertices  $n$ . Some of our results do not require this supposition and can therefore be extended to cases where  $m$  is exponential in  $n$ . Moreover, when  $m = \Theta(2^n)$ , most questions we study in this paper become trivial.
- The probability that a given vertex  $v$  appears in a given hyperedge is independent from the probability that a vertex  $u$  appears in the same hyperedge or from the probability that  $v$  appears in another hyperedge.
- An hypergraph with  $n$  vertices and  $m$  hyperedges can be seen as a binary matrix  $M(\mathcal{H}) = (m_{i,j}(\mathcal{H}))_{i=1..m,j=1..n}$ . Each row in the matrix encodes an hyperedge. The value  $m_{i,j}$  at line  $i$  and column  $j$  is equal to 1 if the vertex  $j$  belongs to the hyperedge encoded in row  $i$ ,  $m_{i,j} = 0$  otherwise.

### 2.1. The Single-Parameter Model

**Definition 1.** **HG( $\mathbf{n}, \mathbf{m}, \mathbf{p}$ ) random model.** *The  $HG(n, m, p)$  model supposes that the family of random variables  $(m_{i,j})_{i=1..m,j=1..n}$  forms an independent and identically distributed family of random variables following the same Bernoulli law of parameter  $p$  ( $0 < p < 1$ ).*

In this model and in the following results, the reader can consider that  $p$  is a constant but announced results are also valid for  $p$  depending on  $n$  and  $p > 1 - e^{-\frac{1}{\ln n}}$ . This bound will be useful in the multiparametric model.

2.1.1. On the average number of transversals.

In the following,  $T_j$  is the number of transversals of size  $j$  on a given hypergraph. Recall that, a given subset  $X \subseteq V$  of vertices is a transversal if for all hyperedges  $H$  of the hypergraph, there exists at least one vertex  $v \in X$  such that  $v \in H$ . We note  $q = 1 - p$ , the probability that a vertex does not appear in a hyperedge and that  $m$  is the number of hyperedges in a hypergraph. The probability for a subset  $X$  of size  $j$  to be a transversal is therefore:

$$\mathbb{P}(X \text{ is a transversal}) = (1 - q^j)^m.$$

The following results are obtained by calculations on this probability.

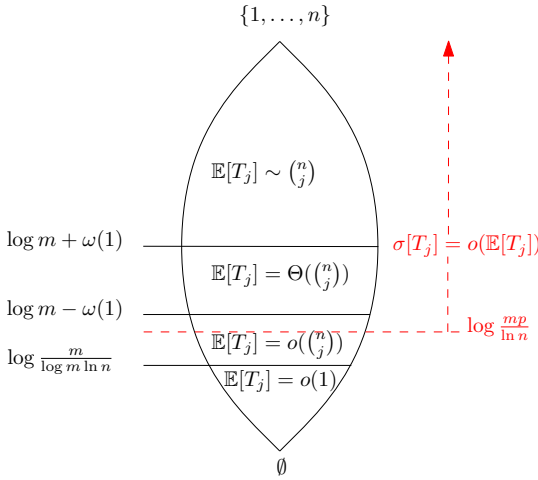


Figure 1: We represented the boolean lattice on the set of vertices. As we will show in the following section, according to the size of the subsets, we can estimate the proportion of transversals. Above a given size, we also have a result on the standard-deviation.

The first result fixes a bound on the size of transversals. The  $\frac{1}{q}$  in the formula are the basis of the logs.

**Lemma 1.** Let  $j_{min} = \log_{\frac{1}{q}} \frac{m}{\log_{\frac{1}{q}} m \ln n}$  be a value chosen for the purpose of calculus. In the  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$  random model, for all  $j < j_{min}$ , the average number of transversals  $\mathbb{E}[T_j]$  tends to 0.

*Proof.* The following inequalities holds:

$$\mathbb{E}[T_j] = \binom{n}{j} (1 - q^j)^m \leq n^j e^{-mq^j} \leq e^{j_{min} \ln n - mq^{j_{min}}}.$$

Now, the expression in the exponential simplifies into

$$j_{min} \ln n - mq^{j_{min}} = -\ln n \times (\log_{\frac{1}{q}} \log_{\frac{1}{q}} m + \log_{\frac{1}{q}} \ln n)$$

which tends to  $-\infty$ . This completes the proof. □

The second lemma gives an order of growth for the average number of transversals whose size is at least logarithmic in  $m$ .

**Lemma 2.** In the  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$  random model, consider  $j$  of the form

$$j = \log_{\frac{1}{q}} m + \log_{\frac{1}{q}} x, \quad x \in ]0, +\infty[.$$

1. if  $x = \Theta(1)$ , only a constant proportion of all sets of  $j$  vertices is a transversal, i.e.,  $\mathbb{E}[T_j] \sim \binom{n}{j} \exp(-1/x)$ .
2. if  $x$  tends to  $+\infty$ , almost all sets of  $j$  vertices is a transversal, i.e.,  $\mathbb{E}[T_j] \sim \binom{n}{j} (1 - \frac{1}{x})$ .
3. if  $x$  tends to 0 with  $x$  lower bounded by  $1/m$ , almost no sets of  $j$  vertices is a transversal, i.e.,  $\mathbb{E}[T_j] = o\left(\binom{n}{j}\right)$ .

*Proof.* Consider a set of vertices  $X$  of size  $j$ . The probability that  $X$  intersects an hyperedge is  $(1 - q^j)$  so that,  $X$  is a transversal with probability

$$\mathbb{P}(X \text{ is a transversal}) = (1 - q^j)^m = \left(1 - \frac{1}{mx}\right)^m$$

for  $j$  as in the lemma. If  $x = \Theta(1)$ , we have the equivalence

$$\left(1 - \frac{1}{mx}\right)^m = \exp\left(-\frac{m}{mx} + O\left(\frac{1}{mx}\right)\right) \sim \exp\left(-\frac{1}{x}\right).$$

If  $x$  tends to  $+\infty$ , the previous equivalence remains true and  $e^{-1/x} \sim 1 - (1/x)$ . To conclude if  $x$  tends to 0, we have

$$\left(1 - \frac{1}{mx}\right)^m \leq \exp\left(-\frac{m}{mx}\right) = \exp\left(-\frac{1}{x}\right)$$

which tends to 0. □

When  $j$  is sufficiently large, the next lemma shows that the standard deviation of  $T_j$  is negligible compared to its mean. Combined with the Markov inequality, this proves that the number of transversals of size  $j$  is almost surely equivalent to the mean number. This result will be fundamental to obtain an almost sure lower bound on the number of minimal transversals.

