

# Closing the Gap Between Short and Long XORs for Model Counting

**Shengjia Zhao**

Computer Science Department  
Tsinghua University  
zhaosj12@mails.tsinghua.edu.cn

**Sorathan Chaturapruek**

Computer Science Department  
Stanford University  
sorathan@cs.stanford.edu

**Ashish Sabharwal**

Allen Institute for AI  
Seattle, WA  
ashishs@allenai.org

**Stefano Ermon**

Computer Science Department  
Stanford University  
ermon@cs.stanford.edu

## Abstract

Many recent algorithms for approximate model counting are based on a reduction to combinatorial searches over random subsets of the space defined by parity or XOR constraints. Long parity constraints (involving many variables) provide strong theoretical guarantees but are computationally difficult. Short parity constraints are easier to solve but have weaker statistical properties. It is currently not known how long these parity constraints need to be. We close the gap by providing matching necessary and sufficient conditions on the required asymptotic length of the parity constraints. Further, we provide a new family of lower bounds and the first non-trivial upper bounds on the model count that are valid for arbitrarily short XORs. We empirically demonstrate the effectiveness of these bounds on model counting benchmarks and in a Satisfiability Modulo Theory (SMT) application motivated by the analysis of contingency tables in statistics.

## Introduction

Model counting is the problem of computing the number of distinct solutions of a given Boolean formula. It is a classical problem that has received considerable attention from a theoretical point of view (Valiant 1979b; Stockmeyer 1985), as well as from a practical perspective (Sang et al. 2004; Gogate and Dechter 2007). Numerous probabilistic inference and decision making tasks, in fact, can be directly translated to (weighted) model counting problems (Richardson and Domingos 2006; Gogate and Domingos 2011). As a generalization of satisfiability testing, the problem is clearly intractable in the worst case. Nevertheless, there has been considerable success in both exact and approximate model counting algorithms, motivated by a number of applications (Sang, Beame, and Kautz 2005).

Recently, approximate model counting techniques based on randomized hashing have emerged as one of the leading approaches (Gomes, Sabharwal, and Selman 2006; Chakraborty, Meel, and Vardi 2013a; Ermon et al. 2014; Ivrii et al. 2015; Achlioptas and Jiang 2015; Belle, Van den Broeck, and Passerini 2015). While approximate, these techniques provide strong guarantees on the accuracy of the results in a probabilistic sense. Further, these methods all reduce model counting to a small number of combinatorial

searches on a randomly projected version of the original formula, obtained by augmenting it with randomly generated parity or XOR constraints. This approach allows one to leverage decades of research and engineering in combinatorial reasoning technology, such as fast satisfiability (SAT) and SMT solvers (Biere et al. 2009).

While modern solvers have witnessed tremendous progress over the past 25 years, model counting techniques based on hashing tend to produce instances that are difficult to solve. In order to achieve strong (probabilistic) accuracy guarantees, existing techniques require each randomly generated parity constraint to be relatively long, involving roughly half of the variables in the original problem. Such constraints, while easily solved in isolation using Gaussian Elimination, are notoriously difficult to handle when conjoined with the original formula (Gomes et al. 2007; Ermon et al. 2014; Ivrii et al. 2015; Achlioptas and Jiang 2015). Shorter parity constraints, i.e., those involving a relative few variables, are friendlier to SAT solvers, but their statistical properties are not well understood.

Ermon et al. (2014) showed that long parity constraints are not strictly necessary, and that one can obtain the *same accuracy guarantees* using shorter XORs, which are computationally much more friendly. They provided a closed form expression, allowing an easy computation of an XOR length that suffices, given various parameters such as the number of problem variables, the number of constraints being added, and the size of the solution space under consideration. It is, however, currently not known how tight their sufficiency condition is, how it scales with various parameters, or whether it is in fact a necessary condition.

*We resolve these open questions by providing an analysis of the optimal asymptotic constraint length* required for obtaining high-confidence approximations to the model count. Surprisingly, for formulas with  $n$  variables, we find that when  $\Theta(n)$  constraints are added, a constraint length of  $\Theta(\log n)$  is both necessary and sufficient. This is a significant improvement over standard long XORs, which have length  $\Theta(n)$ . Constraints of logarithmic length can, for instance, be encoded efficiently with a polynomial number of clauses. We also study upper bounds on the minimum sufficient constraint length, which evolve from  $O(\log n)$  to  $O(n^\gamma \log^2 n)$  to  $n/2$  across various regimes of the number of parity constraints.

As a byproduct of our analysis, we obtain a *new family of probabilistic upper and lower bounds that are valid regardless of the constraint length* used. These upper and lower bounds on the model count reach within a constant factor of each other as the constraint density approaches the aforementioned optimal value. The bounds gracefully degrade as we reduce the constraint length and the corresponding computational budget. While lower bounds for arbitrary XOR lengths were previously known (Gomes et al. 2007; Ermon et al. 2013a), the upper bound we prove in this paper is the first non-trivial upper bound in this setting. Remarkably, even though we rely on random projections and therefore only look at subsets of the entire space (a *local* view, akin to traditional sampling), we are able to say something about the *global* nature of the space, i.e., a probabilistic upper bound on the number of solutions.

We evaluate these new bounds on standard model counting benchmarks and on a new counting application arising from the analysis of contingency tables in statistics. These data sets are common in many scientific domains, from sociological studies to ecology (Sheldon and Dietterich 2011). We provide a new approach based on SMT solvers and a bit-vector arithmetic encoding. Our approach scales very well and produces accurate results on a wide range of benchmarks. It can also handle additional constraints on the tables, which are very common in scientific data analysis problems, where prior domain knowledge translates into constraints on the tables (e.g., certain entries must be zero because the corresponding event is known to be impossible). We demonstrate the capability to handle structural zeroes (Chen 2007) in real experimental data.

## Preliminaries: Counting by Hashing

Let  $x_1, \dots, x_n$  be  $n$  Boolean variables. Let  $S \subseteq \{0, 1\}^n$  be a large, high-dimensional set<sup>1</sup>. We are interested in computing  $|S|$ , the number of elements in  $S$ , when  $S$  is defined succinctly through conditions or constraints that the elements of  $S$  satisfy and membership in  $S$  can be tested using an NP oracle. For example, when  $S$  is the set of solutions of a Boolean formula over  $n$  binary variables, the problem of computing  $|S|$  is known as model counting, which is the canonical  $\#\text{-P}$  complete problem (Valiant 1979b).

In the past few years, there has been growing interest in approximate probabilistic algorithms for model counting. It has been shown (Gomes, Sabharwal, and Selman 2006; Ermon et al. 2013b; Chakraborty, Meel, and Vardi 2013a; Achlioptas and Jiang 2015; Belle, Van den Broeck, and Passerini 2015) that one can reliably estimate  $|S|$ , both in theory and in practice, by repeating the following simple process: randomly partition  $S$  into  $2^m$  cells and select one of these lower-dimensional cells, and compute whether  $S$  has at least 1 element in this cell (this can be accomplished with a query to an NP oracle, e.g., invoking a SAT solver). Somewhat surprisingly, repeating this procedure a small number of times provides a constant factor approximation to  $|S|$  with high probability, even though counting problems (in  $\#\text{-P}$ )

<sup>1</sup>We restrict ourselves to the binary case for the ease of exposition. Our work can be naturally extended to categorical variables.

are believed to be significantly harder than decision problems (in NP).

The correctness of the approach relies crucially on how the space is randomly partitioned into cells. All existing approaches partition the space into cells using parity or XOR constraints. A parity constraint defined on a subset of variables checks whether an odd or even number of the variables take the value 1. Specifically,  $m$  parity (or XOR) constraints are generated, and  $S$  is partitioned into  $2^m$  equivalence classes based on which parity constraints are satisfied.

The way in which these constraints are generated affects the quality of the approximation of  $|S|$  (the model count) obtained. Most methods *randomly* generate  $m$  parity constraints by adding each variable to each constraint with probability  $f \leq 1/2$ . This construction can also be interpreted as defining a hash function, mapping the space  $\{0, 1\}^n$  into  $2^m$  hash bins (cells). Formally,

**Definition 1.** Let  $A \in \{0, 1\}^{m \times n}$  be a random matrix whose entries are Bernoulli i.i.d. random variables of parameter  $f \leq 1/2$ , i.e.,  $\Pr[A_{ij} = 1] = f$ . Let  $b \in \{0, 1\}^m$  be chosen uniformly at random, independently from  $A$ . Then,  $\mathcal{H}_{m \times n}^f = \{h_{A,b} : \{0, 1\}^n \rightarrow \{0, 1\}^m\}$ , where  $h_{A,b}(x) = Ax + b \pmod 2$  and  $h_{A,b} \in_R \mathcal{H}_{m \times n}^f$  is chosen randomly according to this process, is a family of  $f$ -sparse hash functions.

The idea to estimate  $|S|$  is to define progressively smaller cells (by increasing  $m$ , the number of parity constraints used to define  $h$ ), until the cells become so small that no element of  $S$  can be found inside a (randomly) chosen cell. The intuition is that the larger  $|S|$  is, the smaller the cells will have to be, and we can use this information to estimate  $|S|$ .

Based on this intuition, we give a hashing-based counting procedure (Algorithm 1, SPARSE-COUNT), which relies on an NP oracle  $\mathcal{O}_S$  to check whether  $S$  has an element in the cell. It is adapted from the SPARSE-WISH algorithm of Ermon et al. (2014). The algorithm takes as input  $n$  families of  $f$ -sparse hash functions  $\{\mathcal{H}_{i \times n}^{f_i}\}_{i=0}^n$ , used to partition the space into cells. In practice, line 7 is implemented using a SAT solver as an NP-oracle. In a model counting application, this is accomplished by adding to the original formula  $i$  parity constraints generated as in Definition 1 and checking the satisfiability of the augmented formula.

