

On Hashing-Based Approaches to Approximate DNF-Counting*

Kuldeep S. Meel¹, Aditya A. Shrotri², and Moshe Y. Vardi³

- 1 National University of Singapore, Singapore, Singapore
meel@comp.nus.edu.sg
- 2 Rice University, Houston, USA
as128@rice.edu
- 3 Rice University, Houston, USA
vardi@rice.edu

Abstract

Propositional model counting is a fundamental problem in artificial intelligence with a wide variety of applications, such as probabilistic inference, decision making under uncertainty, and probabilistic databases. Consequently, the problem is of theoretical as well as practical interest. When the constraints are expressed as DNF formulas, Monte Carlo-based techniques have been shown to provide a fully polynomial randomized approximation scheme (FPRAS). For CNF constraints, hashing-based approximation techniques have been demonstrated to be highly successful. Furthermore, it was shown that hashing-based techniques also yield an FPRAS for DNF counting without usage of Monte Carlo sampling. Our analysis, however, shows that the proposed hashing-based approach to DNF counting provides poor time complexity compared to the Monte Carlo-based DNF counting techniques. Given the success of hashing-based techniques for CNF constraints, it is natural to ask: Can hashing-based techniques provide an efficient FPRAS for DNF counting? In this paper, we provide a positive answer to this question. To this end, we introduce two novel algorithmic techniques: *Symbolic Hashing* and *Stochastic Cell Counting*, along with a new hash family of *Row-Echelon hash functions*. These innovations allow us to design a hashing-based FPRAS for DNF counting of similar complexity (up to polylog factors) as that of prior works. Furthermore, we expect these techniques to have potential applications beyond DNF counting.

1998 ACM Subject Classification G.1.2 Special Function Approximation, F.4.1 Logic and Constraint Programming

Keywords and phrases Model Counting, Approximation, DNF, Hash Functions

Digital Object Identifier 10.4230/LIPIcs...

1 Introduction

Propositional model counting is a fundamental problem in artificial intelligence with a wide range of applications including probabilistic inference, databases, decision making under uncertainty, and the like [1, 6, 7, 24]. Given a Boolean formula ϕ , the problem of propositional model counting, also referred to as #SAT, is to compute the number of solutions of ϕ [29]. Depending on whether ϕ is expressed as a CNF or DNF formula, the corresponding model counting problems are denoted as #CNF or #DNF, respectively. Both #CNF and #DNF

* The author list has been sorted alphabetically by last name; this should not be used to determine the extent of authors' contributions



have a wide variety of applications. For example, probabilistic-inference queries reduce to solving #CNF instances [1, 21, 22, 24], while evaluation of queries for probabilistic database reduce to #DNF instances [6]. Consequently, both #CNF and #DNF have been of theoretical as well as practical interest over the years [16, 18, 23, 25]. In his seminal paper, Valiant [29] showed that both #CNF and #DNF are #P-complete, a class of problems that are believed to be intractable in general.

Given the intractability of #CNF and #DNF, much of the interest lies in the approximate variants of #CNF and #DNF, wherein for given tolerance and confidence parameters ε and δ , the goal is to compute an estimate C such that C is within a $(1 + \varepsilon)$ multiplicative factor of the true count with confidence at least $1 - \delta$. While both #CNF and #DNF are #P-complete in their exact forms, the approximate variants differ in complexity: approximating #DNF can be accomplished in fully polynomial randomized time [5, 17, 18], but approximate #CNF is NP-hard [25]. Consequently, different techniques have emerged to design scalable approximation techniques for #DNF and #CNF.

In the context of #DNF, the works of Karp, Luby, and Madras [17, 18] led to the development of highly efficient Monte-Carlo based techniques, whose time complexity is linear in the size of the formula. On the other hand, hashing-based techniques have emerged as a scalable approach to the approximate model counting of CNF formulas [3, 4, 10, 13, 25], and are effective even for problems with existing FPRAS such as *network reliability* [8]. These hashing-based techniques employ 2-universal hash functions to partition the space of satisfying solutions of a CNF formula into cells such that a randomly chosen cell contains only a small number of solutions. Furthermore, it is shown that the number of solutions across the cells is *roughly equal* and, therefore, an estimate of the total count can be obtained by counting the number of solutions in a cell and scaling the obtained count by the number of cells. Since the problem of counting the number of solutions in a cell when the number of solutions is *small* can be accomplished efficiently by invoking a SAT solver, the hashing-based techniques can take advantage of the recent progress in the development of efficient SAT solvers. Consequently, algorithms such as ApproxMC [3, 4] have been shown to scale to instances with hundreds of thousands of variables.

While Monte Carlo techniques introduced in the works of Karp et al. have shown to not be applicable in the context of approximate #CNF [18], it was not known whether hashing-based techniques could be employed to obtain efficient algorithms for #DNF. Recently, significant progress in this direction was achieved by Chakraborty, Meel and Vardi [4], who showed that hashing-based framework of ApproxMC could be employed to obtain FPRAS for #DNF counting¹. There is, however, no precise complexity analysis in [4]. In this paper, we provide a complexity analysis of the proposed scheme of Chakraborty et al., which is worse than quartic in the size of formula. In comparison, state-of-the-art approaches achieve complexity linear in the number of variables and cubes for #DNF counting. This begs the question: *How powerful is the hashing-based framework in the context of DNF counting? In particular, can it lead to algorithms competitive in runtime complexity with state-of-the-art?*

In this paper, we provide a positive answer to this question. To achieve such a significant reduction in complexity, we offer three novel algorithmic techniques: (i) A new class of 2-universal hash functions that enable fast enumeration of solutions using Gray Codes, (ii) Symbolic Hashing, and (iii) Stochastic Cell Counting. These techniques allow us to achieve the complexity of $\tilde{O}(mn \log(1/\delta)/\varepsilon^2)$, which is within polylog factors of the complexity achieved

¹ It is worth noting that several hashing-based algorithms based on [10, 28] do not lead to FPRAS schemes for #DNF despite close similarity to Chakraborty et al.'s approach

by Karp et al. [18]. Here, m and n are the number of cubes and variables respectively while ε and δ are the tolerance and confidence of approximation. Furthermore, we believe that these techniques are not restricted to $\#DNF$. Given recent breakthroughs achieved in the development of hashing-based CNF-counting techniques, we believe our techniques have the potential for a wide variety of applications.

The rest of the paper is organized as follows: we introduce notation in section 2 and discuss related work in section 3. We describe our main contributions in section 4, analyze the resulting algorithm in section 5 and discuss future work and conclude in section 6.

2 Preliminaries

DNF Formulas and Counting

We use Greek letters ϕ , θ and ψ to denote boolean formulas. A formula ϕ over boolean variables x_1, x_2, \dots, x_n is in Disjunctive Normal Form (DNF) if it is a disjunction over conjunctions of variables or their negations. We use X to denote the set of variables appearing in the formula. Each occurrence of a variable or its negation is called a *literal*. Disjuncts in the formula are called *cubes* and we denote the i^{th} cube by ϕ^{C_i} . Thus $\phi = \phi^{C_1} \vee \phi^{C_2} \vee \dots \vee \phi^{C_m}$ where each ϕ^{C_i} is a conjunction of *literals*. We will use n and m to denote the number of variables and number of cubes in the input DNF formula, respectively. The number of literals in a cube ϕ^{C_i} is called its *width* and is denoted by $\text{width}[\phi^{C_i}]$.

An assignment to all the variables can be represented by a vector $\mathbf{x} \in \{0, 1\}^n$ with 1 corresponding to *true* and 0 to *false*. $U = \{0, 1\}^n$ is the set of all possible assignments, which we refer to as the *universe* or state space interchangeably. An assignment \mathbf{x} is called a satisfying assignment for a formula ϕ if ϕ evaluates to *true* under \mathbf{x} . In other words \mathbf{x} satisfies ϕ and is denoted as $\mathbf{x} \models \phi$. Note that an assignment \mathbf{x} will satisfy a DNF formula ϕ if $\mathbf{x} \models \phi^{C_i}$ for some i . The DNF-Counting Problem is to count the number of satisfying assignments of a DNF formula.