**Lemma 3.** In the  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$  random model, the standard deviation of the number of transversals of size  $j$  with  $j > \log_{\frac{1}{q}} \frac{mp}{\ln n}$  satisfies

$$\sigma[T_j] = \mathcal{O}\left(\mathbb{E}[T_j] \frac{\ln n}{\sqrt{n}}\right).$$

*Proof.* For any set of vertices  $X$ ,  $\chi_X$  denotes the random equals to 1 if  $X$  is a transversal and 0 otherwise.

Consider  $j \geq 1$ . The variance of  $T_j$  is by definition

$$\mathbb{V}(T_j) = \sum_{\substack{X, Y \subset V \\ |X| = |Y| = j}} \mathbb{P}[\chi_X \chi_Y = 1] - \mathbb{P}[\chi_X = 1] \mathbb{P}[\chi_Y = 1].$$

1  
2  
3  
4  
5  
6  
7  
8  
9 If  $X$  and  $Y$  are disjoint, the random variables  $\chi_X$  and  $\chi_Y$  are independent and the associated term in the previous sum is zero. We are then led to study non disjoint subsets of vertices. If  $X$  and  $Y$  are non disjoint, we define  $I$ ,  $J$  and  $K$  as

$$10 \quad K = X \cap Y, \quad I = X \setminus Y, \quad J = Y \setminus X,$$

11  
12  
13  
14  
15 Note that if  $X$  and  $Y$  have the same cardinality, then the same is true for  $I$  and  $J$ . Therefore we note  $|K| = k$  and  $|I| = |J| = j - k$ . The support of  $K$ , noted  $\mathcal{E}' \subseteq \mathcal{E}$ , is the set of hyperedges that intersect  $K$ . If  $X$  and  $Y$  are transversals, then  $I$  and  $J$  intersects all the hyperedges of  $\mathcal{E} \setminus \mathcal{E}'$ . For a fixed set  $\mathcal{E}'$ , the probability to be  $K$ 's support is  $q^{k \cdot (m - |\mathcal{E}'|)} (1 - q^k)^{|\mathcal{E}'|}$ . The probability that  $I$  and  $J$  intersects all the hyperedges of  $\mathcal{E} \setminus \mathcal{E}'$  is  $(1 - q^{j-k})^{m - |\mathcal{E}'|}$ . Summing over all the possible cardinalities for  $\mathcal{E}'$ , we obtain

$$16 \quad \mathbb{P}[\chi_X = \chi_Y = 1] = \sum_{\ell=0}^m \binom{m}{\ell} (1 - q^k)^\ell q^{k(m-\ell)} (1 - q^{j-k})^{2(m-\ell)}$$

$$17 \quad = (1 - 2q^j + q^{2j-k})^m$$

18  
19 For fixed  $j$  and  $k$ , there are  $\binom{n}{k, j-k, j-k}$  possible choices for the sets  $I$ ,  $J$  and  $K$ . The probability that  $X$  (or  $Y$ ) is a transversal is  $(1 - q^j)^m$ , and summing over all the possible  $k$ , we obtain

$$20 \quad \mathbb{V}(T_j) = \sum_{k=1}^j \binom{n}{k, j-k, j-k} [(1 - 2q^j + q^{2j-k})^m - (1 - q^j)^{2m}].$$

21  
22  
23  
24  
25  
26  
27  
28  
29 Various cases are now possible.

30  
31 **Case (i).**  $j > \log_{\frac{1}{q}} mn$ . Then

$$32 \quad (1 - 2q^j + q^{2j-k})^m - (1 - q^j)^{2m} = \mathcal{O}\left(\frac{1}{n}\right)$$

33  
34 where the constant term in  $\mathcal{O}$  only depends on  $n$  and not on  $j$ . In addition,

$$35 \quad \sum_{k=0}^j \binom{n}{k, j-k, j-k} = \binom{n}{j}^2$$

36  
37 so that the variance satisfies

$$38 \quad \mathbb{V}(T_j) = \left(\frac{1}{n} \binom{n}{j}^2\right) = \mathcal{O}\left(\frac{1}{n} \mathbb{E}_n[T_j]^2\right)$$

39  
40 and the random variable  $T_j$  is concentrated.

1  
2  
3  
4  
5  
6  
7  
8  
9 **Case (ii).** If  $\log_{\frac{1}{q}} \frac{mp}{\ln n} \leq j \leq \log_{\frac{1}{q}} mn$ . The variance satisfies the upper  
10 bounds

$$11 \quad \mathbb{V}(T_j) \leq \sum_{k=1}^j \binom{n}{k, j-k, j-k} (1 - 2q^j + q^{2j-k})^m$$

$$12 \quad \leq \sum_{k=1}^j \binom{n}{k, j-k, j-k} \exp(-2mq^j + mq^{2j-k}).$$

13  
14  
15 Let  $\alpha_k$  denotes the  $k$ -th term of this sum. The ratio of two consecutive  $\alpha_k$  is  
16 upper bounded by

$$17 \quad \frac{\alpha_{k+1}}{\alpha_k} \leq \frac{j^2}{2(n-2j+2)} e^{mpq^{2j}} \leq \frac{\log^2 mn}{2(n-2\log mn+2)} e^{mpq^{\log_{\frac{1}{q}} \left(\frac{mp}{\ln n}\right)^2}}$$

$$18 \quad \leq \frac{e^{\frac{\ln m}{mp}} \log^2 mn}{2(n-2\log^2 mn+2)} = \mathcal{O}\left(\frac{\ln^2 n}{n}\right)$$

19 The variance satisfies  $\mathbb{V}(T_j) \sim \alpha_1$  since  $j\alpha_2 = o(\alpha_1)$ , for  $\log_{\frac{1}{q}} \frac{mp}{\ln n} \leq j \leq$   
20  $\log_{\frac{1}{q}} mn$ . Now, we have

$$21 \quad \mathbb{V}(T_j) \sim e^{-2mq^j} \frac{n^{2j-1}}{(j-1)!^2} \sim \frac{j^2}{n} \mathbb{E}[T_j]^2.$$

22 The result follows for case (ii) and the proof is complete.  $\square$

23 To sum up the results on the average number of transversals:

- 24 1. almost no subset of size less than  $\log_{\frac{1}{q}} \frac{m}{\log_{\frac{1}{q}} m \ln n}$  are transversal,
- 25 2. if the size is  $\log_{\frac{1}{q}} m - \omega(1)$ , then the number of transversals is negligible  
26 compared to the number of subsets,
- 27 3. if the size is equivalent to  $\log_{\frac{1}{q}} m$  (up to an additive constant), then almost  
28 surely a constant proportion of subsets are transversals.
- 29 4. if the size is  $\log_{\frac{1}{q}} m + \omega(1)$ , then almost all subsets are transversals.