Typically,  $\{\mathcal{H}_{i \times n}^{1/2}\}$  is used, corresponding to XORs where each variable is added with probability  $1/2$  (hence with average length  $n/2$ ). We call these **long parity constraints**. In this case, it can be shown that SPARSE-COUNT will output a factor 16 approximation of  $|S|$  with probability at least  $1 - \Delta$  (Ermon et al. 2014). Unfortunately, checking satisfiability (i.e.,  $S(h_{A,b}^i) \geq 1$ , line 7) has been observed to be very difficult when many long parity constraints are added to a formula (Gomes et al. 2007; Ermon et al. 2014; Ivrii et al. 2015; Achlioptas and Jiang 2015). Note, for instance, that while a parity constraint of length one simply freezes a variable right away, a parity constraint of length  $k$  can be propagated only after  $k - 1$  variables have been set. From a theoretical perspective, parity constraints are known to be fundamentally difficult for the resolution proof system

---

**Algorithm 1** SPARSE-COUNT ( $\mathcal{O}_S, \Delta, \alpha, \{\mathcal{H}_{i \times n}^{f_i}\}_{i=0}^n$ )

---

```
1:  $T \leftarrow \left\lceil \frac{\log(1/\Delta)}{\alpha} \log n \right\rceil$ 
2:  $i = 0$ 
3: while  $i \leq n$  do
4:   for  $t = 1, \dots, T$  do
5:      $h_{A,b}^i \leftarrow$  hash function sampled from  $\mathcal{H}_{i \times n}^{f_i}$ 
6:     Let  $S(h_{A,b}^i) = |\{x \in S \mid h_{A,b}^i(x) = 0\}|$ 
7:      $w_i^t \leftarrow \mathbb{I}[S(h_{A,b}^i) \geq 1]$ , invoking  $\mathcal{O}_S$ 
8:   end for
9:   if Median( $w_i^1, \dots, w_i^T$ )  $< 1$  then
10:     break
11:   end if
12:    $i = i + 1$ 
13: end while
14: Return  $\lfloor 2^{i-1} \rfloor$ 
```

---

underlying SAT solvers (cf. exponential scaling of Tseitin tautologies (Tseitin 1968)). A natural question, therefore, is whether **short parity constraints** can be used in SPARSE-COUNT and provide reliable bounds for  $|S|$ .

Intuitively, for the method to work we want the hash functions  $\{\mathcal{H}_{i \times n}^{f_i}\}$  to have a small collision probability. In other words, we want to ensure that when we partition the space into cells, configurations from  $S$  are divided into cells evenly. This gives a direct relationship between the original number of solutions  $|S|$  and the (random) number of solutions in one (randomly chosen) cell,  $S(h)$ . More precisely, we say that the hash family *shatters*  $S$  if the following holds:

**Definition 2.** For  $\epsilon > 0$ , a family of hash functions  $\mathcal{H}_{i \times n}^f$   $\epsilon$ -shatters a set  $S$  if  $\Pr[S(h) \geq 1] \geq 1/2 + \epsilon$  when  $h \in_R \mathcal{H}_{i \times n}^f$ , where  $S(h) = |\{x \in S \mid h(x) = 0\}|$ .

The crucial property we need to obtain reliable estimates is that the hash functions (equivalently, parity constraints) are able to shatter sets  $S$  with arbitrary “shape”. This property is both sufficient and necessary for SPARSE-COUNT to provide accurate model counts with high probability:

**Theorem 1.** (Informal statement) *A necessary and sufficient condition for SPARSE-COUNT to provide a constant factor approximation to  $|S|$  is that each family  $\mathcal{H}_{i \times n}^{f_i}$   $\epsilon$ -shatters all sets  $S'$  of size  $|S'| = 2^{i+c}$  for some  $c \geq 2$ .*

A formal statement, along with all proofs, is provided in a companion technical report (Zhao et al. 2015).

Long parity constraints, i.e.,  $1/2$ -sparse hash functions, are capable of shattering sets of arbitrary shape. When  $h \in_R \mathcal{H}_{i \times n}^{\frac{1}{2}}$ , it can be shown that configurations  $x \in \{0, 1\}^n$  are placed into hash bins (cells) pairwise independently, and this guarantees shattering of sufficiently large sets of arbitrary shape. Recently, Ermon et al. (2014) showed that sparser hash functions can be used for approximate counting as well. In particular,  $f^*$ -sparse hash functions, for sufficiently large  $f^* \lesssim 1/2$ , were shown to have good enough shattering capabilities. It is currently not known whether  $f^*$  is the optimal constraint density.

## Asymptotically Optimal Constraint Density

We analyze the asymptotic behavior of the minimum constraint density  $f$  needed for SPARSE-COUNT to produce correct bounds with high confidence. As noted earlier, the bottleneck lies in ensuring that  $f$  is large enough for a randomly chosen hash bin to receive at least one element of the set  $S$  under consideration, i.e., the hash family shatters  $S$ .

**Definition 3.** Let  $n, m \in \mathbb{N}, n \geq m$ . For any fixed  $\epsilon > 0$ , the *minimum constraint density*  $f = \tilde{f}_\epsilon(m, n)$  is defined as the pointwise smallest function such that for any constant  $c \geq 2$ ,  $\mathcal{H}_{m \times n}^f$   $\epsilon$ -shatters all sets  $S \in \{0, 1\}^n$  of size  $2^{m+c}$ .

We will show (Theorem 2) that for any  $\epsilon > 0$ ,  $\tilde{f}_\epsilon(m, n) = \Omega(\frac{\log m}{m})$ , and this is asymptotically tight when  $\epsilon$  is small enough and  $m = \Theta(n)$ , which in practice is often the computationally hardest regime of  $m$ . Further, for the regime of  $m = \Theta(n^\beta)$  for  $\beta < 1$ , we show that  $\tilde{f}_\epsilon(m, n) = O(\frac{\log^2 m}{m})$ . Combined with the observation that  $\tilde{f}_\epsilon(m, n) = \Theta(1)$  when  $m = \Theta(1)$ , our results thus reveal how the minimum constraint density evolves from a constant to  $\Theta(\frac{\log m}{m})$  as  $m$  increases from a constant to being linearly related to  $n$ .

The *minimum average constraint length*,  $n \cdot \tilde{f}_\epsilon(m, n)$ , correspondingly decreases from  $n/2$  to  $O(n^{1-\beta} \log^2 n)$  to  $\Theta(\log n)$ , showing that in the computationally hardest regime of  $m = \Theta(n)$ , the parity constraints can in fact be represented using only  $2^{\Theta(\log n)}$ , i.e., a polynomial number of SAT clauses.

**Theorem 2.** *Let  $n, m \in \mathbb{N}, n \geq m$ , and  $\kappa > 1$ . The minimum constraint density,  $\tilde{f}_\epsilon(m, n)$ , behaves as follows:*

1. Let  $\epsilon > 0$ . There exists  $M_\kappa$  such that for all  $m \geq M_\kappa$ :

$$\tilde{f}_\epsilon(m, n) > \frac{\log m}{\kappa m}$$

2. Let  $\epsilon \in (0, \frac{3}{10})$ ,  $\alpha \in (0, 1)$ , and  $m = \alpha n$ . There exists  $N$  such that for all  $n \geq N$ :

$$\tilde{f}_\epsilon(m, n) \leq \left(3.6 - \frac{5}{4} \log_2 \alpha\right) \frac{\log m}{m}$$

3. Let  $\epsilon \in (0, \frac{3}{10})$ ,  $\alpha, \beta \in (0, 1)$ , and  $m = \alpha n^\beta$ . There exists  $N_\kappa$  such that for all  $n \geq N_\kappa$ :

$$\tilde{f}_\epsilon(m, n) \leq \frac{\kappa(1-\beta) \log^2 m}{2\beta m}$$

The lower bound in Theorem 2 follows from analyzing the shattering probability of an  $m+c$  dimensional hypercube  $S_c = \{x \mid x_j = 0 \ \forall j > m+c\}$ . Intuitively, random parity constraints of density smaller than  $\frac{\log m}{m}$  do not even touch the  $m+c$  relevant (i.e., non-fixed) dimensions of  $S_c$  with a high enough probability, and thus cannot shatter  $S_c$  (because all elements of  $S_c$  would be mapped to the same hash bin).

For the upper bounds, we exploit the fact that  $\tilde{f}(m, n)$  is at most the  $f^*$  function introduced by Ermon et al. (2014) and provide an upper bound on the latter. Intuitively,  $f^*$  was defined as the minimum function such that the variance of  $S(h)$  is relatively small. The variance was upper bounded by

considering the worst case “shape” of  $|S|$ : points packed together unrealistically tightly, all fitting together within Hamming distance  $w^*$  of a point. For the case of  $m = \alpha n$ , we observe that  $w^*$  must grow as  $\Theta(n)$ , and divide the expression bounding the variance into two parts: terms corresponding to points that are relatively close (within distance  $\lambda n$  for a particular  $\lambda$ ) are shown to contribute a vanishingly small amount to the variance, while terms corresponding to points that are farther apart are shown to behave as if they contribute to  $S(h)$  in a pairwise independent fashion. The  $\frac{\log m}{m}$  bound is somewhat natural and also arises in the analysis of the rank of sparse random matrices and random sparse linear systems (Kolchin 1999). For example, this threshold governs the asymptotic probability that a matrix  $A$  generated as in Definition 1 has full rank (Cooper 2000). The connection arises because, in our setting, the rank of the matrix  $A$  affects the quality of hashing. For example, an all-zero matrix  $A$  (of rank 0) would map all points to the same hash bucket.

### Improved Bounds on the Model Count

In the previous sections, we established the optimal (smallest) constraint density that provides a constant factor approximation on the model count  $|S|$ . However, depending on the size and structure of  $S$ , even adding constraints of density  $f^* \ll 0.5$  can lead to instances that cannot be solved by modern SAT solvers (see Table 1 below).

In this section we show that for  $f < f^*$  we can still obtain probabilistic upper and lower bounds. The bounds constitute a trade off between small  $f$  for easily solved NP queries and  $f$  close to  $f^*$  for guaranteed constant factor approximation.