Next, we formalize the concept of a counting problem. Let $R \subseteq \{0, 1\}^* \times \{0, 1\}^*$ be a relation which is decidable in polynomial time and there is a polynomial p such that for every $(s, t) \in R$ we have $|t| \leq p(|s|)$. The decision problem corresponding to R asks if for a given s there exists a t such that $(s, t) \in R$. Such a problem is in NP. Here, s is called the *problem instance* and t is called the *witness*. We denote the set of all witnesses for a given s by R_s . The *counting problem* corresponding to R is to calculate the size of the witness set $|R_s|$ for a given s . Such a problem is in $\#P$ [29]. The DNF-Counting problem is an example of this formalism: A formula ϕ is a problem instance and a satisfying assignment \mathbf{x} is a witness of ϕ . The set of satisfying assignments or the solution space is denoted R_ϕ and the goal is to compute $|R_\phi|$. It is known that the problem is $\#P$ -Complete, which is believed to be intractable [27]. Therefore, we look at what it means to efficiently and accurately approximate this problem.

A *fully polynomial randomized approximation scheme* (FPRAS) is a randomized algorithm that takes as input a problem instance s , a tolerance $\varepsilon \in (0, 1)$ and confidence parameter $\delta \in (0, 1)$ and outputs a random variable c such that $\Pr[\frac{1}{1+\varepsilon}|R_s| \leq c \leq (1+\varepsilon)|R_s|] \geq 1 - \delta$ and the running time of the algorithm is polynomial in $|s|, 1/\varepsilon, \log(1/\delta)$ [17]. Notably, while exact DNF-counting is inter-reducible with exact CNF-counting, the approximate versions of the two problems are not because multiplicative approximation is not closed under complementation.

Matrix Notation

We use x, y, z, \dots to denote scalar variables. We use subscripts x_1, x_2, \dots as required. In this paper we are dealing with operations over the boolean ring, where the variables are boolean, 'addition' is the XOR operation (\oplus) and 'multiplication' is the AND operation (\wedge). We use the letters i, j, k, l as indices or to denote positions. We denote sets by non-boldface capital letters. We use capital boldface letters $\mathbf{A}, \mathbf{B}, \dots$ to denote matrices, small boldface letters $\mathbf{u}, \mathbf{v}, \mathbf{w}, \dots$ to denote vectors. $\mathbf{A}^{[p \times q]}$ denotes a matrix of p rows and q columns, while $\mathbf{u}^{[q]}$ denotes a vector of length q . $\mathbf{0}^{[q]}$ and $\mathbf{1}^{[q]}$ are the all 0s and all 1s vectors of length q , respectively. We omit the dimensions when clear from context. $\mathbf{x}[i]$ denotes the i^{th} element of \mathbf{x} , while $\mathbf{A}[i, j]$ denotes the element in the i^{th} row and j^{th} column of \mathbf{A} . $\mathbf{A}[r_1 : r_2, c_1 : c_2]$ denotes the sub-matrix of \mathbf{A} between rows r_1 and r_2 excluding r_2 and columns c_1 and c_2 excluding c_2 . Similarly $\mathbf{v}[i : j]$ denotes the sub-vector of \mathbf{v} between index i and index j excluding j . The i^{th} row of \mathbf{A} is denoted $\mathbf{A}[i, :]$ and j^{th} column as $\mathbf{A}[:, j]$. The $p \times (q_1 + q_2)$ matrix formed by concatenating rows of matrices $\mathbf{A}^{[p \times q_1]}$ and $\mathbf{B}^{[p \times q_2]}$ is written in block notation as $[\mathbf{A} \mid \mathbf{B}]$, while $[\frac{\mathbf{A}}{\mathbf{B}}]$ represents concatenation of columns. Similarly the $(q_1 + q_2)$ -length concatenation of vectors $\mathbf{v}^{[q_1]}$ and $\mathbf{w}^{[q_2]}$ is $[\mathbf{v} \mid \mathbf{w}]$. The dot product between matrix \mathbf{A} and vector \mathbf{x} is written as $\mathbf{A} \cdot \mathbf{x}$. The vector formed by element-wise XOR of vectors \mathbf{v} and \mathbf{w} is denoted $\mathbf{v} \oplus \mathbf{w}$.

Hash Functions

A hash function $h : \{0, 1\}^q \rightarrow \{0, 1\}^p$ partitions the elements of the domain $\{0, 1\}^q$ into 2^p cells. $h(\mathbf{x}) = \mathbf{y}$ implies that h maps the assignment \mathbf{x} to the cell \mathbf{y} . $h^{-1}(\mathbf{y}) = \{\mathbf{x} \mid h(\mathbf{x}) = \mathbf{y}\}$ is the set of assignments that map to the cell \mathbf{y} . In the context of counting, 2-universal families of hash functions, denoted by $H(q, p, 2)$, are of particular importance. When h is sampled uniformly at random from $H(q, p, 2)$, 2-universality entails

1. $\Pr[h(\mathbf{x}_1) = h(\mathbf{x}_2)] \leq 2^{-p}$ for all $\mathbf{x}_1 \neq \mathbf{x}_2$
2. $\Pr[h(\mathbf{x}) = \mathbf{y}] = 2^{-p}$ for every $\mathbf{x} \in \{0, 1\}^q$ and $\mathbf{y} \in \{0, 1\}^p$.

Of particular interest is the random XOR family of hash functions, which is defined as $H_{XOR}(q, p) = \{\mathbf{A} \cdot \mathbf{x} \oplus \mathbf{b} \mid \mathbf{A}[i, j] \in \{0, 1\} \text{ and } \mathbf{b}[i] \in \{0, 1\}; 0 \leq i < p, 0 \leq j < q\}$. Selecting $\mathbf{A}[i, j]$ s and $\mathbf{b}[i]$ s randomly from $\{0, 1\}$ is equivalent to drawing uniformly at random from this family. A pair \mathbf{A} and \mathbf{b} now defines a hash function $h_{\mathbf{A}, \mathbf{b}}$ as follows: $h_{\mathbf{A}, \mathbf{b}}(\mathbf{x}) = \mathbf{A} \cdot \mathbf{x} \oplus \mathbf{b}$. This family was shown to be 2-universal in [2]. For a hash function $h \in H_{XOR}(q, p)$, we have that $h(\mathbf{x}) = \mathbf{y}$ is a system of linear equations modulo 2: $\mathbf{A} \cdot \mathbf{x} \oplus \mathbf{b} = \mathbf{y}$. From another perspective, it can be viewed as a boolean formula $\psi = \bigwedge_{i=1}^p (\bigoplus_{j=1}^q (\mathbf{A}[i, j] \wedge \mathbf{x}[j])) \oplus \mathbf{b}[i] = \mathbf{y}[i]$. The solutions to this formula are exactly the elements of the set $h^{-1}(\mathbf{y})$.

Gaussian Elimination

Solving a system of linear equations over q variables and p constraints can be done by row reduction technique known variously as *Gaussian Elimination* or *Gauss-Jordan Elimination*. A matrix is in *Row-Echelon form* if rows with at least one nonzero element are above any rows of all zeros. The matrix is in *Reduced Row-Echelon form* if, in addition, every leading non-zero element in a row is 1 and is the only nonzero entry in its column. We refer to the technique for obtaining the *Reduced Row-Echelon form* of a matrix as Gaussian Elimination. We refer the reader to any standard text on linear algebra (cf., [26]) for details. For a matrix in Reduced Row-Echelon form, the row-rank is simply the number of non-zero rows.

For a system of linear equations $\mathbf{A}\mathbf{x} \oplus \mathbf{b} = \mathbf{y}$, if the row-rank of the augmented matrix is same as row-rank of \mathbf{A} , then the system is consistent and the number of solutions is $2^{q - \text{rowrank}(\mathbf{A})}$ where q is the number of variables in the system of equations. Moreover, if \mathbf{A} is in Reduced Row-Echelon form, then the values of the variables corresponding to leading 1s in each row are completely determined by the values assigned to the remaining variables. The variables corresponding to the leading 1's are called *dependent* variables and the remaining variables are *free*. Let X_F and $X \setminus X_F$ denote the set of free and dependent variables respectively. Let $f = |X_F|$. Clearly $f = q - \text{rowrank}(\mathbf{A})$. For each possible assignment to the free variables we get an assignment to the dependent variables by propagating the values through the augmented matrix in $\mathcal{O}(q^2)$ time. Thus we can enumerate all 2^f satisfying assignments to a system of linear equations $\mathbf{A}\mathbf{x} \oplus \mathbf{b} = \mathbf{y}$ if \mathbf{A} in Reduced Row-Echelon form.

Gray Codes

A Gray code [14] is an ordering of 2^l binary numbers for some $l \geq 1$ with the property that every pair of consecutive numbers in the sequence differ in exactly one bit. Thus starting from $\mathbf{0}^l$ we can iteratively construct the entire Gray code sequence by flipping one bit in each step. We assume access to a procedure *nextGrayBit* that in each call returns the position of the next bit that is to be flipped. Such a procedure can be implemented in constant time by a trivial modification of Algorithm L in [19].