30 Intuitively, the set of minimal transversals will mostly be included in cases  
31 2 and 3. Indeed, in case 4, the probability that a given transversal does not  
32 contain a subset of case 3 which is also transversal, will intuitively be low.

### 33 2.1.2. Average number of irredundants.

34 A subset  $X$  is irredundant if for all vertex  $i \in X$ , there exists an hyper-  
35 edge  $H$  such that  $H \cap X = \{i\}$ . For fixed  $X$  of cardinal  $j$ , fixed  $i$  and  $H$ , the  
36 probability that  $H \cap X = \{i\}$  is  $pq^{j-1}$ . For fixed  $X$  and  $H$ , the probability that  
37  $\nexists i \in X$  such that  $H \cap X = \{i\}$  is equal to  $1 - jpq^{j-1}$ . If  $X$  is a irredundant then



there exists a tuple  $(k_1, \dots, k_j)$  of positive value where  $k_i$  is the number of hyperedges  $H$  such that  $H \cap X = \{i\}$  and the probability that  $X$  is a irredundant set is equal to:

$$\sum_{\substack{\forall i \leq j, k_i \geq 1 \\ k_1 + \dots + k_j = \ell \\ j \leq \ell \leq m}} \binom{m}{k_1, \dots, k_j} (pq^{j-1})^\ell \cdot (1 - jpq^{j-1})^{m-\ell}$$

The next theorem states that the average number of irredundants is quasi-polynomial.

**Theorem 1.** *The average number of irredundants in a random model  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$  is*

$$\mathcal{O} \left( (mn \frac{p}{q^2} \log_{\frac{1}{q}} \sqrt{nm})^{\frac{1}{4}(\log_{\frac{1}{q}} nm - \log_q p - \log_{\frac{1}{q}} \log_{\frac{1}{q}} \sqrt{nm}) + 1} \right)$$

*Proof.* We study a function that bounds the number of irredundants by considering the following necessary condition: let  $X = \{x_1, \dots, x_k\}$ , a *selection* is a set  $\{E_1, \dots, E_k\}$  of hyperedges such that  $E_i \cap X = \{x_i\}$  for all  $i \leq k$ .  $X$  is a irredundant set if and only if there exists a selection in the hypergraph.

For each of the  $\binom{n}{j}$  subsets of size  $j$ , if one can find  $j$  hyperedges amongst  $m$  such that each hyperedge contains one vertex and not the others (the order on hyperedges is not specified), then the condition is satisfied. This occurs with probability  $(pq^{j-1})^j$ . The average number of irredundants of size  $j$  is therefore bounded by the function

$$f(j) = \binom{n}{j} \frac{m!}{(m-j)!} (pq^{j-1})^j,$$

We study the value  $j$  for which  $f(j)$  is maximal.

$$\frac{f(j+1)}{f(j)} = \frac{(m-j)(n-j)}{j+1} pq^{2j}.$$

This ratio is decreasing with  $j$ . The maximum of  $h(j)$  is given by the first integer value of  $j$  such that the ratio is smaller than 1.

$$\iff q^{2j} = \frac{1}{(m-j)(n-j)p} (j+1)$$

$$\iff j = \frac{1}{2} \left( \log_{\frac{1}{q}} (m-j)(n-j) + \log_{\frac{1}{q}} p - \log_{\frac{1}{q}} (j+1) \right)$$

It can be proved that with  $p$  constant (or  $p > 1 - e^{-\frac{1}{\ln n}}$ ) and  $m$  large enough, we have

$$\log_{\frac{1}{q}}(m) - 1 < \log_{\frac{1}{q}}(m-j) \leq \log_{\frac{1}{q}}(m).$$

The same goes for  $\log_{\frac{1}{q}}(n-j)$ . We bootstrap and simplify:

$$j \sim \frac{1}{2} \left( \log_{\frac{1}{q}} mn - \log_q p - \log_{\frac{1}{q}} \log_{\frac{1}{q}} nm \right)$$

The maximum of the function  $f$  is reached for  $\lceil j \rceil$ . We obtain the announced result by computing  $f(j)$  (using Stirling's formula) and multiply by  $n$  (all possible values for  $j$ ).  $\square$

### 2.1.3. Average number of minimal transversals

In the sequel,  $MT$  (resp.  $MT_j$ ) is the random variable equal to the number of minimal transversals (resp. of size  $j$ ). It is known that in the worst case, the number of minimal transversals may be exponential with respect to the size of the input hypergraph. However, for the "naive" uniform distribution on hypergraphs, the number of minimal transversal is almost surely at most output linear. As far as we know, the  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$  model leads to the first non trivial bound on the average number of minimal transversals, as announced by the next theorem.

**Theorem 2.** *Consider the random model  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$  with  $m = \beta n^\alpha$ ,  $\beta > 0$  and  $\alpha > 0$ . There exist a positive constant  $c := c(\alpha, \beta, p)$  such that the average number of minimal transversals is*

$$\mathcal{O}\left(n^{d(\alpha) \log_{\frac{1}{q}} m + c \ln \ln m}\right),$$

with  $d(\alpha) = 1$  if  $\alpha \leq 1$  and  $d(\alpha) = \frac{(\alpha+1)^2}{4\alpha}$  otherwise.

*Proof.* We consider once again a function that bounds the number of minimal transversals. A set  $X$  is a minimal transversal if and only if:

1. there exist a selection  $\mathcal{E}'$  for  $X$ ,
2. for all  $F \in \mathcal{E} \setminus \mathcal{E}'$ ,  $|F \cap X| \geq 1$ .

The probability that Condition 1 holds is  $(pq^{j-1})^j$  whereas the probability that Condition 2 holds is  $(1 - q^j)^{m-j}$ . There are  $\binom{n}{j}$  sets of vertices of size  $j$  and each of them have  $\frac{m!}{(m-j)!}$  possible selections. Then, the average number of minimal transversals of size  $j$  is bounded by  $h(j)$  with

$$h(j) = \binom{n}{j} \frac{m!}{(m-j)!} (pq^{j-1})^j (1 - q^j)^{m-j}.$$

We now determine the  $j$  that maximizes  $h(j)$ . The ratio  $\frac{h(j+1)}{h(j)}$  satisfies

$$\frac{h(j+1)}{h(j)} = \frac{(m-j)(n-j)}{j+1} \frac{p}{1-q^j} q^{2j} \left(1 + \frac{pq^j}{1-q^j}\right)^{m-j-1}$$

As in the previous section, this ratio is also decreasing. Two cases are now possible according to the value of  $\alpha$ .