To facilitate discussion, we define  $S(h) = |\{x \in S \mid h(x) = 0\}| = |S \cap h^{-1}(0)|$  to be the random variable indicating how many elements of  $S$  survive  $h$ , when  $h$  is randomly chosen from  $\mathcal{H}_{m \times n}^f$  as in Definition 1. Let  $\mu_S = \mathbb{E}[S(h)]$  and  $\sigma_S^2 = \text{Var}[S(h)]$ . Then, it is easy to verify that irrespective of the value of  $f$ ,  $\mu_S = |S|2^{-m}$ .  $\text{Var}[S(h)]$  and  $\Pr[S(h) \geq 1]$ , however, do depend on  $f$ .

### Tighter Lower Bound on $|S|$

Our lower bound is based on Markov’s inequality and the fact that the mean of  $S(h)$ ,  $\mu_S = |S|2^{-m}$ , has a simple linear relationship to  $|S|$ . Previous probabilistic lower bounds (Gomes et al. 2007; Ermon et al. 2013a) are based on the following observation: it is very unlikely for at least half of  $T$  repetitions of applying  $h$  to  $S$  to result in some element of  $S$  surviving unless there are at least  $2^{m-2}$  solutions in  $S$ . Otherwise,  $\mu_S$  would be too small (specifically,  $\leq 1/4$ ), making it unlikely for solutions to survive often.

Unlike previous methods, we not only check whether the estimated  $\Pr[S(h) \geq 1]$  is at least  $1/2$ , but also consider an empirical estimate of  $\Pr[S(h) \geq 1]$ . This results in a tighter lower bound, with a probabilistic correctness guarantee derived using Chernoff’s bound.

**Lemma 1.** *Let  $S \subseteq \{0, 1\}^n$ ,  $f \in [0, 1/2]$ , and for each  $m \in \{1, 2, \dots, n\}$ , let  $h_m \in \mathcal{R}_{m \times n}^f$ . Then,*

$$|S| \geq \max_{m=1}^n 2^m \Pr[S(h_m) \geq 1]. \quad (1)$$

Our theoretical lower bound is  $L = 2^m \Pr[S(h) \geq 1]$ , which satisfies  $|S| \geq L$  by the previous Lemma. In practice, we cannot compute  $\Pr[S(h) \geq 1]$  exactly, so our practical lower bound  $\hat{L}$  will be based on an empirical estimate  $\Pr_{\text{est}}[S(h) \geq 1]$  derived from samples. Because  $\hat{L}$  is a random variable, we would like to have a statement of the form  $\Pr[|\hat{L}| \geq L] \geq 1 - \delta$ , where the probability is with respect to  $\hat{L}$ . This is formalized by the following Theorem.

**Theorem 3.** *Let  $m > 0, T > 0$ , and  $h_m^1, \dots, h_m^T$  be hash functions sampled independently from  $\mathcal{H}_{m \times n}^f$ . Let  $Y_k = \mathbb{I}[S(h_m^k) \geq 1]$ ,  $Y = \sum_{k=1}^T Y_k$ , and  $\Pr_{\text{est}}[S(h) \geq 1] = Y/T$  be random variables. Let  $\kappa > 0, c > 0$ . Define a random variable  $\mathcal{B} = \mathcal{B}(S, h_m^1, \dots, h_m^T)$  as  $\frac{2^m c}{(1+\kappa)}$  if  $\Pr_{\text{est}}[S(h) \geq 1] \geq c$  and 0 otherwise. Then*

$$\Pr[|S| \geq \mathcal{B}] \geq 1 - \exp\left(-\frac{\kappa^2 c T}{(1+\kappa)(2+\kappa)}\right).$$

### New Upper Bound for $|S|$

The upper bound expression for  $f$  that we derive next is based on the contrapositive of the observation of Ermon et al. (2014) that the larger  $|S|$  is, the smaller an  $f$  suffices to shatter it.

For  $n, m, f$ , and  $\epsilon(n, m, q, f)$  from (Ermon et al. 2014), define:

$$v(q) = \frac{q}{2^m} \left(1 + \epsilon(n, m, q, f) \cdot (q-1) - \frac{q}{2^m}\right) \quad (2)$$

This quantity is an upper bound on the variance  $\text{Var}[S(h)]$  of the number of surviving solutions as a function of the size of the set  $q$ , the number of constraints used  $m$ , and the statistical quality of the hash functions, which is controlled by the constraint density  $f$ . The following Lemma characterizes the asymptotic behavior of this upper bound on the variance:

**Lemma 2.**  *$q^2/v(q)$  is an increasing function of  $q$ .*

Using Lemma 2, we are ready to obtain an upper bound on the size of the set  $S$  in terms of the probability  $\Pr[S(h) \geq 1]$  that at least one configuration from  $S$  survives after adding the randomly generated constraints.

**Lemma 3.** *Let  $S \subseteq \{0, 1\}^n$  and  $h \in \mathcal{R}_{m \times n}^f$ . Then*

$$|S| \leq \min \left\{ z \mid \frac{1}{1 + 2^{2m} v(z)/z^2} > \Pr[S(h) \geq 1] \right\}$$

The probability  $\Pr[S(h) \geq 1]$  is unknown, but can be estimated from samples. In particular, we can draw independent samples of the hash functions and get accurate estimates using Chernoff style bounds. This yields the following theorem:

**Theorem 4.** *Let  $S \subseteq \{0, 1\}^n$ . Let  $\Delta \in (0, 1)$ . Suppose we draw  $T \geq 24 \ln \frac{1}{\Delta}$  hash functions  $h_1, \dots, h_T$  randomly from  $\mathcal{H}_{m \times n}^f$ . Let*

$$U(n, m, f) = \min \left\{ z \mid \frac{1}{1 + 2^{2m} v(z)/z^2} \geq \frac{3}{4} \right\}$$

Let  $\mathcal{A}(S, h_1, \dots, h_T)$  be a random variable that equals  $U(n, m, f)$  if  $\text{Median}(\mathbb{I}[S(h_1) = 0], \dots, \mathbb{I}[S(h_T) = 0]) = 1$ , and  $2^n$  otherwise. Then

$$\Pr[|S| \leq \mathcal{A}(S, h_1, \dots, h_T)] \geq 1 - \Delta \quad (3)$$

Note that the upper bound on  $|S|$  given by Theorem 4 is vacuous unless  $\text{Median}(\mathbb{I}[S(h_1) = 0], \dots, \mathbb{I}[S(h_T) = 0]) = 1$ . For brevity, let  $\mathcal{E}$  denote this random event.  $\mathcal{E}$  can be extremely unlikely for certain values of  $m$ ,  $|S|$ , and  $f$ . For instance, if  $|S| > 0$  and  $m = 0$ , event  $\mathcal{E}$  is impossible, and the upper bound is necessarily vacuous. However, for  $|S| > 0$ , as  $m$  grows (i.e., more and more independent parity constraints are added),  $\Pr[S(h) = 0]$  increases and eventually approaches 1, making  $\mathcal{E}$  increasingly likely. It is possible to show via Markov’s inequality that, regardless of the value of  $f$ ,  $\mathcal{E}$  is likely to occur when  $m \geq \log |S| + c$  for some constant  $c > 0$ . Thus, for large enough  $m$ , one expects Theorem 4 to provide a non-trivial upper bound on  $|S|$ .

When  $f = 0.5$ ,  $v(z) = z/2^m(1 - 1/2^m)$ . This means that  $U(n, m, 0.5) \approx 3 \cdot 2^m$ . For  $m^* \approx \log |S| + c$  and  $f = 0.5$ , Theorem 4 provides a good upper bound on  $|S|$  (with multiplicative error bounded by a constant factor). This is essentially the standard analysis for counting using hash functions in the pairwise independent case (Trevisan 2004).

When  $f < 0.5$ , however, it is possible that  $U(n, m, f) \gg 2^m$ , resulting in potentially loose upper bounds on  $|S|$  even for small values of  $m$ .

We also note that when  $f = 0.5$ , event  $\mathcal{E}$  is unlikely to occur when  $m \leq \log |S| - c$  for some constant  $c$ . When  $f < 0.5$ , this need not be the case.

## Experimental Evaluation

### Model Counting Benchmarks <sup>2</sup>

We evaluate the quality of our new bounds on a standard model counting benchmark (ANOR2011) from Kroc, Sabharwal, and Selman (2011). Both lower and upper bounds presented in the previous section are parametric: they depend both on  $m$ , the number of constraints, and  $f$ , the constraint density. Increasing  $f$  is always beneficial, but can substantially increase the runtime. The dependence on  $m$  is more complex, and we explore it empirically. To evaluate our new bounds, we consider a range of values for  $f \in [0.01, 0.5]$ , and use a heuristic approach to first identify a promising value for  $m$  using a small number of samples  $T$ , and then collect more samples for that  $m$  to reliably estimate  $P[S(h) \geq 1]$  and improve the bounds on  $|S|$ .

We primarily compare our results to ApproxMC (Chakraborty, Meel, and Vardi 2013b), which can compute a constant factor approximation to user specified precision with arbitrary confidence (at the cost of more computation time). ApproxMC is similar to SPARSE-COUNT, and uses long parity constraints. For both methods, Cryptominisat version 2.9.4 (Soos, Nohl, and Castelluccia 2009) is used as the NP-oracle  $\mathcal{O}_S$  (with Gaussian elimination enabled), and the confidence parameter is set to 0.95, so that bounds reported hold with 95% probability.