3 Related Work

Propositional model counting has been of theoretical as well as practical interest over the years [16, 17, 23, 27]. Early investigations showed that both #CNF and #DNF are #P-complete [29]. Consequently, approximation algorithms have been explored for both problems. A major breakthrough for approximate #DNF was achieved by the seminal work of Karp and Luby [17], which provided a Monte Carlo-based FPRAS for #DNF. The proposed FPRAS was improved by follow-up work of Karp, Luby and Madras [18] and Dagum et al. [5], achieving the best known complexity of $\mathcal{O}(mn \log(1/\delta)/\varepsilon^2)$. In this work, we bring certain ideas of Karp et al. into the hashing framework with significant adaptations.

For #CNF, early work on approximate counting resulted in hashing-based schemes that required polynomially many calls to an NP-oracle [25, 28]. No practical algorithms materialized from these schemes due to the impracticality of the underlying NP queries. Subsequent attempts to circumvent hardness led to the development of several hashing and sampling-based approaches that achieved scalability but provided very weak or no guarantees [13, 11]. Due to recent breakthroughs in the design of hashing-based techniques, several tools have been developed recently that can handle formulas involving hundreds of thousands of variables while providing rigorous formal guarantees. Overall, these tools can be broadly classified by their underlying hashing-based technique as: (i) obtain a constant factor approximation and then use identical copies of the input formula to obtain ε approximations [10], or (ii) directly obtain ε guarantees [3, 4]. The first technique when applied to DNF formulas is not an FPRAS. In contrast, Chakraborty, Meel and Vardi [4] recently showed that tools based on the latter approach, such as **ApproxMC2**, do provide FPRAS for #DNF counting. Chakraborty et al. did not analyze the complexity of the algorithm in their work. We now provide a precise complexity analysis of **ApproxMC2** for #DNF. To that end, we first describe the **ApproxMC** framework on which **ApproxMC2** is built.

3.1 ApproxMC Framework

Chakraborty et al. introduced in [3] a hashing-based framework called ApproxMC that requires linear (in n) number of SAT calls. Subsequently in ApproxMC2, the number of SAT calls was reduced from linear to logarithmic (in n). The core idea of ApproxMC is to employ 2-universal hash functions to partition the solution space into *roughly equal small* cells, wherein a cell is called *small* if it has less than or equal to $hiThresh$ solutions, such that $hiThresh$ is a function of ε . A SAT solver is employed to check if a cell is small by enumerating solutions one-by-one until either there are no more solutions or we have already enumerated $hiThresh + 1$ solutions. Following the terminology of [3], we refer to the above described procedure as BSAT (bounded SAT). To determine the number of cells, ApproxMC performs a search that requires $\mathcal{O}(\log n)$ steps and the estimate is returned as the count of the solutions in a randomly picked small cell scaled by the total number of cells. To amplify confidence to the desired levels of $1 - \delta$, ApproxMC invokes the estimation routine $\mathcal{O}(\log \frac{1}{\delta})$ times and reports the median of all such estimates. Hence, the number of BSAT invocations is $\mathcal{O}(\log n \log(\frac{1}{\delta}))$.

FPRAS for #DNF

The key insight of Chakraborty et al. [4] is that the BSAT procedure can be done in polynomial time when the input formula to ApproxMC is in DNF. In particular, the input to every invocation of BSAT is a formula that is a conjunction of the input DNF formula and a set of XOR constraints derived from the hash function. Chakraborty et al. observed that one can iterate over all the cubes of the input formula, substitute each cube into the set of XOR constraints separately, and employ Gaussian Elimination to enumerate the solutions of the simplified XOR constraints. Note that at no step would one have to enumerate more than $hiThresh$ solutions. Since Gaussian Elimination is a polynomial-time procedure, BSAT can be accomplished in polynomial time as well. Chakraborty et al. did not provide a precise complexity analysis of BSAT. We now provide such an analysis. We defer all proofs to the appendix. The following lemma states the time complexity of the BSAT routine.

► **Lemma 1.** *The complexity of BSAT when the input formula to ApproxMC2 is in DNF is $\mathcal{O}(mn^3 + mn^2/\varepsilon^2)$.* ◀

We can now complete the complexity analysis:

► **Lemma 2.** *The complexity of ApproxMC2 is $\mathcal{O}((mn^3 + mn^2/\varepsilon^2) \log n \log(1/\delta))$ when the input formula is in DNF.* ◀

4 Efficient Hashing-based DNF Counter

We now present three key novel algorithmic innovations that allow us to design hashing-based FPRAS for #DNF with complexity similar to Monte Carlo-based state-of-the-art techniques. We first introduce a new family of 2-universal hash functions that allow us to circumvent the need for expensive Gaussian Elimination. We then discuss the concept of Symbolic Hashing, which allows us to design hash functions over a space different than the assignment space, allowing us to achieve significant reduction in the complexity of search procedure for the number of the cells. Finally, we show that BSAT can be replaced by an efficient stochastic estimator. These three techniques allow us to achieve significant reduction in the complexity of hashing-based DNF counter without loss of theoretical guarantees.

Algorithm 1 enumNextREX($\mathbf{D}, \mathbf{u}, \mathbf{v}, k$)

```

1:  $\mathbf{v}' \leftarrow \mathbf{v}$ ;
2:  $\mathbf{v}'[k] \leftarrow \neg \mathbf{v}[k]$ ;
3:  $\mathbf{u}' \leftarrow \mathbf{u} \oplus \mathbf{D}[:, k]$ ;
4: return  $(\mathbf{u}', \mathbf{v}')$ 

```

4.1 Row-Echelon XOR Hash Functions

The complexity analysis presented in Section 3 shows that the expensive Gaussian Elimination contributes significantly to poor time complexity of ApproxMC2. Since the need for Gaussian Elimination originates from the usage of H_{XOR} , we seek a family of 2-universal hash functions that circumvents this need. We now introduce a Row-Echelon XOR family of hash functions defined as $H_{REX}(\mathbf{q}, \mathbf{p}) = \{\mathbf{A} \cdot \mathbf{x} \oplus \mathbf{b} \mid \mathbf{A}^{[p \times q]} = [\mathbf{I}^{[p \times p]} \ : \ \mathbf{D}^{[p \times (q-p)]]}\}$ where \mathbf{I} is the identity matrix, \mathbf{D} and \mathbf{b} are random 0/1 matrix and vector respectively. In particular, we ensure that for every $\mathbf{D}[i, j]$ and $\mathbf{b}[i]$ we have $\Pr[\mathbf{D}[i, j] = 1] = \Pr[\mathbf{D}[i, j] = 0] = 0.5$ and also $\Pr[\mathbf{b}[i] = 1] = \Pr[\mathbf{b}[i] = 0] = 0.5$. Note that \mathbf{D} and \mathbf{b} completely define a hash function from H_{REX} . The following theorem establishes the desired properties of universality for H_{REX} . The proof is deferred to Appendix.

► **Theorem 3.** H_{REX} is 2-universal. ◀

The naive way of enumerating satisfying assignments for a given $\mathbf{D}^{[p \times (q-p)]}$, $\mathbf{b}^{[p]}$, and $\mathbf{y}^{[p]}$ is to iterate over all 2^f assignments to the free variables in sequence starting from $\mathbf{0}^{[f]}$ to $\mathbf{1}^{[f]}$, where $f = (q - p)$. For each assignment $\mathbf{v}^{[f]}$ to the free variables, the corresponding assignment to the dependent variables $\mathbf{u}^{[q-f]}$ can be calculated as $\mathbf{u} = (\mathbf{D} \cdot \mathbf{v}) \oplus \mathbf{b} \oplus \mathbf{y}$, which requires $O(pq)$ time. Can we do better?

We answer the above question positively by iterating over the 2^f assignments to the free variables out of sequence. In particular, we iterate using the Gray code sequence for f bits. The procedure is outlined in enumNextREX (Algorithm 1). The algorithm takes the hash matrix \mathbf{D} , an assignment to the free variables \mathbf{v} , and an assignment to the dependent variables \mathbf{u} as inputs, and outputs the next free-variable assignment \mathbf{v}' in the Gray sequence and the corresponding assignment \mathbf{u}' to the dependent variables. k represents the position of the bit that is changed between \mathbf{v} and \mathbf{v}' . Thus enumNextREX constructs a satisfying assignment to a Row-Echelon XOR hash function in each invocation in $\mathcal{O}(q)$ time.