*Case  $\alpha \leq 1$ .*  $h(j)$  is maximized by  $j_{min}$  with:

- $j_{min} = \log_{\frac{1}{q}} \frac{m(1-\alpha)}{p \ln m}$  when  $\alpha < 1$ ,

- $j_{min} = \log_{\frac{1}{q}} \frac{mp}{\ln \ln m}$  when  $\alpha = 1$ .

Precisely, for some real  $x > 0$  (resp.  $x < 0$ ), the ratio  $h(j_{min} + 1 + x)/h(j_{min} + x)$  tends to 0 (resp.  $+\infty$ ) when  $n$  grows so that the average number of minimal transversals satisfies

$$\mathbb{E}[M] \leq \sum_{j=1}^n h(j) = O(h(j_{min} + 1)).$$

Now, asymptotic computations give

$$h(j_{min}) = \exp\left(\ln n \log_{\frac{1}{q}} m + O(\ln m \ln \ln m)\right)$$

which completes the proof in case  $\alpha \leq 1$ .

*Case  $\alpha > 1$ .* Consider  $j_{min} = \frac{1}{2} \log_{\frac{1}{q}} mn - \frac{1}{2} \log_{\frac{1}{q}} \ln m$ . Note that  $j_{min}$  may not maximize  $h(j)$  but if  $x$  tends to  $+\infty$  (resp.  $-\infty$ ) with  $|x| = O(\ln \ln \ln m)$ , the ratio  $h(j_{min} + 1 + x)/h(j_{min} + x)$  tends to 0 (resp.  $+\infty$ ). Then, the average number of minimal transversals satisfies

$$\begin{aligned} \mathbb{E}[M] &\leq \sum_{j=1}^n h(j) \sim \sum_{j=j_{min}-\ln \ln \ln m}^{j_{min}+\ln \ln \ln m} h(j) \\ &= O\left(\ln \ln \ln m \max_{|j-j_{min}| \leq \ln \ln \ln m} h(j)\right). \end{aligned}$$

The result follows from the fact that

$$\max_{|j-j_{min}| \leq \ln \ln \ln m} h(j) = \exp\left(d(\alpha) \ln n \log_{\frac{1}{q}} m + O(\ln m \ln \ln m)\right).$$

□

#### 2.1.4. Almost sure lower bound for the number of minimal transversals

A generic lower bound of a random variable (i.e. a lower bound which is true with probability close to 1) is often obtained by studying the moments of higher order or the variance. We did not succeed in studying the higher moments of the number of minimal transversals. However, we relate the number of minimal transversals to the number of transversals and use the concentration property given in Lemma 3.

The number  $MT$  of minimal transversals is lower bounded by  $M_j$ , the number of minimal transversals of size  $j$ . Among the  $T_j$  transversals of cardinal  $j$ ,  $MT_j$  are irredundant and  $T_j - MT_j$  are supersets of transversals of cardinal  $j - 1$ . By definition, there are  $T_{j-1}$  transversals with cardinal  $j - 1$  and each of these transversals can be completed in at most  $n - j + 1$  transversals of cardinal  $j$ . We deduce the inequalities

$$MT \geq MT_j \quad \text{and} \quad T_j - (n - j + 1)T_{j-1} \leq MT_j \leq T_j. \quad (1)$$

The lower bound entails that for all  $0 < \epsilon < 1$ ,

$$\mathbb{P}(MT_j \leq \epsilon \mathbb{E}[T_j]) \leq \mathbb{P}(T_j - (n - j + 1)T_{j-1} \leq \epsilon \mathbb{E}[T_j])$$

$$\mathbb{P}(MT_j < \epsilon \mathbb{E}[T_j]) \leq \mathbb{P}(T_{j-1} < \epsilon \mathbb{E}[T_{j-1}]) + \mathbb{P}(T_j < \epsilon(\mathbb{E}[T_j] + (n - j + 1)\mathbb{E}[T_{j-1}])).$$

Consider  $l = \log_{\frac{1}{q}} \frac{mp}{\ln n} + 1$  (in order to use Lemma 3). We have

$$\mathbb{E}[T_l] + (n - l + 1)\mathbb{E}[T_{l-1}] = \mathbb{E}[T_l] \left( 1 + O\left(\frac{\ln n}{n}\right) \right)$$

and the Bienaymé-Tchebychev Inequality with Lemma 3 lead to the following proposition.

**Proposition 1.** *Consider  $\epsilon$  with  $0 < \epsilon < 1$ . In the random model  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$ , the number  $MT$  of minimal transversals satisfies*

$$\mathbb{P}(MT < \epsilon \mathbb{E}[T_l]) = O\left(\frac{\ln^2 mn}{(1 - \epsilon)^2 n}\right)$$

where  $T_l$  is the set of transversals of size  $l = \log_{\frac{1}{q}} m - \log_{\frac{1}{q}} \ln n + \log_{\frac{1}{q}} \frac{p}{q}$ .

**Corollary 1.** *In the random model  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$ , the number of minimal transversals is almost surely greater than*

$$n^{\log_{\frac{1}{q}} m + \mathcal{O}(\ln \ln m)}$$

*Proof.* The idea is to compute  $\mathbb{E}[T_l]$  and use Lemma 3. □

## 2.2. Multiparametric Model

In this section, we no longer consider that the vertices occur in an hyperedge with the same probability. As we will show in Section 4, this is much more consistent with real-case datasets.

**Definition 2** ( $HG(n, m, g)$  random model). *A hypergraph  $\mathcal{H}$  with  $n$  hyperedges and  $m$  vertices is seen as a binary matrix  $M(\mathcal{H}) = (m_{i,j}(\mathcal{H}))_{i=1..m, j=1..n}$ . The  $HG(n, m, g)$  model supposes that the family of random variables  $(m_{i,j})_{i=1..m, j=1..n}$  forms an independent family of random variables. In addition, for all  $i, j$ , the random variable  $m_{i,j}$  follows a Bernoulli law of parameter  $p_i = g(i)$  (and  $q_i = 1 - g(i)$ ) with  $g : \mathbb{N} \rightarrow [0, 1]$ .*

We partition the set  $V$  into 3 subsets:

- The set  $U$  of ubiquitous vertices. Let  $x$  be a fixed constant. For all vertices  $u \in U$ , we have that  $q_u < \frac{x}{m}$  with  $q_u = 1 - p_u$ .
- The set  $R$  of rare events. For all vertices  $r \in R$ , we have that  $p_r < 1 - e^{-\frac{1}{\ln n}}$ . Note that this implies that  $p_r < \frac{1}{\ln n}$ . The latter bound slowly tends to zero and is relevant on experimental data (see Section 4), the former simplifies calculus.
- The set  $O$  of other events, that is  $O = V \setminus \{U \cup R\}$ .