<sup>2</sup>Source code for this experiment can be found at <https://github.com/ShengjiaZhao/XORModelCount>

Our results on the *Langford12* instance (576 variables, 13584 clauses,  $10^5 \approx \exp(11.5)$  solutions) are shown in Figure 1. The pattern is representative of all other instances we tested. The tradeoff between quality of the bounds and runtime, which is governed by  $f$ , clearly emerges. Instances with small  $f$  values can be solved orders of magnitude faster than with full length XORs ( $f \approx 0.5$ ), but provide looser bounds. Interestingly, lower bounds are not very sensitive to  $f$ , and we empirically obtain good bounds even for very small values of  $f$ . We also evaluate ApproxMC (horizontal and vertical lines) with parameter setting of  $\epsilon = 0.75$  and  $\delta = 0.05$ , obtaining an 8-approximation with probability at least 0.95. The runtime is 47042 seconds. It can be seen that ApproxMC and our bounds offer comparable model counts for dense  $f \approx 0.5$ . However, our method allows to trade off computation time against the quality of the bounds. We obtain non-trivial upper bounds using as little as 0.1% of the computational resources required with long parity constraints, a flexibility not offered by any other method.

Table 1 summarizes our results on other instances from the benchmark and compares them with ApproxMC with a 12 hour timeout. We see that for the instances on which ApproxMC is successful, our method obtains approximate model counts of comparable quality and is generally faster. While ApproxMC requires searching for tens or even hundreds of solutions during each iteration, our method needs only one solution per iteration. Further, we see that long parity constraints can lead to very difficult instances that cannot be solved, thereby reinforcing the benefit of provable upper and lower bounds using sparse constraints (small  $f$ ). Our method is designed to produce non-trivial bounds even when the computational budget is significantly limited, with bounds degrading gracefully with runtime.

### SMT Models for Contingency Table Analysis

In statistics, a *contingency table* is a matrix that captures the (multivariate) frequency distribution of two variables with  $r$  and  $c$  possible values, resp. For example, if the variables are gender (male or female) and handedness (left- or right-handed), then the corresponding  $2 \times 2$  contingency table contains frequencies for the four possible value combinations, and is associated with  $r$  row sums and  $c$  column sums, also known as the row and column marginals.

Fisher’s exact test (Fisher 1954) tests contingency tables for homogeneity of proportion. Given fixed row and column marginals, it computes the probability of observing the entries found in the table under the null hypothesis that the distribution across rows is independent of the distribution across columns. This exact test, however, quickly becomes intractable as  $r$  and  $c$  grow. Statisticians, therefore, often resort to approximate significance tests, such as the chi-squared test.

The associated counting question is: *Given  $r$  row marginals  $R_i$  and  $c$  column marginals  $C_j$ , how many  $r \times c$  integer matrices  $M$  are there with these row and column marginals?* When entries of  $M$  are restricted to be in  $\{0, 1\}$ , the corresponding contingency table is called binary. We are interested in counting both binary and integer tables.

This counting problem for integer tables is #P-complete

Table 1: Comparison of run time and bound quality on the ANOR2011 dataset. All true counts and bounds are in log scale.

SAT Instance	Num. Vars.	True Count	Upper Bound			Lower Bound			ApproxMC	
			UB	Runtime(s)	$f$	LB	Runtime (s)	$f$	Estimate	Runtime (s)
lang12	576	–	14.96	19000	0.5	11.17	1280	0.1	12.25	47042
wff.3.150.525	150	32.57	35.06	3800	0.5	31.69	3600	0.1	32.58	43571
2bitmax-6	252	67.52	69.72	355	0.5	67.17	1000	0.5	–	timeout
lang15	1024	–	352.5	240	0.02	19.58	6400	0.02	–	timeout
ls8-normalized	301	27.01	201.5	3800	0.02	25.34	3600	0.02	–	timeout
wff.3.100.150	100	48.94	51.00	2100	0.5	48.30	3600	0.5	–	timeout
wff.4.100.500	100	–	69.31	1100	0.05	37.90	3600	0.05	–	timeout

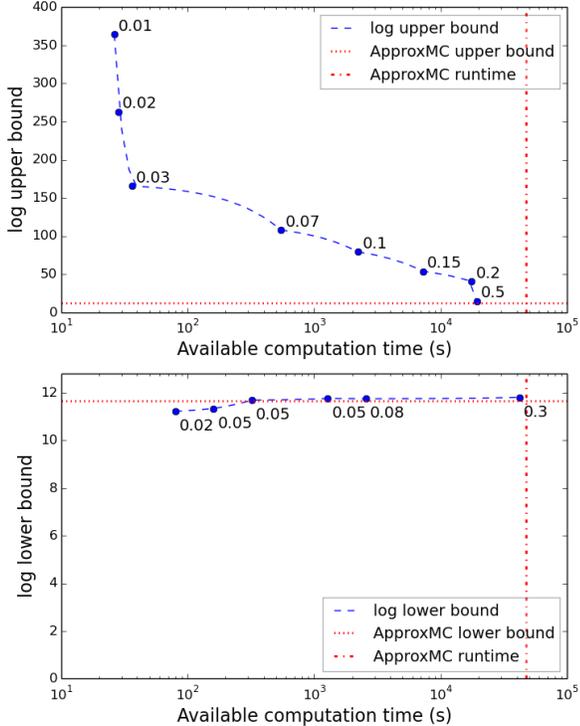


Figure 1: Log upper and lower bounds vs. computation time. Point labels indicate the value of  $f$  used.

even when  $r$  is 2 or  $c$  is 2 (Dyer, Kannan, and Mount 1997). Further, for binary contingency tables with so-called structural zeros (Chen 2007) (i.e., certain entries of  $M$  are required to be 0), we observe that the counting problem is still #P-complete. This can be shown via a reduction from the well-known *permanent* computation problem, which is #P-complete even for 0/1 matrices (Valiant 1979a).

Model counting for contingency tables is formulated most naturally using integer variables and arithmetic constraints for capturing the row and column sums. While integer linear programming (ILP) appears to be a natural fit, ILP solvers do not scale very well on this problem as they are designed to solve optimization problems and not feasibility queries. We therefore propose to encode the problem using a Satisfiability Modulo Theory (SMT) framework (Barrett, Stump, and

Tinelli 2010), which extends propositional logic with other underlying theories, such as bitvectors and real arithmetic. We choose a bitvector encoding where each entry  $a_{ij}$  of  $M$  is represented as a bitvector of size  $\lceil \log_2 \min(R_i, C_j) \rceil$ . The parity constraints are then randomly generated over the individual bits of each bitvector, and natively encoded into the model as XORs. As a solver, we use Z3 (De Moura and Bjørner 2008).

We evaluate our bounds on six datasets:

**Darwin’s Finches (df).** The marginals for this binary contingency table dataset are from Chen et al. (2005). This is one of the few datasets with known ground truth:  $\log_2 |S| \approx 55.8982$ , found using a clever divide-and-conquer algorithm of David desJardins. The 0-1 label in cell  $(x, y)$  indicates the presence or absence of one of 13 finch bird species  $x$  at one of 17 locations  $y$  in the Galápagos Islands. To avoid trivialities, we drop one of the species that appears in every island, resulting in  $12 \times 17 = 204$  binary variables.

**Climate Change Perceptions (icons).** This  $9 \times 6$  non-binary dataset is taken from the `alymerR` package (West and Hankin 2008). It concerns lay perception of climate change. The dataset is based on a study reported by O’Neil (2008) in which human subjects are asked to identify which icons (such as polar bears) they find the most concerning. There are 18 structural zeros representing that not all icons were shown to all subjects.

**Social Anthropology (purum).** This  $5 \times 5$  non-binary dataset (West and Hankin 2008) concerns marriage rules of an isolated tribe in India called the Purums, which is subdivided into 5 sibs. Structured zeros represent marriage rules that prevent some sibs from marrying other sibs.

**Industrial Quality Control (iqd).** This  $4 \times 7$  non-binary dataset (West and Hankin 2008) captures an industrial quality control setting. Cell  $(x, y)$  is the number of defects in the  $x$ -th run attributable to machine  $y$ . It has 9 structured zeros, representing machines switched off for certain runs.

**Synthetic Data (synth).** This  $n \times n$  binary dataset contains *blocked matrices* (Golshan, Byers, and Terzi 2013). The row and column marginals are both  $\{1, n - 1, \dots, n - 1\}$ . It can be seen that a blocked matrix has either a value of 1 in entry  $(1, 1)$  or it has two distinct entries with value 1 in the first row and the first column, cell  $(1, 1)$  excluded. Instantiating the first row and the first column completely determines the rest of the table. It is also easy to verify that the desired count is  $1 + (n - 1)^2$ .

Table 2 summarizes the obtained lower and upper bounds

Table 2: Lower (LB) and upper (UB) bounds on  $\log_2 |S|$ . The trivial upper bound (Trv. UB) is the number of binary variables.  $f^*$  denotes the best previously known minimum  $f$  (Ermon et al. 2014) required for provable upper bounds.

Dataset	Table Size	$f^*$	LB ( $f$ )	$\log_2  S $	UB	Trv. UB
df	$12 \times 17$	0.18	<b>53</b> (0.03)	55.90	<b>150</b>	204
icons	$9 \times 6$	0.19	<b>58</b> (0.04)	-	<b>183</b>	236
purum	$5 \times 5$	0.26	<b>29</b> (0.13)	-	<b>52</b>	125
iqd	$4 \times 7$	0.34	<b>15</b> (0.10)	-	<b>17</b>	76
synth_8	$8 \times 8$	0.41	<b>5</b> (0.30)	5.64	<b>16</b>	64
synth_20	$20 \times 20$	0.42	<b>8</b> (0.40)	8.49	<b>14</b>	400

on the number of contingency tables, with a 10 minute timeout. For the datasets with ground truth, we see that very sparse parity constraints (e.g.,  $f = 0.03$  for the Darwin finches dataset, as opposed to a theoretical minimum of  $f^* = 0.18$ ) often suffice in practice to obtain very accurate lower bounds. For the iqd dataset, we obtain upper and lower bounds within a small constant factor. For other datasets, there is a wider gap between the upper and lower bounds. However, the upper bounds we obtain are orders of magnitude tighter than the trivial log-upper bounds, which is the number of variables in a binary encoding of the problem.