4.2 Symbolic Hashing

For DNF formulas, R_ϕ can be exponentially sparse compared to U , which is undesirable². It is possible, however, to transform U to another space U' and the solution space R_ϕ to R'_ϕ such that the ratio $|U'|/|R'_\phi|$ is polynomially bounded and $|R_\phi| = |R'_\phi|$. For DNF formulas, the new universe U' is defined as $U' = \{(\mathbf{x}, \phi^{C^i}) \mid \mathbf{x} \models \phi^{C^i}\}$. Thus, corresponding to each $\mathbf{x} \models \phi$ that satisfies cubes $\phi^{C^{i_1}}, \dots, \phi^{C^{i_\infty}}$ in ϕ , we have the states $\{(\mathbf{x}, \phi^{C^{i_1}}), (\mathbf{x}, \phi^{C^{i_2}}), \dots, (\mathbf{x}, \phi^{C^{i_\infty}})\}$ in U' . Next, the solution space is defined as $R'_\phi = \{(\mathbf{x}, \phi^{C^i}) \mid \mathbf{x} \models \phi^{C^i} \text{ and } \forall j < i, \mathbf{x} \not\models \phi^{C^j}\}$ for a fixed ordering of the cubes. The definition of R'_ϕ ensures that $|R_\phi| = |R'_\phi|$. This transformation is due to Karp and Luby [17].

The key idea of Symbolic Hashing is to perform 2-universal hashing symbolically over the transformed space. In particular, the sampled hash function partitions the space U' instead

² Number of steps of ApproxMC2 search procedure increases with sparsity

Algorithm 2 SymbolicDNFAproxMCCore(ϕ , hiThresh)

```

1:  $w \leftarrow$  width of cubes;
2:  $q \leftarrow n - w + \log m$ ;
3:  $sI \leftarrow n - w - \log \text{hiThresh}$ ;
4:  $(\hat{\mathbf{D}}^{[(q-1) \times (q-sI)]}, \hat{\mathbf{b}}, \hat{\mathbf{y}}) \leftarrow \text{SampleBase}(q, sI)$ ;
5:  $\mathbf{p} \leftarrow \text{LogSATSearch}(\phi, \hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}, \text{hiThresh}, sI, q - 1)$ ;
6:  $\text{solCount} \leftarrow \text{BSAT}(\phi, \hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}, \text{hiThresh}, \mathbf{p}, q, sI)$ ;
7: return  $(2^{\mathbf{p}}, \text{solCount})$ 

```

of U . Therefore, we employ hash functions from $H_{\text{REX}}(q, \mathbf{p})$ over $q = n - w + \log m$ variables instead of n variables. Note that the variables of a satisfying assignment $\mathbf{z} \in \{0, 1\}^q$ to the hash function are now different from the variables to a satisfying assignment $\mathbf{x} \in \{0, 1\}^n$ of the input formula ϕ . We interpret \mathbf{z} as follows: the last $\log m$ bits of \mathbf{z} are converted to a number i such that $1 \leq i \leq m$. Now ϕ^{C^i} corresponds to a partial assignment of $\text{width}[\phi^{C^i}]$ variables in that cube. For simplicity, we assume that each cube is of the same width w .³ The remaining $n - w$ bits of \mathbf{z} are interpreted to be the assignment to the $n - w$ variables not in ϕ^{C^i} giving a complete assignment \mathbf{x} . Thus we get a pair (\mathbf{x}, ϕ^{C^i}) from \mathbf{z} such that $\mathbf{x} \models \phi^{C^i}$. For a fixed ordering of variables and cubes we see that there is a bijection between (\mathbf{x}, ϕ^{C^i}) and \mathbf{z} and hence the 2-universality guarantee holds over the partitioned space of U' .

4.3 Stochastic Cell-Counting

To estimate the number of solutions in a cell, we need to check for every tuple (\mathbf{x}, ϕ^{C^i}) generated using symbolic hash function as described above: if $(\mathbf{x}, \phi^{C^i}) \in R'_\phi$. Such a check would require iteration over cubes ϕ^{C^j} for $1 \leq j \leq (i - 1)$ and returning no if $\mathbf{x} \models \phi^{C^j}$ for some j and yes otherwise. This would result in procedure with $O(mn)$ complexity.

Our key observation is that a precise count of the number of solutions in a cell is not required and therefore, one can employ a stochastic estimator for the number of solutions in a cell. We proceed as follows: we define the *coverage* of an assignment \mathbf{x} as $\text{cov}(\mathbf{x}) = \{j \mid \mathbf{x} \models \phi^{C^j}\}$. Note that $\sum_{(\mathbf{x}, \phi^{C^i}) \in U'} \frac{1}{|\text{cov}(\mathbf{x})|} = |R_\phi|$.

We define a random variable $\mathbf{c}_\mathbf{x}$ as the number of steps taken to uniformly and independently sample from $\{1, 2, \dots, m\}$, a number j such that $\mathbf{x} \models \phi^{C^j}$. For a randomly chosen j , the probability $\Pr[\mathbf{x} \models \phi^{C^j}] = |\text{cov}(\mathbf{x})|/m$, which follows the Bernoulli distribution. The random variable $\mathbf{c}_\mathbf{x}$ is the number of Bernoulli trials for the first success, which follows the geometric distribution. Therefore, $\mathbf{E}[\mathbf{c}_\mathbf{x}] = m/|\text{cov}(\mathbf{x})|$, and $\mathbf{E}[\mathbf{c}_\mathbf{x}/m] = 1/|\text{cov}(\mathbf{x})|$. The estimator $\mathbf{c}_\mathbf{x}/m$ has been previously employed by Karp et al. [18]. Here, we show that it can also be used for Stochastic Cell-Counting: we define the estimator for the number of solutions in a cell as $\Omega_\mathbf{y} = \sum_{(\mathbf{x}, \phi^{C^i}) \in h^{-1}(\mathbf{y})} \mathbf{c}_\mathbf{x}/m$.

4.4 The Full Algorithm

We now incorporate the above techniques into ApproxMC2 and call the revised algorithm SymbolicDNFAproxMC, which is presented as Algorithm 3. First, note that expression for hiThresh is twice that for ApproxMC2. Then, in line 4, a matrix $\hat{\mathbf{D}}$ and vectors $\hat{\mathbf{b}}$ and

³ We can handle non-uniform width cubes by sampling ϕ^{C^i} with probability $\frac{2^{n-\text{width}[\phi^{C^i}]}}{\sum_{j=1}^m 2^{n-\text{width}[\phi^{C^j}]}}$ instead of uniformly

Algorithm 3 SymbolicDNFAproxMC($\phi, \varepsilon, \delta$)

```

1: hiThresh  $\leftarrow 2 * (1 + 9.84 \left(1 + \frac{\varepsilon}{1+\varepsilon}\right) \left(1 + \frac{1}{\varepsilon}\right)^2)$ ;
2:  $t \leftarrow \lceil 17 \log_2(3/\delta) \rceil$ ;
3: EstimateList  $\leftarrow$  emptyList; iter  $\leftarrow$  0;
4: repeat
5:   iter  $\leftarrow$  iter + 1;
6:   (cellCount, solCount)  $\leftarrow$  SymbolicDNFAproxMCCore( $\phi$ , hiThresh);
7:   if (cellCount  $\neq \perp$ ) then AddToList(EstimateList, solCount  $\times$  cellCount);
8: until (iter  $\geq t$ );
9: finalEstimate  $\leftarrow$  FindMedian(EstimateList);
10: return finalEstimate

```

Algorithm 4 SampleBase(\mathbf{q}, sI)

```

1: Sample  $\mathbf{G}$  uniformly from  $\{0, 1\}^{\lceil sI \times (\mathbf{q} - sI) \rceil}$ ;
2: Sample uniformly an upper triangular matrix  $\mathbf{E}^{\lceil (\mathbf{q} - sI - 1) \times (\mathbf{q} - sI) \rceil}$  with  $\mathbf{E}[i, i] = 1$  for all  $i$ .
3:  $\hat{\mathbf{D}} \leftarrow \begin{bmatrix} \mathbf{G} \\ \mathbf{E} \end{bmatrix}$ ;
4: Sample  $\hat{\mathbf{b}}$  and  $\hat{\mathbf{y}}$  uniformly from  $\{0, 1\}^{\mathbf{q}-1}$ ;
5: return  $\hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}$ 

```

$\hat{\mathbf{y}}$ are obtained, which are employed to construct an appropriate hash function and cell during the search procedure of SymbolicDNFAproxMCCore. SymbolicDNFAproxMC makes $t = O(\log(1/\delta))$ calls to SymbolicDNFAproxMCCore (line 4-8) and returns median of all the estimates (lines 9-10) to boost the probability of success to $1 - \delta$.