We mainly use this decomposition to study the average number of minimal transversals in the  $HG(n, m, g)$  model.

1  
2  
3  
4  
5  
6  
7  
8  
9 *2.2.1. A lower bound on the average number of transversals*

10 In the  $HG(n, m, g)$  model, let  $\mu$  the average value of the random variables  
11  $q_i$ . In other words:

$$12 \mu = \sum_{i=1}^n \frac{q_i}{n}$$

13  
14  
15 **Lemma 4.** *In the  $HG(n, m, g)$  model, the average number of transversals is*  
16  $\Omega(2^n - m(1 + \mu)^n)$ .  
17

18 *Proof.* The average number of transversals is given by the following formula:  
19

$$20 \mathbb{E}[T] = \sum_{j=0}^n \sum_{\substack{X \subset V \\ |X|=j}} (1 - \prod_{i \in X} q_i)^m.$$

21  
22 Using the Bernoulli inequality, we have:  
23

$$24 \mathbb{E}[T] \geq \sum_{j=0}^n \sum_{\substack{X \subset V \\ |X|=j}} (1 - m \prod_{i \in X} q_i) \geq 2^n - m \sum_{j=0}^n \sum_{\substack{X \subset V \\ |X|=j}} \prod_{i \in X} q_i = 2^n - m \prod_{i=1}^n (1 + q_i)$$

25  
26 Then, according to the geometric inequality we have  
27

$$28 \prod_{i=1}^n (1 + q_i) \leq \left( \sum_{i=1}^n \frac{1 + q_i}{n} \right)^n$$

29 and therefore  
30

$$31 \mathbb{E}[T] \geq 2^n - m(1 + \mu)^n$$

32 which concludes the proof.  $\square$   
33

34  
35 **Corollary 2.** *The average number of irredundants (and then minimal transversals)*  
36 *is bounded by  $\mathcal{O}(mn(1 + \mu)^n)$ .*  
37

38  
39 *Proof.* From Lemma 4 we obtain an upper bound on the average number of  
40 subset that are not transversals, that is to say  $\mathcal{O}(m(1 + \mu)^n)$ . In the worst case,  
41 all those sets are irredundant and if we add any vertex to one of those subsets,  
42 it becomes a minimal transversal.  
43  $\square$   
44

45 Although this is not precise, it has a certain advantage: starting from a  
46 complex model with a large number of parameters, we now have an estimation  
47 that relies on only 3 parameters. It can also be easily interpreted: if  $\mu$  tends  
48 to 0 then almost all vertices appears in almost all hyperedges, hyperedges are  
49 all similar. In this case, there is few minimal transversals. If  $\mu$  tends to 1 then  
50 almost each vertex appears in few hyperedges, hyperedges are all really different,  
51 that is to say the pairwise intersection of hyperedges is always small and there  
52 is an exponential number of minimal transversals.  
53  
54  
55  
56  
57  
58

2.2.2. On the average number of irredundants

**Lemma 5.** *The average number of irredundants containing only vertices in  $O \cup U$  is*

$$\mathcal{O}((nm \ln \sqrt{nm})^{\frac{1}{4}} \ln n (\ln nm - 2 \ln \ln n - \ln \ln \sqrt{nm}))$$

*Proof.* In order to obtain an upper bound, we use the result in Theorem 1. In the  $HG(n, m, p)$  model, the upper bound on the average number of irredundants decreases as  $p$  increases. Recall that for all vertices  $a \in O \cup U$ , we have  $p_a \geq 1 - e^{-\frac{1}{\ln n}}$ . The average number of irredundants in the  $HG(n, m, 1 - e^{-\frac{1}{\ln n}})$  model is an upper bound of the average number of irredundants containing only vertices in  $O$ . If  $q = e^{-\frac{1}{\ln n}}$ , for all  $y$  we have  $\log_{1/q} y = \ln y \times \ln n$ . Using this simplification, we obtain the announced upper bound.  $\square$

Let  $Irr_{j,l}$  be the number of irredundants of size  $j$  containing exactly  $l$  rare vertices. We have:

$$Irr_{j,l} \leq Irr_{j-l,0} \times Irr_{l,l}$$

where  $Irr_{j-l,0}$  is exactly the number of irredundants containing only vertices in  $O \cup U$  and  $Irr_{l,l}$  is the number of irredundants containing only vertices in  $R$ . Note that Lemma 5 gives an upper bound on  $Irr_{j-l,0}$ . We obtain a bound on  $Irr_{l,l}$  by adapting the function  $f(j)$  from Section 2.1.2. Since for all vertices  $r \in R$ , we have  $p < 1 - e^{-\frac{1}{\ln n}} < \frac{1}{\ln n}$ , we have:

$$Irr_{j,l} \leq \binom{|R|}{l} \frac{m!}{(m-l)!} \left(\frac{1}{\ln n}\right)^l Irr_{j-l,0}$$

Now we obtain an upper bound on the average number of irredundant in the  $HG(n, m, g)$  model.

$$Irr_j = \sum_{l=0}^{\min\{j, |R|\}} \binom{|R|}{l} \left(\frac{m}{\ln n}\right)^l Irr_{j-l,0}$$

From this formula, it can be deduced that if  $|R|$  is not too large in the size of  $n$ , then the average number of irredundants is at most quasi-polynomial.