## Conclusions

We introduced a novel analysis of the randomized hashing schemes used by numerous recent approximate model counters and probabilistic inference algorithms. We close a theoretical gap, providing a tight asymptotic estimate for the minimal constraint density required. Our analysis also shows, for the first time, that even very short parity constraints can be used to generate non-trivial upper bounds on model counts. Thanks to this finding, we proposed a new scheme for computing anytime upper and lower bounds on the model count. Asymptotically, these bounds are guaranteed to become tight (up to a constant factor) as the constraint density grows. Empirically, given very limited computational resources, we are able to obtain new upper bounds on a variety of benchmarks, including a novel application for the analysis of statistical contingency tables.

A promising direction for future research is the analysis of related ensembles of random parity constraints, such as low-density parity check codes (Achlioptas and Jiang 2015).

## Acknowledgments

This work was supported by the Future of Life Institute (grant 2015-143902).

## References

Achlioptas, D., and Jiang, P. 2015. Stochastic integration via error-correcting codes. In *Proc. Uncertainty in Artificial Intelligence*.

Angluin, D., and Valiant, L. 1979. Fast probabilistic algorithms for hamiltonian circuits and matchings. *Journal of Computer and System Sciences* 18(2):155–193.

Barrett, C.; Stump, A.; and Tinelli, C. 2010. The Satisfiability Modulo Theories Library (SMT-LIB). [www.SMT-LIB.org](http://www.SMT-LIB.org).

Belle, V.; Van den Broeck, G.; and Passerini, A. 2015. Hashing-based approximate probabilistic inference in hybrid domains. In *Proceedings of the 31st Conference on Uncertainty in Artificial Intelligence (UAI)*.

Biere, A.; Heule, M.; van Maaren, H.; and Walsh, T. 2009. Handbook of satisfiability. frontiers in artificial intelligence and applications, vol. 185.

Calabro, C. 2009. *The Exponential Complexity of Satisfiability Problems*. Ph.D. Dissertation, University of California, San Diego.

Chakraborty, S.; Meel, K.; and Vardi, M. 2013a. A scalable and nearly uniform generator of SAT witnesses. In *Proc. of the 25th International Conference on Computer Aided Verification (CAV)*.

Chakraborty, S.; Meel, K.; and Vardi, M. 2013b. A scalable approximate model counter. In *Proc. of the 19th International Conference on Principles and Practice of Constraint Programming (CP)*, 200–216.

Chen, Y.; Diaconis, P.; Holmes, S. P.; and Liu, J. S. 2005. Sequential monte carlo methods for statistical analysis of tables. *Journal of the American Statistical Association* 100(469):109–120.

Chen, Y. 2007. Conditional inference on tables with structural zeros. *Journal of Computational and Graphical Statistics* 16(2).

Cooper, C. 2000. On the rank of random matrices. *Random Structures & Algorithms* 16(2):209–232.

De Moura, L., and Bjørner, N. 2008. Z3: An efficient smt solver. In *Tools and Algorithms for the Construction and Analysis of Systems*. Springer. 337–340.

Dyer, M.; Kannan, R.; and Mount, J. 1997. Sampling contingency tables. *Random Structures and Algorithms* 10(4):487–506.

Ermon, S.; Gomes, C. P.; Sabharwal, A.; and Selman, B. 2013a. Optimization with parity constraints: From binary codes to discrete integration. In *Proc. of the 29th Conference on Uncertainty in Artificial Intelligence (UAI)*.

Ermon, S.; Gomes, C. P.; Sabharwal, A.; and Selman, B. 2013b. Taming the curse of dimensionality: Discrete integration by hashing and optimization. In *Proc. of the 30th International Conference on Machine Learning (ICML)*.

Ermon, S.; Gomes, C. P.; Sabharwal, A.; and Selman, B. 2014. Low-density parity constraints for hashing-based discrete integration. In *Proc. of the 31st International Conference on Machine Learning (ICML)*, 271–279.

Fisher, R. 1954. *Statistical Methods for Research Workers*. Oliver and Boyd.

Gogate, V., and Dechter, R. 2007. Approximate counting by sampling the backtrack-free search space. In *Proc. of the 22nd National Conference on Artificial Intelligence (AAAI)*, volume 22, 198–203.

Gogate, V., and Domingos, P. 2011. Probabilistic theorem proving. In *Uncertainty in Artificial Intelligence*.

Golshan, B.; Byers, J.; and Terzi, E. 2013. What do row and column marginals reveal about your dataset? In *Advances in Neural Information Processing Systems*, 2166–2174.

Gomes, C. P.; Hoffmann, J.; Sabharwal, A.; and Selman, B. 2007. Short XORs for model counting: From theory to practice. In *Theory and Applications of Satisfiability Testing (SAT)*, 100–106.

Gomes, C. P.; Sabharwal, A.; and Selman, B. 2006. Model counting: A new strategy for obtaining good bounds. In *Proc. of the 21st National Conference on Artificial Intelligence (AAAI)*, 54–61.

Guruswami, V. 2010. Introduction to coding theory - lecture notes.

Ivrii, A.; Malik, S.; Meel, K. S.; and Vardi, M. Y. 2015. On computing minimal independent support and its applications to sampling and counting. *Constraints* 1–18.

Kolchin, V. F. 1999. *Random graphs*. Number 53 in Encyclopedia of Mathematics and its Applications. Cambridge University Press.

Kroc, L.; Sabharwal, A.; and Selman, B. 2011. Leveraging belief propagation, backtrack search, and statistics for model counting. *Annals of Operations Research* 184(1):209–231.

O’Neil, S. 2008. *An Iconic Approach to Communicating Climate Change*. Ph.D. Dissertation, School of Environmental Science, University of East Anglia.

Richardson, M., and Domingos, P. 2006. Markov logic networks. *Machine Learning* 62(1):107–136.

Sang, T.; Beame, P.; and Kautz, H. 2005. Solving Bayesian networks by weighted model counting. In *Proc. of the 20th National Conference on Artificial Intelligence (AAAI)*, volume 1, 475–482.

Sang, T.; Bacchus, F.; Beame, P.; Kautz, H.; and Pitassi, T. 2004. Combining component caching and clause learning for effective model counting. In *Theory and Applications of Satisfiability Testing (SAT)*.

Sheldon, D. R., and Dietterich, T. G. 2011. Collective graphical models. In *Advances in Neural Information Processing Systems*, 1161–1169.

Sinclair, A. 2011. Randomness and computation - lecture notes.

Soos, M.; Nohl, K.; and Castelluccia, C. 2009. Extending SAT solvers to cryptographic problems. In *Theory and Applications of Satisfiability Testing (SAT)*.

Stockmeyer, L. 1985. On approximation algorithms for #P. *SIAM Journal on Computing* 14(4):849–861.

Trevisan, L. 2004. Lecture notes on computational complexity.

Tseitin, G. S. 1968. On the complexity of derivation in the propositional calculus. In Slisenko, A. O., ed., *Studies in Constructive Mathematics and Mathematical Logic, Part II*.

Valiant, L. G. 1979a. The complexity of computing the permanent. *Theoretical computer science* 8(2):189–201.

Valiant, L. 1979b. The complexity of enumeration and reliability problems. *SIAM Journal on Computing* 8(3):410–421.

West, L. J., and Hankin, R. K. 2008. Exact tests for two-way contingency tables with structural zeros. *Journal of Statistical Software* 28(11):1–19.

Zhao, S.; Chaturapruek, S.; Sabharwal, A.; and Ermon, S. 2015. Closing the gap between short and long xors for model counting. Technical report, Stanford University.

## Appendix: Proofs

**Theorem 5** (Formal statement of Theorem 1). *Let  $\{\mathcal{H}_{i \times n}^{f_i}\}_{i=0}^n$  be families of  $f_i$ -sparse hash functions.*

- (Sufficiency) *If there exist  $c \geq 2$  and  $\epsilon > 0$  such that for all  $i$ ,  $\mathcal{H}_{i \times n}^{f_i}$   $\epsilon$ -shatters all sets  $S' \subseteq \{0, 1\}^n$  of size  $|S'| = 2^{i+c}$ , then for any set  $S$ ,  $0 < \Delta < 1$ , and  $\alpha \leq 2 (\min(\epsilon, 1/2 - 1/2^c))^2 \ln 2$ ,  $\text{SPARSE-COUNT}(\mathcal{O}_S, \Delta, \alpha, \{\mathcal{H}_{i \times n}^{f_i}\})$  outputs a  $2^{c+1}$  approximation of  $|S|$  with probability at least  $1 - \Delta$ .*
- (Necessity) *If there exists an  $i$  and a set  $S$  of size  $2^{i+c}$  such that for all  $\epsilon > 0$   $\mathcal{H}_{i \times n}^{f_i}$  does not  $\epsilon$ -shatter  $S$ , then for any choice of  $\alpha > 0$  and  $0 < \Delta < 1$ ,  $\text{SPARSE-COUNT}(\mathcal{O}_S, \Delta, \alpha, \{\mathcal{H}_{i \times n}^{f_i}\})$  outputs a  $2^c$  approximation of  $|S|$  with probability at most  $1/2$ .*

*Proof.* For the sufficiency part, we show that there exist  $c > 0$  and  $\delta > 2$  such that for all  $i$  two conditions hold:

(a) for all sets  $S \subseteq \{0, 1\}^n$  of size  $|S| \leq 2^{i-c}$

$$\Pr[S(h) = 0] \geq 1 - \frac{1}{\delta}$$

when  $h$  is chosen from  $\mathcal{H}_{i \times n}^{f_i}$

(b) for all sets  $S \subseteq \{0, 1\}^n$  of size  $|S| \geq 2^{i+c}$

$$\Pr[S(h) \geq 1] \geq 1 - \frac{1}{\delta}$$

when  $h$  is chosen from  $\mathcal{H}_{i \times n}^{f_i}$ . Standard analysis following Ermon et al. (2014) then implies that for any set  $S$  and  $0 < \Delta < 1$ , if  $\alpha \leq 2(1 - \frac{1}{\delta} - \frac{1}{2})^2 \ln 2$ , we have

$$\frac{|S|}{2^{c+1}} \leq \text{SPARSE-COUNT}(\mathcal{O}_S, \Delta, \alpha, \{\mathcal{H}_{i \times n}^{f_i}\}) \leq |S|2^c$$

with probability at least  $1 - \Delta$ .