We now discuss the subroutine SymbolicDNFAproxMCCore, which is an adaptation of ApproxMC2Core but with significant differences. First, for DNF formulas with cube width w , the number of solutions is lower bounded by 2^{n-w} . Therefore, instead of starting with 1 hash constraint, we can safely start with $sI = n - w - \log \text{hiThresh}$ constraints (lines 3-4). Thereafter, SymbolicDNFAproxMCCore calls LogSATSearch in line 5 to find the right number \mathbf{p} of constraints. The cell count with \mathbf{p} constraints is calculated in line 6 and the estimate $(2^{\mathbf{p}}, \text{solCount})$ is returned in line 7.

SampleBase algorithm constructs the base matrix $\hat{\mathbf{D}}$ and base vectors $\hat{\mathbf{b}}$ and $\hat{\mathbf{y}}$ required for sampling from H_{REX} family. \mathbf{G} is a random matrix of dimension $sI \times (\mathbf{q} - sI)$ and \mathbf{E} is a random upper triangular matrix of dimension $(\mathbf{q} - sI - 1) \times (\mathbf{q} - sI)$ with all diagonal elements 1. In line 3, $\hat{\mathbf{D}}$ is constructed as the vertical concatenation $\begin{bmatrix} \mathbf{G} \\ \mathbf{E} \end{bmatrix}$.

LogSATSearch (algorithm 5) performs a binary search to find the number of constraints \mathbf{p} at which the cell count falls below hiThresh. For DNF formula with cube width of w , since the number of solutions is bounded between 2^{n-w} and $\mathbf{m} * 2^{n-w}$, we need to perform search for \mathbf{p} between $n - w$ and $n - w + \log \mathbf{m}$. Therefore, binary search can take at most $O(\log \log \mathbf{m})$ steps to find correct \mathbf{p} .

Symbolic Hashing is implemented in Algorithm 6 (BSAT). In line 2, we obtain a hash function from $H_{REX}(\mathbf{q}, \mathbf{p})$ over $\mathbf{q} = n - w + \log \mathbf{m}$ variables by calling Extract. We assume access to a procedure *nextGrayBit* in line 10 that returns the position of the bit that is flipped between two consecutive assignments. A satisfying assignment \mathbf{z} to the hash function is constructed in line 6. \mathbf{z} is interpreted to generate a pair (\mathbf{x}, ϕ^{C^i}) in line 7 which is checked for satisfiability in line 8. The final cell count is returned in line 12.

In CheckSAT (algorithm 7), we implement the stochastic cell counting procedure. The

Algorithm 5 $\text{LogSATSearch}(\phi, \hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}, \text{hiThresh}, \text{low}, \text{hi})$

```

1: lowerFib  $\leftarrow$  0; upperFib  $\leftarrow$   $hi - low + 1$ ;  $p \leftarrow low$ ;
2: FailRecord[0]  $\leftarrow$  1; FailRecord[ $hi - low + 1$ ]  $\leftarrow$  0;
3: FailRecord[ $i$ ]  $\leftarrow$   $\perp$  for all  $i$  other than 0 and  $hi - low + 1$ ;
4: while true do
5:    $C_{BSAT} \leftarrow \text{BSAT}(\phi, \hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}, \text{hiThresh}, p, q, sI)$ ;
6:   if ( $C_{BSAT} \geq \text{hiThresh}$ ) then
7:     if (FailRecord[ $p + 1 - low + 1$ ] = 0) then return  $p + 1$ ;
8:     FailRecord[ $i$ ]  $\leftarrow$  1 for all  $i \in \{1, \dots, p - low + 1\}$ ;
9:     lowerFib  $\leftarrow$   $p - low + 1$ ;
10:     $p \leftarrow (\text{upperFib} + \text{lowerFib})/2$ ;
11:   else
12:     if (FailRecord[ $p - 1 - low + 1$ ] = 1) then return  $p$ ;
13:     FailRecord[ $i$ ]  $\leftarrow$  0 for all  $i \in \{p, \dots, hi - low + 1\}$ ;
14:     upperFib  $\leftarrow$   $p - low + 1$ ;
15:      $p \leftarrow (\text{upperFib} + \text{lowerFib})/2$ ;

```

Algorithm 6 $\text{BSAT}(\phi, \hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}, \text{hiThresh}, p, q, sI)$

```

1:  $count \leftarrow$  0;
2:  $\mathbf{D}, \mathbf{b}, \mathbf{y} \leftarrow \text{Extract}(\hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}, p, q, sI)$ ;
3:  $\mathbf{u}^p \leftarrow \mathbf{b} \oplus \mathbf{y}$ ;
4:  $\mathbf{v}^{q-p} \leftarrow \mathbf{0}^{q-p}$ ;
5: for ( $j = 0$ ;  $j < 2^{q-p}$ ;  $j++$ ) do
6:    $\mathbf{z} \leftarrow [\mathbf{u} : \mathbf{v}]$ ;
7:    $(\mathbf{x}, \phi^{C^j}) = \text{interpret}(\mathbf{z})$ ;
8:    $count = count + \text{CheckSAT}(\mathbf{x}, \phi^{C^j}, count, \text{hiThresh})$ ;
9:   if  $count \geq \text{hiThresh}$  then return  $hiThresh$ ;
10:   $k \leftarrow \text{nextGrayBit}(q - p, j)$ ;
11:   $(\mathbf{u}, \mathbf{v}) \leftarrow \text{enumNextREX}(\mathbf{D}, \mathbf{u}, \mathbf{v}, k)$ ;
12: return  $count$ 

```

key idea is to sample cubes uniformly at random from $\{1, 2, \dots, m\}$ till a cube ϕ^{C^j} is found such that $\mathbf{x} \models \phi^{C^j}$ (lines 2-5). The number of cubes sampled $c_{\mathbf{x}}$ divided by total number of cubes m is the estimate returned (line 6).

Procedure LogSATSearch in $\text{SymbolicDNFApproxMC}$ is based upon the LogSATSearch in ApproxMC2 [4]. As noted in the analysis of ApproxMC2 , such a logarithmic search procedure requires that the solution space for a hash function with $p + 1$ hash constraints is a subset of the solution space with p hash constraints. Furthermore, we want to preserve Row-Echelon nature of the resulting hash constraints. To this end, we first construct $\mathbf{D}^{[q \times (q-p)]}$ and $\mathbf{b}^{[q-1]}$ as follows:

To seed the construction procedure, in SampleBase (algorithm 4) we first randomly sample a 0/1 vector $\hat{\mathbf{b}}$ of size $q - 1$ which is the maximum number of hash constraints possible. We then construct a 0/1 matrix $\hat{\mathbf{D}}$ as follows: $\hat{\mathbf{D}}^{[(q-1) \times (q-sI)]} = \begin{bmatrix} \mathbf{G} \\ \mathbf{E} \end{bmatrix}$ where matrix $\mathbf{G}^{[sI \times (q-sI)]}$ is a random 0/1 matrix with sI rows, and matrix $\mathbf{E}^{[(q-sI) \times (q-sI)]}$ is defined as as follows:

- $\mathbf{E}[i, j] = 1$ if $i = j$
- $\mathbf{E}[i, j] = 0$ if $i > j$

Algorithm 7 CheckSAT($x, \phi^{C^i}, count, hiThresh$)

```

1:  $c_x \leftarrow 0$ ;
2: while  $count + c_x/m < hiThresh$  do
3:   Uniformly sample  $j$  from  $\{1, 2, \dots, m\}$ ;
4:    $c_x \leftarrow c_x + 1$ ;
5:   if  $x \models \phi^{C^j}$  then
6:     return  $c_x/m$ ;
7: return  $c_x/m$ 

```

Algorithm 8 Extract($\hat{\mathbf{D}}, \hat{\mathbf{b}}, \hat{\mathbf{y}}, p, q, sI$)

```

1:  $\hat{\mathbf{D}}' \leftarrow \hat{\mathbf{D}}[0 : p, 0 : q - sI]$ ;  $\mathbf{b} \leftarrow \hat{\mathbf{b}}[0 : p]$ ;
2:  $\mathbf{y} \leftarrow \hat{\mathbf{y}}[0 : p]$ ;
3: for ( $i = sI$ ;  $i < p$ ;  $i++$ ) do
4:   for ( $j = 0$ ;  $j < i$ ;  $j++$ ) do
5:     if  $\hat{\mathbf{D}}'[j, (i - sI)] == 1$  then
6:        $\hat{\mathbf{D}}'[j, \cdot] \leftarrow \hat{\mathbf{D}}'[j, \cdot] \oplus \hat{\mathbf{D}}'[i, \cdot]$ ;
7:        $\mathbf{b}[j] \leftarrow \mathbf{b}[j] \oplus \mathbf{b}[i]$ ;
8:        $\mathbf{y}[j] \leftarrow \mathbf{y}[j] \oplus \mathbf{y}[i]$ ;
9:  $\mathbf{D} \leftarrow \hat{\mathbf{D}}'[0 : p, p - sI : q - sI]$ ;
10: return  $\mathbf{D}, \mathbf{b}, \mathbf{y}$ 

```

■ $\Pr[\mathbf{E}[i, j] = 1] = \Pr[\mathbf{E}[i, j] = 0] = 0.5$ if $i < j$

The reason for this definition of $\hat{\mathbf{D}}$ is that for DNF counting we have a good lower bound on the number of hash constraints we can start with. The number of rows in \mathbf{G} corresponds to this lower bound. The definition of \mathbf{E} ensures that the rows of \mathbf{E} are linearly independent which results in a monotonically shrinking solution space.