**Lemma 6.** *In the  $HG(n, m, g)$  model, if  $|R| = \mathcal{O}((\ln n)^c)$  where  $c$  is a constant, then the number of irredundants is quasi-polynomial.*

*Proof.* From Lemma 5, we know that  $Irr_{j-l,0}$  is at most quasi-polynomial. Since  $l = \mathcal{O}(|R|) = \mathcal{O}((\ln n)^c)$ , it is clear that  $\binom{|R|}{l} \left(\frac{m}{\ln n}\right)^l$  is also quasi-polynomial.  $\square$

**Lemma 7.** *In the  $HG(n, m, g)$  model, the probability to have a polynomial number of irredundants containing ubiquitous vertices tends to 1.*

1  
2  
3  
4  
5  
6  
7  
8  
9 *Proof.* Let  $Irr(\mathcal{H})$  denote the set of irredundants of an hypergraph  $\mathcal{H}$ . For a  
10 fixed vertex  $e \in V$ , let  $Irr(\mathcal{H}, e)$  the set of irredundants containing  $e$ . Let  $N(e)$   
11 be the set of hyperedges not containing the vertex  $e$ . We have:

$$12 \quad Irr(H, e) \subseteq \{X \cup \{e\} \mid X \in Irr(N(e))\}.$$

13  
14 Hence we have:

$$15 \quad |Irr(H, e)| \leq |Irr(N(e))|,$$

16 that is to say the number of irredundants containing a given vertex  $e$  is bounded  
17 by the number of irredundants of the hypergraph reduced to the hyperedges that  
18 do not contain  $e$ . The number of hyperedges that do not contain an ubiquitous  
19 vertex is bounded by  $x$  and using Poisson paradigm we know that in a random  
20 hypergraph an ubiquitous vertex is in at least  $m - x\sqrt{x}$  hyperedges with proba-  
21 bility tending to 1. The number of irredundants of an hypergraph with at most  
22  $x\sqrt{x}$  hyperedges is polynomial since  $x$  is a constant.  $\square$

23  
24  
25  
26 *2.2.3. On the average number of minimal transversals.*

27 Recall that  $M$  denotes the number of minimal transversals. Since the upper  
28 bounds obtained in Lemma 6, 5 and 7 also holds for minimal transversals, we  
29 obtain the following theorem.

30  
31 **Theorem 3.** *In the  $HG(n, m, g)$  model, we have the following:*

- 32 • If  $|O \cup R| = \mathcal{O}(\ln n)$ , then  $\mathbb{E}[M]$  is at most polynomial.
- 33 • If  $|R| = \mathcal{O}((\ln n)^c)$  where  $c$  is a constant, then  $\mathbb{E}[M]$  is at most is quasi-  
34 polynomial.
- 35 • If  $|R| = \Theta(n)$ , then  $\mathbb{E}[M]$  is at most exponential on  $|R|$ .

36  
37  
38  
39 The first point comes from Lemma 7 and the fact that the number of minimal  
40 transversals of a set of size  $c \ln n$  is  $n^c$ .

### 41 42 43 3. Algorithm analysis

44  
45 In this section we study the average complexity of the MT-Miner Algorithm.  
46 The worst-case input complexity of this algorithm is a polynomial in  $n$  times  
47 the number of irredundants. The analysis we perform and the results we ob-  
48 tain is also valid on any algorithm whose search space is bounded by the set  
49 of irredundants: Apriori, Dong-Li algorithm [3], Kavvadias-Stravropoulos [16],  
50 Uno-Murakami [18]...

```

Input: an hypergraph  $\mathcal{H}$  with  $n$  hyperedges
Output: the minimal transversals of  $\mathcal{H}$ 

 $MT := \{\{v\} \mid v \in V, |\text{Supp}_{\mathcal{H}}(\{v\})| = m\}$ 
 $N_1 := \{\{v\} \mid v \in V, n > \text{Supp}_{\mathcal{H}}(\{v\}) \neq \emptyset\}$ 
 $j = 1$ 
While  $N_j \neq \emptyset$  do
  for all prefix  $V$  with  $V \cup \{v_1\}$  and  $V \cup \{v_2\}$ 
  in  $N_j \times N_j$  do
     $W = V \cup \{v_1\} \cup \{v_2\}$ 
    if  $W$  is irredundant then
      if  $\text{Supp}_{\mathcal{H}}(W) = n$  then add  $W$  to  $MT$ 
      else add  $W$  to  $N_{j+1}$  end if
    end if
  end for
   $j=j+1$ 
end While
return  $MT$ .

```

Figure 2: The MTMINER-algorithm. Here  $\text{Supp}_{\mathcal{H}}(W)$  is the number of hyperedges that intersect  $W$ .

### 3.1. Average complexity of the MT-Miner Algorithm.

The MTMINER algorithm was described by Hébert, Bretto and Crémilleux in [14]. The algorithm computes all the minimal transversals of a given hypergraph using a levelwise strategy. Precisely at the  $j^{\text{th}}$  level, the algorithm computes the *irredundants* formed with  $j$  vertices. Among the irredundants, some are minimal transversals and are stored in a data structure. The others are not minimal transversals but they might be part of one and they are used to build irredundants of size  $j + 1$ .

Each irredundant of size  $j$  can be extended in at most  $n - j$  sets of size  $j + 1$  and the minimality of each candidate set can be tested in polynomial time (w.r.t. the input size). MTMINER uses a prefix tree to optimize this generation step but even with the naive method (generate all the possible extensions), the complexity of MTMINER is  $O(\text{Poly}(m, n)N)$  where  $N$  is the number of irredundants. The algorithm is described in Figure 3.1.

A non trivial upper bound on the average complexity of MTMINER follows from Theorem 1.

**Proposition 2.** *In the  $HG(n, m, p)$  model, there exist some positive constant  $c$  such that the average complexity of MTMINER is*

$$\mathcal{O} \left( \left( mn \log_{\frac{1}{q}} \sqrt{nm} \frac{p}{q^2} \right)^{\frac{1}{4} (\log_{\frac{1}{q}} nm - \log_q p - \log_{\frac{1}{q}} \log_{\frac{1}{q}} \sqrt{nm}) + c} \right).$$

An equivalent result can be obtained in the  $HG(n, m, g)$  model using Lemma 6.



**Proposition 3.** *In the  $HG(n, m, g)$  model, if  $|R| = \mathcal{O}((\ln n)^c)$  where  $c$  is a constant, then the average input-complexity of MTMINER is at most quasi-polynomial.*

Since MTMINER generates at least all the minimal transversals, its complexity is lower bounded by  $M$ . The generic lower bound given by Proposition 2 entails a generic lower bound on the complexity of MTMINER.

### 3.2. Generic-case complexity of the THG-problem

The notion of generic-case complexity is defined in [15]. The idea is to study the worst-case complexity of an algorithm on a generic-subset of inputs. In a given random model, a subset  $E$  of inputs is said to be generic if the probability that a random input is in  $E$  tends to 1. The study of generic complexity is particularly interesting when no polynomial method is known to solve a problem in the general case ( $NP$ -complete problems, undecidable problems [19]), whereas there seem to be efficient methods in practice. The theoretical complexity of the THG-problem was discussed in the introduction. In particular, it is not known whether the problem is output-polynomial in the worst-case. The following theorem states that with probability that tends to 1 in the single-parameter model, the algorithm MTMINER is output-polynomial. In other words, the set of inputs for which the algorithm is output-polynomial is a generic set.