The second condition (b) is implied by the shattering properties of  $h$  in the assumptions for some  $c = c'$  and  $\epsilon = 1/2 - \frac{1}{\delta}$ .

The first condition (a) is trivially satisfied for any  $c \geq 2$  and  $\delta = 2^c$ . Formally, we have

$$\begin{aligned} \Pr[S(h) > 0] &= \Pr[S(h) \geq 1] = \Pr[S(h) \geq 2^c \mu_S] \\ &\leq \frac{1}{2^c} \end{aligned}$$

from Markov’s inequality. Conditions (a) and (b) are therefore simultaneously met choosing  $c = c'$  and  $\delta = \min(2^c, \frac{1}{1/2 - \epsilon})$ .

For the necessity part, let  $S$  be a set of size  $2^{i+c}$  as in the statement of the Theorem, i.e., not shattered by  $\mathcal{H}_{i \times n}^{f_i}$ . Let us condition on the event that the outer loop of  $\text{SPARSE-COUNT}(\mathcal{O}_S, \Delta, \alpha, \{\mathcal{H}_{i \times n}^{f_i}\})$  reaches iteration  $i$ . For any  $T \geq 1$  (therefore, for any choice of  $\Delta$  and  $\alpha$ ), the while loop breaks at iteration  $i$  with probability at least  $1/2$  because by assumption  $\Pr[S(h) \geq 1] \leq 1/2$  when  $h$  is chosen from  $\mathcal{H}_{i \times n}^{f_i}$ . Otherwise,  $\mathcal{H}_{i \times n}^{f_i}$  would  $\epsilon$ -shatter  $S$  for some  $\epsilon > 0$ . Therefore, the output satisfies  $\frac{|S|}{2^c} \leq \text{SPARSE-COUNT}(\mathcal{O}_S, \Delta, \alpha, \{\mathcal{H}_{i \times n}^{f_i}\})$  with probability at most  $1/2$ . This also bounds the probability that the output is a  $2^c$  approximation of  $|S|$ .  $\square$

## Proofs of New Upper and Lower Bounds

We continue with proofs for the bounds on  $|S|$  for arbitrary constraint density  $f$ . Several of the following proofs will rely on the following notion:

**Definition 4** (Ermon et al. 2014). Let  $m, n \in \mathbb{N}, m \leq$

$n, f \leq \frac{1}{2}$ , and  $q \leq 2^n + 1$ . Then:

$$w^*(n, q) = \max \left\{ w \mid \sum_{j=1}^w \binom{n}{j} \leq q - 1 \right\} \quad (4)$$

$$r(n, q) = \left( q - 1 - \sum_{w=1}^{w^*(n, q)} \binom{n}{w} \right) \quad (5)$$

$$\begin{aligned} \epsilon(n, m, q, f) = \frac{1}{q-1} & \left[ \sum_{w=1}^{w^*(n, q)} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m \right. \\ & \left. + \frac{r}{2^m} (1 + (1-2f)^{w^*(n, q)+1})^m \right] \quad (6) \end{aligned}$$

We observe that  $w^*(n, q)$  is always at most  $n$ . We will often be interested in the case where  $q = 2^{m+c}$  for  $c \geq 0$ .

### New Lower Bound

*Proof of Lemma 1.* Let hash function  $h_m$  be drawn randomly from  $\mathcal{H}_{m \times n}^f$ . Recall that our random variable  $S(h_m)$  takes on a non-negative integer value. Then, for any  $m$ , by Markov's inequality,

$$\Pr[S(h_m) \geq 1] \leq \mathbb{E}[S(h_m)] = \frac{|S|}{2^m}.$$

Hence, for any  $m$ ,  $|S| \geq 2^m \Pr[S(h_m) \geq 1]$ . Taking the maximum over all choices of  $m$  finishes the proof.  $\square$

*Proof of Theorem 3.* We observe that  $\mathbb{E}[Y] = T\mathbb{E}[Y_1] = T\Pr[S(h_m^1) \geq 1]$ . There are two cases: either  $\mathbb{E}[Y] > cT/(1+\kappa)$  or  $\mathbb{E}[Y] \leq cT/(1+\kappa)$ .

Case 1:  $\mathbb{E}[Y] > cT/(1+\kappa)$ . This implies  $\Pr[S(h_m^1) \geq 1] > c/(1+\kappa)$ . From Lemma 1,  $|S| \geq 2^m c/(1+\kappa)$  which is no smaller than  $\mathcal{B}$  for all realizations of the random variables  $h_m^1, \dots, h_m^T$ . Thus, in this case,  $\Pr[|S| \geq \mathcal{B}] = 1$ .

Case 2:  $\mathbb{E}[Y] \leq cT/(1+\kappa)$ . Define  $Z = \sum_{k=1}^T Z_k$  where  $Z_k$  are i.i.d. Bernoulli variables with probability  $c/(1+\kappa)$  of being 1. Then  $\mathbb{E}[Z] = cT/(1+\kappa)$ . Since  $Y_k$  and  $Z_k$  are i.i.d. Bernoulli random variables with  $\mathbb{E}[Y_k] \leq \mathbb{E}[Z_k]$ , we have  $\Pr[Y \geq d] \leq \Pr[Z \geq d]$  for any  $d$ . Thus:

$$\begin{aligned} \Pr \left[ \Pr_{\text{est}}[S(h) \geq 1] \geq c \right] &= \Pr[Y \geq cT] \\ &\leq \Pr[Z \geq cT] \\ &= \Pr[Z \geq (1+\kappa)\mathbb{E}[Z]] \\ &\leq \exp \left( -\frac{\kappa^2 cT}{2+\kappa} \frac{cT}{1+\kappa} \right) \end{aligned}$$

where the last inequality follows from Chernoff's bound (Angluin and Valiant 1979; Sinclair 2011). Hence, with a probability of at least  $1 - \exp \left( -\frac{\kappa^2 cT}{(1+\kappa)(2+\kappa)} \right)$ ,  $\mathcal{B}$  must be 0. We therefore have:

$$\Pr[|S| \geq \mathcal{B}] \geq 1 - \exp \left( -\frac{\kappa^2 cT}{(1+\kappa)(2+\kappa)} \right)$$

This finishes the proof.  $\square$

### New Upper Bound

*Proof of Lemma 2.* Let  $f(q) = q^2/v(q)$ , where

$$v(q) = \frac{q}{2^m} \left( 1 + \epsilon(n, m, q, f) \cdot (q-1) - \frac{q}{2^m} \right)$$

is defined as in Definition 2. We show that  $v(q) < v(q+1)$  for all  $q$ . Removing constant terms in  $v(q)$  we see that it suffices to show that

$$g(q) = \frac{q}{B_1 + B_2 - q}, \quad (7)$$

where  $B_1 = 2^m + \sum_{w=1}^{w^*(n, q)} \binom{n}{w} (1+x^w)^m$  and  $B_2 = \left( q - 1 - \sum_{w=1}^{w^*(n, q)} \binom{n}{w} \right) (1+x^{w^*(n, q)+1})^m$ , is an increasing function of  $q$ . The relevant quantities are defined in Definition 4, and  $x = 1 - 2f$  for brevity. We show that  $g(q) < g(q+1)$ . We note that

$$w^*(n, q+1) = w^*(n, q) + h(q), \quad (8)$$

where  $h(q) \in \{0, 1\}$  and it is 1 only when  $\sum_{j=1}^{w^*(n, q)+1} \binom{n}{j} = q$  (by looking at the definition). Define

$$\begin{aligned} t(q) = B_1 + B_2 - q &= 2^m + \sum_{w=1}^{w^*(n, q)} \binom{n}{w} (1+x^w)^m + \\ & \left( q - 1 - \sum_{w=1}^{w^*(n, q)} \binom{n}{w} \right) (1+x^{w^*(n, q)+1})^m - q. \end{aligned} \quad (9)$$

**Case 1:**  $h(q) = 0$ . We have  $g(q) < g(q+1)$  if and only if

$$\frac{q}{t(q)} < \frac{q+1}{t(q) + (1+x^{w^*(n, q)+1})^m - 1}, \quad (10)$$

which is true if and only if

$$(1+x^{w^*(n, q)+1})^m q - q < t(q). \quad (11)$$

Expanding the definition of  $t(q)$ , we get that the above inequality is true if and only if

$$\begin{aligned} 0 < \left( 2^m - (1+x^{w^*(n, q)+1})^m \right) + \\ & \sum_{w=1}^{w^*(n, q)} \binom{n}{w} \left( (1+x^w)^m - (1+x^{w^*(n, q)+1})^m \right), \end{aligned} \quad (12)$$

which is true because both terms are positive.