The Extract procedure (algorithm 8) takes $\hat{\mathbf{D}}, \hat{\mathbf{b}}$ and $\hat{\mathbf{y}}$ and a number p as input and returns \mathbf{D}, \mathbf{b} and cell \mathbf{y} such that (\mathbf{D}, \mathbf{b}) represents a hash function from H_{REX} with p constraints and \mathbf{y} represents a cell. A precondition for Extract is $sI \leq p \leq q - 1$. In lines 1 and 2, the first p rows of $\hat{\mathbf{D}}$ and first p elements of $\hat{\mathbf{b}}$ and $\hat{\mathbf{y}}$ are selected as $\hat{\mathbf{D}}', \mathbf{b}$ and \mathbf{y} respectively. The first sI rows of $\hat{\mathbf{D}}'$ form the matrix \mathbf{G} in the definition of $\hat{\mathbf{D}}$ and the remaining $p - sI$ rows of $\hat{\mathbf{D}}'$ are the first $p - sI$ rows of matrix \mathbf{E} . Each row from sI to p is used to reduce the preceding rows in lines 5 to 8 so that the only non-zero elements of the first $p - sI$ columns are the leading 1s in rows sI to p . Thus Extract ensures that for a given $\hat{\mathbf{D}}, \mathbf{b}$ and $\hat{\mathbf{y}}$, the solution space of $\mathbf{D}^{[p \times (q-p)]}, \mathbf{b}^{[p]}$ and $\mathbf{y}^{[p]}$ is a superset of solution space of $\mathbf{D}^{[(p+1) \times (q-p-1)]}, \mathbf{b}^{[p+1]}$ and $\mathbf{y}^{[p+1]}$ for all p .

5 Analysis

In order to prove the correctness of SymbolicDNFAproxMC, we first state the following helper lemma. We defer the proofs to the Appendix.

► **Lemma 4.** For every $1 \leq p \leq q$ and let $\mu_p = |R_\phi|/2^p$. For every $\beta > 0$ and $0 < \varepsilon < 1$ we have

1. $\Pr[|\Omega_{\mathbf{y}} - \mu_p| > \frac{\varepsilon}{(1+\varepsilon)}\mu_p] \leq \frac{2}{(1+\varepsilon)^2\mu_p}$
2. $\Pr[\Omega_{\mathbf{y}} \leq \beta\mu_p] \leq \frac{2}{2+(1-\beta^2)\mu_p}$



The difference in lemma 4 and lemma 1 in [4] is that the probability bounds differ by a factor of 2. We account for this difference by making `hiThresh` in `SymbolicDNFAproxMC` twice the value of `hiThresh` in `ApproxMC2`. Therefore the rest of the proof of Theorem 7 (below) is exactly the same as the proof of Theorem 4 of [4]. For completeness, we first restate lemmas 2 and 3 from [4] below.

In the following, T_p denotes the event $(\Omega_{\mathbf{y}} < \text{hiThresh})$, and L_p and U_p denote the events $(\Omega_{\mathbf{y}} < \frac{|R_\phi|}{(1+\varepsilon)2^p})$ and $(\Omega_{\mathbf{y}} > \frac{|R_\phi|}{2^p}(1 + \frac{\varepsilon}{1+\varepsilon}))$ respectively. p^* denotes the integer $\lceil \log_2 |R_\phi| - \log_2(4.92(1 + \frac{1}{\varepsilon})^2) \rceil$

► **Lemma 5.** *The following bounds hold:*

1. $\Pr[T_{p^*-3}] \leq \frac{1}{62.5}$
2. $\Pr[L_{p^*-2}] \leq \frac{1}{20.68}$
3. $\Pr[L_{p^*-1}] \leq \frac{1}{10.84}$
4. $\Pr[L_{p^*} \cup U_{p^*}] \leq \frac{1}{4.92}$



Let B denote the event that `SymbolicDNFAproxMC` returns a pair $(2^p, nSols)$ such that $2^p * nSols$ does not lie in the interval $[\frac{|R_\phi|}{1+\varepsilon}, |R_\phi|(1 + \varepsilon)]$.

► **Lemma 6.** $\Pr[B] \leq 0.36$



► **Theorem 7.** *Let `SymbolicDNFAproxMC`($\phi, \varepsilon, \delta$) return count c . Then $\Pr[|R_\phi|/(1 + \varepsilon) \leq c \leq (1 + \varepsilon)|R_\phi|] \geq 1 - \delta$.*



Theorem 7 follows from lemmas 4, 5 and 6 and noting that `SymbolicDNFAproxMC` boosts the probability of correctness of the count returned by `SymbolicDNFAproxMCCore` to $1 - \delta$ by using median of $t = O(\log(1/\delta))$ calls.

► **Theorem 8.** `SymbolicDNFAproxMC` runs in $\tilde{O}(mn \log(1/\delta)/\varepsilon^2)$ time.⁴



6 Conclusion

Hashing-based techniques have emerged as a promising approach to obtain counting algorithms and tools that scale to large instances while providing strong theoretical guarantees. This has led to an interest in designing hashing-based algorithms for counting problems that are known to be amenable to fully polynomial randomized approximation schemes. The prior hashing-based approach [4] provided FPRAS for DNF but with complexity much worse than state-of-the-art techniques. In this work, we introduced (i) Symbolic Hashing, (ii) Stochastic Cell-Counting, and (iii) a new 2-universal family of hash functions, and obtained a hashing-based FPRAS for #DNF with complexity similar to state-of-the-art.

Given the recent interest in hashing-based techniques and generality of our contributions, we believe concepts introduced in this paper can lead to design of hashing-based techniques for other classes of constraints. For example, all prior versions of `ApproxMC` relied on deterministic SAT solvers for exactly counting the solutions in a cell for #CNF. The technique of Stochastic Cell-Counting opens up the door for the usage of probabilistic SAT solvers for #CNF. Furthermore, a salient feature of the H_{REX} family is the sparsity of its hash functions.

⁴ We say $f(n) \in \tilde{O}(g(n))$ if $\exists k : f(n) \in \mathcal{O}(g(n) \log^k(g(n)))$

In fact, the sparsity increases with the addition of constraints. Sparse hash functions have been shown to be desirable for efficiently solving CNF+XOR constraints [15, 9, 12]. An interesting direction for future work is to test H_{REX} family with CNF formulas.

Acknowledgements The authors thank Jeffrey Dudek, Supratik Chakraborty and Dror Fried for valuable discussions. Work supported in part by NSF grants CCF-1319459 and IIS-1527668, by NSF Expeditions in Computing project "ExCAPE: Expeditions in Computer Augmented Program Engineering", and by an IBM Graduate Fellowship. Kuldeep S. Meel is supported by the IBM PhD Fellowship and the Lodieska Stockbridge Vaughn Fellowship.