**Theorem 4.** *Consider the random model  $\mathbf{HG}(\mathbf{n}, \mathbf{m}, \mathbf{p})$  with  $m = \beta n^\alpha$ ,  $\beta > 0$  and  $\alpha > 0$ . Under this model, the generic complexity of the THG-problem is output-polynomial. Precisely, there exist an algorithm (MTMINER) such that for all  $\epsilon > 0$ , the algorithm computes the minimal transversals of an input hypergraph in time  $M^{\epsilon + \frac{(\alpha+1)^2}{4\alpha}}$  with probability asymptotically 1 and where  $M$  is the number of minimal transversals.*

**Proof.** To simplify the notations, we write  $\gamma = \epsilon + \frac{(\alpha+1)^2}{4\alpha}$  and  $a = \frac{1}{2}\mathbb{E}[T_j]$  with  $j$  as in Proposition 1. We have

$$\begin{aligned} & \mathbb{P}(D > M^\gamma) \\ &= \mathbb{P}([D > M^\gamma] \cap [M < a]) + \mathbb{P}([D > M^\gamma] \cap [M \geq a]) \\ &\leq \mathbb{P}(M < a) + \mathbb{P}(D \geq a^\gamma) \\ &= O\left(\frac{\ln^2 mn}{n}\right) + \frac{\mathbb{E}[D]}{a^\gamma} \end{aligned}$$

Alternative expressions for  $a^\gamma$  and  $\mathbb{E}[D]$  are

$$\begin{aligned} \mathbb{E}[D] &= n^{\frac{(\alpha+1)^2}{4\alpha} \log \frac{1}{q} m + O(\ln \ln m)}, \\ a^\gamma &= n^{(\epsilon + \frac{(\alpha+1)^2}{4\alpha}) \log \frac{1}{q} m + O(\ln \ln m)}. \end{aligned}$$

In particular,  $\mathbb{E}[D]/a^\gamma$  is  $O(n^{-\frac{\epsilon}{2} \log \frac{1}{q} m})$  when  $\gamma > \epsilon + \frac{(\alpha+1)^2}{4\alpha}$  and the  $\mathcal{O}$ -term tends to 0. This completes the proof.  $\square$

1  
2  
3  
4  
5  
6  
7  
8  
9 **4. Comparison with real data**

10 The following benchmark were made using datasets from the *Frequent Item-*  
11 *set Mining Dataset Repository*<sup>1</sup>, that are often used by the datamining commu-  
12 nity. The experimental results that we exhibit might therefore be well known  
13 by specialists. The objective of this section is to exhibit the links between real  
14 datasets and our probabilistic models.  
15

16 We only present a selection of experimental results, that is already able to  
17 capture the diversity of hypergraphs in various contexts. It is important to note  
18 that for each example presented in this paper, one can find several datasets  
19 satisfying the same properties.

20 From Figure 3 to Figure 6, the histograms on the left side represent the num-  
21 ber of hyperedges in which each vertex appears. Vertices are sorted according  
22 to the decreasing order of the number of hyperedges in which they appear. His-  
23 tograms on the right side represents the size of hyperedges in each hypergraph.  
24

- 25 • *mushroom.dat*: contains a few ubiquitous vertices and a few rare vertices.  
26 All the hyperedges have the same size. Those kinds of datasets validate  
27 the choice of most study to focus on  $k$ -uniform hypergraphs.  
28
- 29 • *accident.dat*: contains a few ubiquitous vertices and a lot of rare vertices.  
30 The size of hyperedges seems to follow a Gaussian distribution.  
31
- 32 • *pumsbstar.dat*: does not contain ubiquitous vertices and most vertices are  
33 rare.  
34
- 35 • *T10I4D100K.dat*: all vertices are rare. Again, the size of hyperedges seems  
36 to follow a Gaussian distribution.  
37

38 Our probability models seem to be more appropriate on the second and  
39 the fourth example, as the size of the hyperedges seems to follow a Gaussian  
40 distribution.

41 In our experiments, the generation of the minimal transversals were efficient  
42 on databases that have similar distribution to *mushroom.dat*. This result com-  
43 forts us in the belief that the presence of rare events is the main parameter  
44 (by opposition to being just an important parameter amongst others) to decide  
45 whether the number of the minimal transversals is going to explode. Figure 7  
46 shows the distribution of the minimal transversals in *mushroom.dat*. As we can  
47 see, the number of minimal transversals  $T_{min,j}$  of a given size  $j$  is maximal when  
48  $j = 16$  and  $T_{min,16} \sim 6 \times 10^6$ . Using the model  $H(n, m, g)$  and Theorem 3, we  
49 could foretell that in cases *accidents.dat*, *pumsbstar.dat*, *T10I4D100K.dat* the  
50 number of minimal transversals explodes, which seems indeed to be the case.  
51 Even after a long time execution (more than a week) on a regular computer,  
52 only a small proportion of the search space had been visited, whereas the num-  
53 ber of minimal transversals was tremendously huge. We also tried to write the  
54  
55

---

56 <sup>1</sup><http://fimi.ua.ac.be/data/>  
57  
58

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

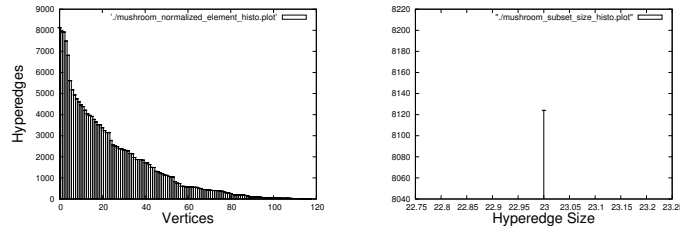


Figure 3: *mushroom.dat* (119 vertices, 8124 hyperedges)

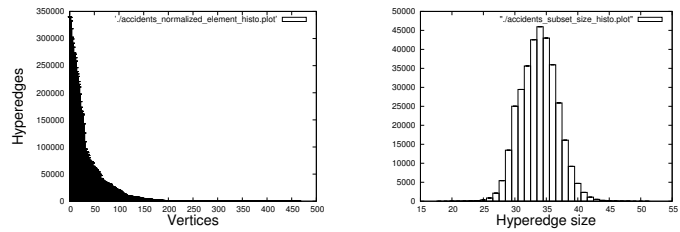


Figure 4: *accidents.dat* (468 vertices, 340183 hyperedges)