**Case 2:**  $h(q) = 1$ . This implies  $\sum_{j=1}^{w^*(n, q)+1} \binom{n}{j} = q$ . We have  $g(q) < g(q+1)$  if and only if

$$\frac{q}{t(q)} < \frac{q+1}{t(q+1)}, \quad (13)$$

and we have

$$t(q+1) = 2^m + \sum_{w=1}^{w^*(n, q)+1} \binom{n}{w} (1+x^w)^m - q - 1. \quad (14)$$

Expanding the definition of  $t(q)$  and  $t(q+1)$ , we get that the above inequality is true if and only if

$$0 < \left(2^m - (1 + x^{w^*(n,q)+1})^m\right) + \sum_{w=1}^{w^*(n,q)} \binom{n}{w} \left((1 + x^w)^m - (1 + x^{w^*(n,q)+1})^m\right), \quad (15)$$

which is true because both terms are positive (note this is the same inequality as before).  $\square$

*Proof of Lemma 3.* Let  $Q \subseteq \{0, 1\}^n$  be any set of size exactly  $q$  and  $h \in_R \mathcal{H}_{m \times n}^f$ . Following Ermon et al. (2014), we can get a worst-case bound for the variance of  $Q(h)$  as a function of  $q$ . Regardless of the structure of  $Q$ , we have

$$\sigma^2(Q) \leq v(q) = \frac{q}{2^m} \left(1 + \epsilon(n, m, q, f)(q-1) - \frac{q}{2^m}\right)$$

where  $\sigma^2(Q) = \text{Var}[\mathbb{I}\{x \in Q \mid h(x) = 0\}]$  is the variance of the random variable  $Q(h)$ , and  $\epsilon(n, m, q, f)$  is from Definition 4. From Cantelli's inequality

$$\Pr[Q(h) > 0] \geq 1 - \frac{\sigma^2(Q)}{\sigma^2(Q) + \left(\frac{|Q|}{2^m}\right)^2} \geq 1 - \frac{v(q)}{v(q) + \left(\frac{q}{2^m}\right)^2}$$

We claim that

$$\frac{v(q)}{v(q) + \left(\frac{q}{2^m}\right)^2}$$

which gives a lower bound on the shattering probability, is a decreasing function of  $q$ . By dividing numerator and denominator by  $v(q)$ , it is sufficient to show that

$$\frac{\left(\frac{q}{2^m}\right)^2}{v(q)}$$

is increasing in  $q$ , which follows from Lemma 2. To prove the Lemma, suppose by contradiction that

$$|S| > \min \left\{ x \mid 1 - \frac{v(x)}{v(x) + \left(\frac{x}{2^m}\right)^2} > \Pr[S(h) > 0] \right\}$$

Since  $\frac{v(q)}{v(q) + \left(\frac{q}{2^m}\right)^2}$  is a decreasing function of  $q$ , and  $|S|$  is assumed to be larger than the smallest element in the set above, it holds that

$$1 - \frac{v(|S|)}{v(|S|) + \left(\frac{|S|}{2^m}\right)^2} > \Pr[S(h) > 0] \quad (16)$$

From Cantelli's inequality

$$\Pr[S(h) > 0] \geq 1 - \frac{\sigma^2(S)}{\sigma^2(S) + \left(\frac{|S|}{2^m}\right)^2} \geq 1 - \frac{v(|S|)}{v(|S|) + \left(\frac{|S|}{2^m}\right)^2}$$

where the second inequality holds because  $v(|S|)$  upper bounds the true variance  $\sigma^2(S)$ . The last inequality contradicts eq. (16).  $\square$

*Proof of Theorem 4.* There are two possibilities for the unknown value  $|S|$ : either  $|S| \leq U(n, m, f)$  or  $|S| > U(n, m, f)$ . Note that, by definition,  $|S| \leq 2^n$ .

Case 1:  $|S| \leq U(n, m, f)$ . In this case,  $|S| \leq \min\{U(n, m, f), 2^n\} \leq \mathcal{A}(S, h_1, \dots, h_T)$  for all realizations of the random variables  $h_1, \dots, h_T$ . Thus, eq. (3) holds trivially for any  $\Delta \geq 0$ .

Case 2:  $|S| > U(n, m, f)$ . In this case, rearranging terms, we obtain

$$|S| > \min \left\{ z \mid 1 - \frac{v(z)}{v(z) + \left(\frac{z}{2^m}\right)^2} \geq \frac{3}{4} \right\}$$

By the monotonicity of  $\frac{v(q)}{v(q) + \left(\frac{q}{2^m}\right)^2}$  shown in the proof of Lemma 3, it holds that

$$1 - \frac{v(|S|)}{v(|S|) + \left(\frac{|S|}{2^m}\right)^2} \geq \frac{3}{4}$$

From Cantelli's inequality

$$\begin{aligned} \Pr[S(h) > 0] &\geq 1 - \frac{\text{Var}(S(h))}{\text{Var}(S(h)) + \left(\frac{|S|}{2^m}\right)^2} \\ &\geq 1 - \frac{v(|S|)}{v(|S|) + \left(\frac{|S|}{2^m}\right)^2} \\ &\geq \frac{3}{4} \end{aligned}$$

where the second inequality is derived observing that  $v(|S|)$  is an upper bound on the true variance  $\text{Var}(S(h))$ . Thus  $\Pr[S(h) = 0] \leq \frac{1}{4}$ , and from Chernoff's bound

$$\Pr[\text{Median}(\mathbb{I}[S(h_1) = 0], \dots, \mathbb{I}[S(h_T) = 0]) = 1] \leq \exp(-T/24) \leq \Delta$$

By the definition of  $\mathcal{A}(S, h_1, \dots, h_T)$ , we therefore have

$$\Pr[\mathcal{A}(S, h_1, \dots, h_T) = 2^n] \geq 1 - \Delta$$

Since  $|S| \leq 2^n$ , eq. (3) follows immediately.  $\square$

## Proof of Theorem 2

Theorem 2 contains three statements:

1. Let  $\epsilon > 0$ .  $\kappa > 1$ . There exists  $M_\kappa$  such that for all  $m \geq M_\kappa$ :

$$\tilde{f}_\epsilon(m, n) > \frac{\log m}{\kappa m}$$

2. Let  $\epsilon \in (0, \frac{3}{10})$ ,  $\alpha \in (0, 1)$ , and  $m = \alpha n$ . There exists  $N$  such that for all  $n \geq N$ :

$$\tilde{f}_\epsilon(m, n) \leq \left(3.6 - \frac{5}{4} \log_2 \alpha\right) \frac{\log m}{m}$$

3. Let  $\epsilon \in (0, \frac{3}{10})$ ,  $\alpha, \kappa > 1$ ,  $\beta \in (0, 1)$ , and  $m = \alpha n^\beta$ . There exists  $N_\kappa$  such that for all  $n \geq N_\kappa$ :

$$\tilde{f}_\epsilon(m, n) \leq \frac{\kappa(1-\beta) \log^2 m}{2\beta m}$$

We will prove them in turn in the following subsections. The arguments will often use the inequality  $1 + x \leq \exp(x)$  which holds for any  $x \in \mathbb{R}$ . We also use the expression "for sufficiently large  $n$ " to mean the more formal statement

$$\exists N > 0, \forall n > N$$

### Part I: Lower Bound

*Proof of Theorem 2 (part 1).* Since the minimum constraint density must work for every set  $S$ , it must also work for the hypercube  $S_c = \{0, 1\}^{m+c} \times \{0\}^{n-m-c}$  with  $2^{m+c}$  elements. This is a set where the first  $m+c$  variables are "free", while the remaining  $n - m - c$  are fixed to 0.

Let  $h$  be a hash function drawn from  $\mathcal{H}_{m \times n}^f$ . Let's consider a parity constraint of the form  $a_{i1}x_1 \oplus \dots \oplus a_{in}x_n = b_i$  as in Definition 1. If  $a_{i1} = a_{i2} = \dots = a_{i(m+c)} = 0$  and  $b_i = 1$ , then  $\{x \in S_c : ax = b \bmod 2\} = \emptyset$ . If the constraint is constructed as in Definition 1, this happens with probability  $\frac{1}{2}(1-f)^{m+c}$ . Accumulating this probability over  $m$  independent parity constraints and setting  $f = \frac{\log m}{\kappa m}$  for any  $\kappa > 1$ , we obtain:

$$\begin{aligned} \Pr[S_c(h) > 0] &\leq \left(1 - \frac{1}{2}(1-f)^{m+c}\right)^m \\ &= \left(1 - \frac{1}{2}\left(1 - \frac{\log m}{\kappa m}\right)^{m+c}\right)^m \end{aligned}$$

For any  $\lambda > 1$ , it can be verified from the Taylor expansion of the exponential function that for any small enough  $x > 0$ ,  $1 - x \geq \exp(-\lambda x)$ . Observe that for any fixed  $\kappa > 1$ ,  $1 - \frac{\log m}{\kappa m} > 0$  as long as  $m$  is large enough. It follows that for any  $\gamma > 1$ , there exists an  $M_{\kappa, \gamma}$  such that for all  $m \geq M_{\kappa, \gamma}$ , the above expression is upper bounded by:

$$\begin{aligned} &\left(1 - \frac{1}{2} \exp\left(-\gamma \frac{\log m}{\kappa m} (m+c)\right)\right)^m \\ &= \left(1 - \frac{1}{2} m^{-\gamma(m+c)/(\kappa m)}\right)^m \\ &\leq \exp\left(-\frac{1}{2} m^{1-\gamma(m+c)/(\kappa m)}\right) \end{aligned}$$

where the last inequality follows from  $1+x \leq \exp(x)$ . Since  $\kappa > 1$ , we can choose  $\gamma$  such that  $1 < \gamma < \kappa$ . In this case, for large enough  $m$ , the last expression above is less than  $1/2$ . In other words, there exists an  $M_\kappa$  such that for all  $m \geq M_\kappa$ ,  $\Pr[S_c(h) > 0] < 1/2$ . It follows that for all such  $m$ , the minimum constraint density,  $\tilde{f}(m, n)$ , must be larger than  $\frac{\log m}{\kappa m}$ , finishing the proof.  $\square$

### Part II: Upper bound when $m = \Theta(n)$

To prove the upper bound for  $m = \Theta(n)$ , we will need to first establish a few lemmas. In all the proofs below we will assume  $m = \alpha n$ , with  $\alpha$  constant with respect to  $n$ . Let us denote the binary entropy function as  $H(p) \triangleq -p \log_2 p - (1-p) \log_2(1-p)$ . It is well known that  $H(0) = 0$ ,  $H(\frac{1}{2}) = 1$ , and it is monotonically increasing in the interval  $[0, \frac{1}{2}]$ . We use the following relationship between the sum of binomials and the binary entropy function:

**Proposition 1** (Guruswami 2010, Lemma 5). *For any  $n \in \mathbb{N}$  and  $\lambda \in [0, \frac{1}{2}]$ ,*

$$\sum_{j=0}^{\lambda n} \binom{n}{j} \leq 2^{H(\lambda)n}$$

**Lemma 4.** *Let  $\alpha \in (0, 1)$ , there exists a unique  $\lambda^* < \frac{1}{2}$  such that  $H(\lambda^*) = \alpha$ , where  $H$  is the binary entropy function. For all  $\lambda < \lambda^*$ ,*

$$\lim_{n \rightarrow \infty} \frac{\sum_{j=1}^{\lambda n} \binom{n}{j}}{2^{\alpha n}} = 0$$

*Proof.* We can always find a unique  $\lambda^* < \frac{1}{2}$  such that  $H(\lambda^*) = \alpha$ . This is because  $H(\lambda)$  increases monotonically from 0 to 1 as  $\lambda$  increases from 0 to  $\frac{1}{2}$ , so  $H^{-1}(\alpha)$  takes one and only one value in the range  $(0, 1/2)$ .