References

- 1 F. Bacchus, S. Dalmao, and T. Pitassi. Algorithms and complexity results for #SAT and Bayesian inference. In *Proc. of FOCS*, pages 340–351, 2003. URL: <http://dl.acm.org/citation.cfm?id=946243.946291>.
- 2 J Lawrence Carter and Mark N Wegman. Universal classes of hash functions. In *Proceedings of the ninth annual ACM symposium on Theory of computing*, pages 106–112. ACM, 1977.
- 3 S. Chakraborty, K. S. Meel, and M. Y. Vardi. A scalable approximate model counter. In *Proc. of CP*, pages 200–216, 2013.
- 4 S. Chakraborty, K. S. Meel, and M. Y. Vardi. Algorithmic improvements in approximate counting for probabilistic inference: From linear to logarithmic SAT calls. In *Proc. of IJCAI*, 2016.
- 5 Paul Dagum, Richard Karp, Michael Luby, and Sheldon Ross. An optimal algorithm for monte carlo estimation. *SIAM Journal on computing*, 29(5):1484–1496, 2000.
- 6 Nilesh Dalvi and Dan Suciu. Efficient query evaluation on probabilistic databases. *The VLDB Journal—The International Journal on Very Large Data Bases*, 16(4):523–544, 2007.
- 7 C. Domshlak and J. Hoffmann. Probabilistic planning via heuristic forward search and weighted model counting. *Journal of Artificial Intelligence Research*, 30(1):565–620, 2007.
- 8 Leonardo Duenas-Osorio, Kuldeep S Meel, Roger Paredes, and Moshe Y Vardi. Counting-based reliability estimation for power-transmission grids. In *AAAI*, pages 4488–4494, 2017.
- 9 S. Ermon, C. P. Gomes, A. Sabharwal, and B. Selman. Low-density parity constraints for hashing-based discrete integration. In *Proc. of ICML*, pages 271–279, 2014.
- 10 Stefano Ermon, Carla P. Gomes, Ashish Sabharwal, and Bart Selman. Taming the curse of dimensionality: Discrete integration by hashing and optimization. In *Proc. of ICML*, pages 334–342, 2013.
- 11 V. Gogate and R. Dechter. Approximate counting by sampling the backtrack-free search space. In *Proc. of the AAAI*, volume 22, page 198, 2007.
- 12 C. P. Gomes, J. Hoffmann, A. Sabharwal, and B. Selman. Short XORs for Model Counting; From Theory to Practice. In *SAT*, pages 100–106, 2007.
- 13 C. P. Gomes, A. Sabharwal, and B. Selman. Model counting: A new strategy for obtaining good bounds. In *Proc. of AAAI*, volume 21, pages 54–61, 2006.
- 14 Frank Gray. Pulse code communication, March 17 1953. US Patent 2,632,058.
- 15 Alexander Ivrii, Sharad Malik, Kuldeep S. Meel, and Moshe Y. Vardi. On computing minimal independent support and its applications to sampling and counting. *Constraints*, pages 1–18, 2015. URL: <http://dx.doi.org/10.1007/s10601-015-9204-z>, doi:10.1007/s10601-015-9204-z.
- 16 M.R. Jerrum, L.G. Valiant, and V.V. Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theoretical Computer Science*, 43(2-3):169–188, 1986. URL: <http://portal.acm.org/citation.cfm?id=11534.11537>.

- 17 R.M. Karp and M. Luby. Monte-carlo algorithms for enumeration and reliability problems. *Proc. of FOCS*, 1983.
- 18 R.M. Karp, M. Luby, and N. Madras. Monte-Carlo approximation algorithms for enumeration problems. *Journal of Algorithms*, 10(3):429–448, 1989.
- 19 Donald E Knuth. Generating all n-tuples. *The Art of Computer Programming*, 4, 2004.
- 20 George Markowsky, J Lawrence Carter, and Mark Wegman. Analysis of a universal class of hash functions. *Mathematical Foundations of Computer Science 1978*, pages 345–354, 1978.
- 21 James D Park and Adnan Darwiche. Complexity results and approximation strategies for map explanations. *Journal of Artificial Intelligence Research*, pages 101–133, 2006.
- 22 D. Roth. On the hardness of approximate reasoning. *Artificial Intelligence*, 82(1):273–302, 1996. URL: [http://dx.doi.org/10.1016/0004-3702\(94\)00092-1](http://dx.doi.org/10.1016/0004-3702(94)00092-1), doi:10.1016/0004-3702(94)00092-1.
- 23 T. Sang, F. Bacchus, P. Beame, H. Kautz, and T. Pitassi. Combining component caching and clause learning for effective model counting. In *Proc. of SAT*, 2004.
- 24 T. Sang, P. Beame, and H. Kautz. Performing bayesian inference by weighted model counting. In *Prof. of AAAI*, pages 475–481, 2005.
- 25 L. Stockmeyer. The complexity of approximate counting. In *Proc. of STOC*, pages 118–126, 1983.
- 26 Gilbert Strang. *Introduction to linear algebra*, volume 3. Wellesley-Cambridge Press Wellesley, MA, 1993.
- 27 S. Toda. On the computational power of PP and (+)P. In *Proc. of FOCS*, pages 514–519. IEEE, 1989.
- 28 L. Trevisan. Lecture notes on computational complexity. *Notes written in Fall*, 2002. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.71.9877&rep=rep1&type=pdf>.
- 29 L.G. Valiant. The complexity of enumeration and reliability problems. *SIAM Journal on Computing*, 8(3):410–421, 1979.

Appendix

Proof of Lemma 1

Proof. When the input formula ϕ to ApproxMC2 is in DNF, BSAT is invoked with a formula of the form $\phi \wedge \psi$ where ψ is a conjunction of XOR constraints. For each cube ϕ^{C_i} , BSAT proceeds by performing Gaussian Elimination on $\phi^{C_i} \wedge \psi$. Since the number of XOR constraints can be $\mathcal{O}(n)$, Gaussian elimination can take $\mathcal{O}(n^3)$ time resulting in a cumulative complexity of $\mathcal{O}(mn^3)$ for all cubes.

At most $\text{hiThresh} = \mathcal{O}(1/\varepsilon^2)$ solutions to each $\phi^{C_i} \wedge \psi$ may have to be enumerated and each enumeration requires $\mathcal{O}(n^2)$ time. Therefore the complexity of enumeration is $\mathcal{O}(mn^2/\varepsilon^2)$. Thus the BSAT runs in $\mathcal{O}(mn^3 + mn^2/\varepsilon^2)$ time in the worst case. ◀ ◀

Proof of Lemma 2

Proof. Theorem 4 in [4] shows that ApproxMC2 makes $\mathcal{O}(\log n \log(1/\delta))$ calls to BSAT. Substituting the complexity of BSAT from lemma 1, we get $\mathcal{O}((mn^3 + mn^2/\varepsilon^2) \log n \log(1/\delta))$ complexity for ApproxMC2. ◀ ◀

Proof of Theorem 3:

Proof. Let h be a random hash function from $H_{\text{REX}}(\mathbf{q}, \mathbf{p})$ with $\mathbf{A}^{[\mathbf{p} \times \mathbf{q}]}$ as its matrix. Let $\mathbf{x}_1, \mathbf{x}_2$ be any two assignments such that $\mathbf{x}_1 \neq \mathbf{x}_2$. To prove 2-universality of H_{REX} , we need to show that for all \mathbf{x} and \mathbf{y} :

$$\Pr[h(\mathbf{x}_1) = h(\mathbf{x}_2)] \leq \frac{1}{2^{\mathbf{p}}} \quad (1)$$

$$\Pr[h(\mathbf{x}) = \mathbf{y}] = \frac{1}{2^{\mathbf{p}}} \quad (2)$$

Equation 2 follows from the random choice of \mathbf{b} . In particular, for every \mathbf{x} and \mathbf{y} for a chosen \mathbf{A} , there is a unique \mathbf{b} such that $\mathbf{A} \cdot \mathbf{x} \oplus \mathbf{b} = \mathbf{y}$. Since there are $2^{\mathbf{p}}$ possible choices for \mathbf{b} , we have $\Pr[\mathbf{A} \cdot \mathbf{x} \oplus \mathbf{b} = \mathbf{y}] = \frac{1}{2^{\mathbf{p}}}$.