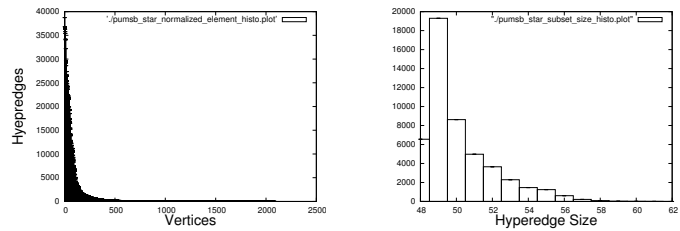


Figure 5: *pumsbstar.dat* (7116 vertices, 49046 hyperedges)

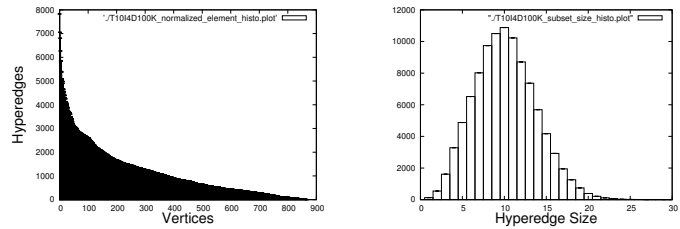


Figure 6: *T10I4D100K.dat* (999 vertices, 100000 hyperedges)

minimal transversals in a file: within a day, the program stopped because our 500 Go hard drive was full.

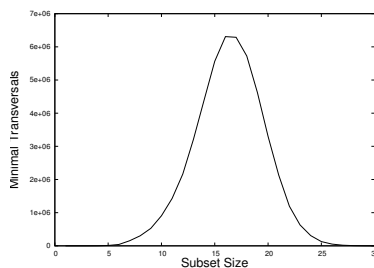


Figure 7: Minimal Transversals of *mushroom.dat*

## 5. Conclusion

The models we have studied already give a partial information on the average number of minimal transversals in real context and on the average complexity of the algorithms. Indeed, our models predict the order of growth of the size of minimal transversals. Hence, we are able to tell whether the computation can be made in reasonable time and space. Though, the upper bounds we have obtained on the number of minimal transversals still seems too large compared to real-data examples.

## References

- [1] Dimitris Achlioptas and Cristopher Moore. On the 2-colorability of random hypergraphs. In *Proc. 6th RANDOM, 7890*, pages 78–90. Springer-Verlag.
- [2] Peter Damaschke. Parameterized enumeration, transversals, and imperfect phylogeny reconstruction. *Theor. Comput. Sci.*, 351(3):337–350, 2006.
- [3] Guozhu Dong and Jinyan Li. Mining border descriptions of emerging patterns from dataset pairs. *Knowledge and Information Systems*, 8:178–202, 2005.
- [4] Andrzej Dudek and Alan Frieze. Loose hamilton cycles in random k-uniform hypergraphs, 2010.
- [5] Thomas Eiter and Georg Gottlob. Identifying the minimal transversals of a hypergraph and related problems. *SIAM J. Comput.*, 24(6):1278–1304, 1995.
- [6] Thomas Eiter and Georg Gottlob. Hypergraph transversal computation and related problems in logic and ai. In Sergio Flesca, Sergio Greco, Nicola Leone, and Giovambattista Ianni, editors, *JELIA*, volume 2424 of *Lecture Notes in Computer Science*, pages 549–564. Springer, 2002.

- 1  
2  
3  
4  
5  
6  
7  
8  
9 [7] Thomas Eiter, Georg Gottlob, and Kazuhisa Makino. New results on mono-  
10 tone dualization and generating hypergraph transversals. *SIAM J. Com-*  
11 *put.*, 32(2):514–537, 2003.
- 12 [8] Thomas Eiter, Kazuhisa Makino, and Georg Gottlob. Computational as-  
13 pects of monotone dualization: A brief survey. *Discrete Applied Mathematics*,  
14 156(11):2035–2049, 2008.
- 15 [9] P. Erdős and A. Rényi. On random graphs. I. *Publ. Math. Debrecen*,  
16 6:290–297, 1959.
- 17 [10] P. Erdős and A. Rényi. On the evolution of random graphs. In *Publication*  
18 *of the Mathematical Institute of the Hungarian Academy of Sciences*, pages  
19 17–61, 1960.
- 20 [11] Michael L. Fredman and Leonid Khachiyan. On the complexity of dualiza-  
21 tion of monotone disjunctive normal forms. *J. Algorithms*, 21(3):618–628,  
22 1996.
- 23 [12] Dimitrios Gunopulos, Roni Khardon, Heikki Mannila, and Hannu Toivo-  
24 nen. Data mining, hypergraph transversals, and machine learning. In  
25 *PODS*, pages 209–216. ACM Press, 1997.
- 26 [13] Matthias Hagen. *Algorithmic and Computational Complexity Issues*  
27 *of MONET*. Dissertation, Institut für Informatik, Friedrich-Schiller-  
28 Universität Jena, December 2008.
- 29 [14] Céline Hébert, Alain Bretto, and Bruno Crémilleux. A data mining formal-  
30 ization to improve hypergraph minimal transversal computation. *Fundam.*  
31 *Inf.*, 80:415–433, December 2007.
- 32 [15] Ilya Kapovich, Alexei Myasnikov, Paul Schupp, and Vladimir Shpilrain.  
33 Generic-case complexity, decision problems in group theory and random  
34 walks. *J. Algebra*, 264:665–694, 2003.
- 35 [16] Dimitris J. Kavvadias and Elias C. Stavropoulos. An efficient algorithm  
36 for the transversal hypergraph generation. *J. Graph Algorithms Appl.*,  
37 9(2):239–264, 2005.
- 38 [17] M. Lelarge. A new approach to the orientation of random hypergraphs.  
39 In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on*  
40 *Discrete Algorithms*, SODA '12, pages 251–264. SIAM, 2012.
- 41 [18] Keisuke Murakami and Takeaki Uno. Efficient algorithms for dualizing  
42 large-scale hypergraphs. *Discrete Applied Mathematics*, 70:83–94, 2014.
- 43 [19] Alexei G. Myasnikov and Alexander N. Rybalov. Generic complexity of  
44 undecidable problems. *J. Symbolic Logic*, 73(2):656–673, 2008.
- 45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4  
5  
6  
7  
8  
9 [20] Sarkar Saswati and Sivarajan Kumar N. Hypergraph models for cellular  
10 mobile communication systems. *IEEE Transactions on Vehicular Technol-*  
11 *ogy*, 47(2):460–471, 1998.  
12  
13 [21] Ilya Shmulevich, Aleksey D. Korshunov, and Jaakko Astola. Almost all  
14 monotone boolean functions are polynomially learnable using membership  
15 queries. *Inf. Process. Lett.*, 79(5):211–213, September 2001.  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65