Let $\mathbf{x} = \mathbf{x}_1 \oplus \mathbf{x}_2$. Note that $\mathbf{x} \neq \mathbf{0}^{[\mathbf{q}]}$ since $\mathbf{x}_1 \neq \mathbf{x}_2$. We prove that $\Pr[h(\mathbf{x}) = \mathbf{0}^{[\mathbf{q}]}] \leq 1/2^{\mathbf{p}}$ which is equivalent to 1. This is the same as showing $\Pr[\mathbf{A} \cdot \mathbf{x} = \mathbf{0}^{[\mathbf{q}]}] \leq 1/2^{\mathbf{p}}$. We can write $\mathbf{A} = [\mathbf{I}^{[\mathbf{p} \times \mathbf{p}]} : \mathbf{D}^{[\mathbf{p} \times (\mathbf{q} - \mathbf{p})]}]$ and \mathbf{x} as $\mathbf{x} = [\mathbf{u}^{[\mathbf{p}]} : \mathbf{v}^{[(\mathbf{q} - \mathbf{p})]}]$ in block notation. Then we have $\mathbf{A} \cdot \mathbf{x} = \mathbf{I} \cdot \mathbf{u} \oplus \mathbf{D} \cdot \mathbf{v}$. Since $\mathbf{x} \neq \mathbf{0}$, either $\mathbf{u} \neq \mathbf{0}$ or $\mathbf{v} \neq \mathbf{0}$ leading to the following three cases:

Case 1: If $\mathbf{u} \neq \mathbf{0}^{[\mathbf{p}]}$ and $\mathbf{v} = \mathbf{0}^{[\mathbf{q} - \mathbf{p}]}$, we get $\mathbf{I} \cdot \mathbf{u} \neq \mathbf{0}^{[\mathbf{p}]}$ and $\mathbf{D} \cdot \mathbf{v} = \mathbf{0}^{[\mathbf{p}]}$. Therefore $\mathbf{A} \cdot \mathbf{x} = \mathbf{I} \cdot \mathbf{u} \oplus \mathbf{D} \cdot \mathbf{v} \neq \mathbf{0}^{[\mathbf{p}]}$.

Case 2: If $\mathbf{u} = \mathbf{0}^{[\mathbf{p}]}$ and $\mathbf{v} \neq \mathbf{0}^{[\mathbf{q} - \mathbf{p}]}$, we get $\mathbf{I} \cdot \mathbf{u} = \mathbf{0}^{[\mathbf{p}]}$. From the proof of Theorem 1 in [20] we have $\Pr[\mathbf{D} \cdot \mathbf{v} = \mathbf{0}^{[\mathbf{p}]}] = \frac{1}{2^{\mathbf{p}}}$. Therefore $\Pr[\mathbf{A} \cdot \mathbf{x} = \mathbf{0}^{[\mathbf{p}]}] = \frac{1}{2^{\mathbf{p}}}$.

Case 3: If $\mathbf{u} \neq \mathbf{0}^{[\mathbf{p}]}$ and $\mathbf{v} \neq \mathbf{0}^{[\mathbf{q} - \mathbf{p}]}$, we get $\mathbf{I} \cdot \mathbf{u} \neq \mathbf{0}^{[\mathbf{p}]}$ and $\Pr[\mathbf{D} \cdot \mathbf{v} = \mathbf{0}^{[\mathbf{p}]}] = \frac{1}{2^{\mathbf{p}}}$. Therefore $\Pr[\mathbf{A} \cdot \mathbf{x} = \mathbf{0}^{[\mathbf{p}]}] = \frac{1}{2^{\mathbf{p}}}$.

◀

▶

Proof of Lemma 4:

Proof. The *coverage* of an assignment is $\text{cov}(\mathbf{x}) = \{j | \mathbf{x} \models \phi^{Cj}\}$. We have $\Pr[\mathbf{x} \models \phi^{Cj}] = |\text{cov}(\mathbf{x})|/m$ when j is drawn uniformly at random from $\{1, 2, \dots, m\}$. Also, $\sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \frac{1}{|\text{cov}(\mathbf{x})|} = |R_\phi|$, and $\mathbb{E}[\mathbf{c}_x] = m/|\text{cov}(\mathbf{x})|$ and $\mathbb{E}[\mathbf{c}_x^2] = \frac{2m^2}{|\text{cov}(\mathbf{x})|^2} - \frac{m}{|\text{cov}(\mathbf{x})|}$.

Let $\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}$ be a random variable such that $\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}} = \mathbf{c}_x/m$ if $h(\mathbf{x}) = \mathbf{y}$ and $(\mathbf{x}, \phi^{Cj}) = \text{interpret}(\mathbf{z})$, and $\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}} = 0$ otherwise, where the number of constraints in h is \mathbf{p} . Let $\eta = \Pr[h(\mathbf{x}) = \mathbf{y}] = 1/2^{\mathbf{p}}$. Then $\mathbb{E}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}] = \eta * \mathbb{E}[\frac{\mathbf{c}_x}{m}] = \frac{\eta}{|\text{cov}(\mathbf{x})|}$.

Now, $\mathbb{V}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}] = \mathbb{E}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}^2] - (\mathbb{E}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}])^2$. But $\mathbb{E}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}^2] = \eta * \mathbb{E}[(\mathbf{c}_x/m)^2] = \eta(\frac{2}{|\text{cov}(\mathbf{x})|^2} - \frac{1}{m * |\text{cov}(\mathbf{x})|})$ and $(\mathbb{E}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}])^2 = \frac{\eta^2}{|\text{cov}(\mathbf{x})|^2}$. Substituting back we get $\mathbb{V}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}] = \frac{2\eta - \eta^2}{|\text{cov}(\mathbf{x})|^2} - \frac{\eta}{m * |\text{cov}(\mathbf{x})|}$.

Define $\tau_{\mathbf{y}} = \sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}$. Clearly $\tau_{\mathbf{y}} = |\Omega_{\mathbf{y}}|$. We have $\mathbb{E}[\tau_{\mathbf{y}}] = \sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \mathbb{E}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}] = \sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \frac{\eta}{|\text{cov}(\mathbf{x})|} = \eta |R_\phi|$. Also since $\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}$ s are 2-universal, $\mathbb{V}[\tau_{\mathbf{y}}] \leq \sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \mathbb{V}[\gamma_{(\mathbf{x}, \phi^{Cj}), \mathbf{y}}] = (2\eta - \eta^2) \sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \frac{1}{|\text{cov}(\mathbf{x})|^2} - \frac{\eta}{m} \sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \frac{1}{|\text{cov}(\mathbf{x})|}$. But $\sum_{(\mathbf{x}, \phi^{Cj}) \in U'} \frac{1}{|\text{cov}(\mathbf{x})|^2} \leq |R_\phi|$. Therefore $\mathbb{V}[\tau_{\mathbf{y}}] \leq 2\eta |R_\phi|$ which implies $\mathbb{V}[\tau_{\mathbf{y}}] \leq 2\mathbb{E}[\tau_{\mathbf{y}}]$.

Applying Chebyshev's inequality we get $\Pr[|\Omega_{\mathbf{y}} - \mu_{\mathbf{p}}| \geq \frac{\epsilon}{1+\epsilon} \mu_{\mathbf{p}}] \leq \frac{\mathbb{V}[\tau_{\mathbf{y}}]}{\frac{\epsilon^2}{(1+\epsilon)^2} \mu_{\mathbf{p}}^2}$. Rearranging the terms and simplifying, we get the first part of the lemma.

Using Paley-Zygmund inequality, we get $\Pr[\Omega_{\mathbf{y}} \leq \beta \mu_{\mathbf{p}}] \leq 1 - \frac{(1-\beta)^2 \mu_{\mathbf{p}}^2}{\mathbb{V}[\tau_{\mathbf{y}}] + (1-\beta^2) \mu_{\mathbf{p}}^2}$ from which we get the second part of the lemma. ◀ ▶

Proof of Theorem 8

Proof. SymbolicDNFAproxMC makes $\mathcal{O}(\log 1/\delta)$ calls to SymbolicDNFAproxMCCore.

SymbolicDNFAproxMCCore samples a hash matrix $\hat{\mathbf{D}}$ in $\mathcal{O}(n(\log m + \log(1/\epsilon^2)))$ time and makes one call each to LogSATSearch and BSAT. LogSATSearch search in turn makes up to $\mathcal{O}(\log \log m)$ calls to BSAT. BSAT first reduces the matrix $\hat{\mathbf{D}}$ to \mathbf{D} in $\mathcal{O}(n(\log m + \log(1/\epsilon^2)))$

XX:16 On Hashing-Based Approaches to Approximate DNF-Counting

time by calling `Extract`. Next, the check in line 9 of `BSAT` ensures that at most $m \times \text{hiThresh}$ calls are made to the check in line 5 of `CheckSAT` during one execution of `BSAT`. Every assignment is enumerated in $\mathcal{O}(n)$ and the check in line 5 of `CheckSAT` is performed in $\mathcal{O}(n)$ time. Therefore, the time complexity of `BSAT` is $\mathcal{O}(n \times m \times \text{hiThresh})$. Hence, the time complexity of each invocation of `SymbolicDNFAproxMCCore` is $\tilde{O}(mn/\varepsilon^2)$, which implies that the time complexity of `SymbolicDNFAproxMC` is $\tilde{O}(mn \log(1/\delta)/\varepsilon^2)$. ◀ ◀