Furthermore, due to monotonicity,  $H(\lambda) < \alpha = H(\lambda^*)$  for all  $\lambda < \lambda^*$ . From Proposition 1, for any  $\lambda < \frac{1}{2}$ , the sum of binomials in the numerator of the desired quantity is at most  $2^{H(\lambda)n}$ . Hence, the fraction is at most  $2^{(H(\lambda)-\alpha)n}$ , which approaches 0 as  $n$  increases because  $H(\lambda) < \alpha$ . Since numerator and denominator are non-negative, the limit is zero and this concludes the proof.  $\square$

**Corollary 1.** *Let  $\alpha \in (0, 1)$ ,  $c \geq 2$ ,  $w^*(n, q)$  be as in Definition 4, and  $\lambda^* < \frac{1}{2}$  be such that  $H(\lambda^*) = \alpha$ . Then for all  $\lambda < \lambda^*$ , and any  $n$  sufficiently large*

$$w^*(n, 2^{m+c}) = w^*(n, 2^{\alpha n+c}) \geq \lambda n$$

*Proof of Corollary 1.* By Lemma 4, for all  $\lambda < H^{-1}(\alpha)$ , when  $n$  is sufficiently large,

$$\sum_{j=1}^{\lambda n} \binom{n}{j} < 2^{\alpha n} < 2^{m+c} - 1$$

Thus, it follows immediately from the definition of  $w^*$  that for sufficiently large  $n$ ,  $\lambda n \leq w^*(n, 2^{m+c})$ .  $\square$

**Remark 1.** Corollary 1, together with the trivial fact that  $w^*(n, q) \leq n$ , implies  $w^*(n, q) = \Theta(n)$  when  $m = \alpha n$  and  $q = 2^{m+c}$ .

**Lemma 5.** *For all  $\delta > 0$  and  $w \in \mathbb{R}$ , the function  $f_\delta(w) = \log(1 + \delta^w)$  is convex.*

*Proof.* We will show that the second derivative of  $f_\delta(w)$  is non-negative:

$$\begin{aligned} f'_\delta(w) &= \frac{\delta^w \log \delta}{1 + \delta^w} \\ f''_\delta(w) &= \frac{\delta^w (1 + \delta^w) (\log \delta)^2 - \delta^{2w} (\log \delta)^2}{(1 + \delta^w)^2} \\ &= \frac{\delta^w (\log \delta)^2}{(1 + \delta^w)^2} \geq 0 \end{aligned}$$

It follows that  $f_\delta(w)$  is convex.  $\square$

**Lemma 6.** *Let  $t > 0$ ,  $0 < \delta < 1$ ,  $k > -t \frac{\log(\frac{2}{1+\delta}-1)}{\log(1+\delta)}$ , and  $w \geq 0$ . Then for all  $m$  sufficiently large,*

$$\frac{\left(\frac{1}{2} + \frac{1}{2} \left(1 - \frac{k \log m}{m}\right)^w\right)^m}{m^{-tw} + (1 + \delta)^{-m}} < 1 \quad (17)$$

Lemma 6 is an attempt to simplify the expression below, which we will call  $\zeta(w)$  and make its dependence on  $k$  and  $m$  implicit.

$$\zeta(w) = \zeta(w, k, m) = \left( \frac{1}{2} + \frac{1}{2} \left( 1 - \frac{k \log m}{m} \right)^w \right)^m$$

Note that for large enough  $m$  such that  $\frac{k \log m}{m} \leq 1$ ,  $\zeta(w)$  is monotonically non-increasing in  $w$ , a property we will use. This term is too complex to study in detail, therefore, we upper bound it by the sum of two simpler expressions. The intuition for this bound is shown in Figure 2.

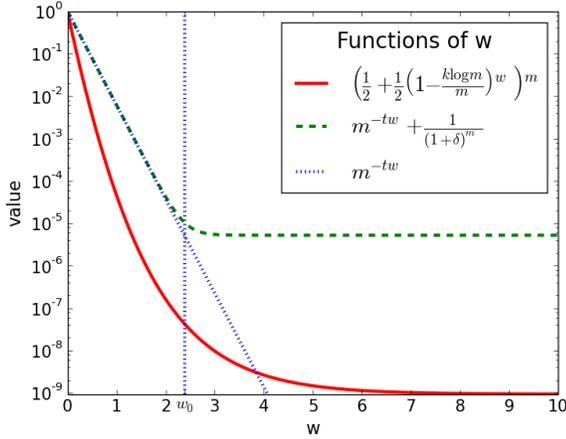


Figure 2: Graphical representation of Lemma 6

Lemma 6 can thus be restated as claiming  $\zeta(w) \leq m^{-tw} + (1 + \delta)^{-m}$  for  $m$  sufficiently large. Towards this end, when  $w < w_0$ , we show that  $m^{-tw}$  is the dominant term and that  $\zeta(w) < m^{-tw}$ . When  $w > w_0$ , we show that the term  $(1 + \delta)^{-m}$  dominates and that  $\zeta(w) < (1 + \delta)^{-m}$ . Combining these two regimes, we deduce that  $\zeta(w)$  must be upper bounded by their sum for all values of  $w$ . A formal proof follows.

*Proof.* We will show that for  $m$  sufficiently large,  $\zeta(w) \leq m^{-tw} + (1 + \delta)^{-m}$ . Assume w.l.o.g. that  $m$  is large enough to satisfy:

$$1 - k \frac{\log m}{m} \geq 0 \quad (18)$$

Next we consider the location  $w_0$  where the dominant term of  $m^{-tw} + (1 + \delta)^{-m}$  switches from  $m^{-tw}$  to  $(1 + \delta)^{-m}$ . This is where

$$m^{-tw_0} = (1 + \delta)^{-m}$$

which gives us  $w_0 = \frac{m \log(1+\delta)}{t \log m}$ . At this  $w_0$  we have

$$\begin{aligned} & \frac{\zeta(w_0)}{(1 + \delta)^{-m}} \\ & \leq \frac{(1 + \delta)^m}{2^m} \left( 1 + \left( 1 - \frac{k \log m}{m} \right)^{w_0} \right)^m \\ & \leq \frac{(1 + \delta)^m}{2^m} \left( 1 + \exp \left( -\frac{k \log m}{m} \frac{m \log(1 + \delta)}{t \log m} \right) \right)^m \\ & = \frac{(1 + \delta)^m}{2^m} \left( 1 + \exp \left( -\frac{k}{t} \log(1 + \delta) \right) \right)^m \\ & = \frac{(1 + \delta)^m}{2^m} \left( 1 + (1 + \delta)^{-\frac{k}{t}} \right)^m \\ & = \left( \frac{(1 + \delta)(1 + (1 + \delta)^{-\frac{k}{t}})}{2} \right)^m \end{aligned}$$

Clearly if we choose  $k$  such that  $\frac{(1+\delta)(1+(1+\delta)^{-\frac{k}{t}})}{2} < 1$ , then the entire expression is less than 1. This condition is satisfied if:

$$k > -t \frac{\log \left( \frac{2}{1+\delta} - 1 \right)}{\log(1 + \delta)}$$

Recall our earlier observation that  $\zeta(w)$  is monotonically non-increasing in  $w$  for large enough  $m$ . Thus, for  $m$  sufficiently large and any  $w \geq w_0$ , we have  $\zeta(w) \leq \zeta(w_0) < (1 + \delta)^{-m} \leq m^{-tw} + (1 + \delta)^{-m}$ .

Finally, let's consider the case where  $w < w_0$ , again assuming  $m$  sufficiently large so that  $\frac{k \log m}{m} < 1$ . Notice that

$$\log \zeta(w) = m \log \left( \frac{1}{2} + \frac{1}{2} \left( 1 - \frac{k \log m}{m} \right)^w \right)$$

is convex with respect to  $w$  because of Lemma 5. We have that for all positive  $m$ :

$$\log \zeta(0) = \log(1) = \log [m^{-t0}]$$

and for all  $m$  sufficiently large:

$$\log \zeta(w_0) < \log [(1 + \delta)^{-m}] = \log [m^{-tw_0}]$$

where the inequality is from the above analysis and the equality is by definition of  $w_0$ .

For  $w \in [0, w_0]$ , we can write  $w = (1 - \lambda)0 + \lambda w_0$  for some  $\lambda \geq 0$ . Therefore, for such  $w$  and for  $m$  sufficiently large, by convexity of  $\log \zeta(w)$ :

$$\begin{aligned} \log \zeta(w) & \leq (1 - \lambda) \log \zeta(0) + \lambda \log \zeta(w_0) \\ & \leq (1 - \lambda) \log [m^{-t0}] + \lambda \log [m^{-tw_0}] \\ & = \log [m^{-tw}] \end{aligned}$$

i.e.,  $\zeta(w) \leq m^{-tw} < m^{-tw} + (1 + \delta)^{-m}$ , as desired.  $\square$

**Lemma 7.** Let  $\alpha, \delta \in (0, 1)$ ,  $k > -\frac{\log \left( \frac{2}{1+\delta} - 1 \right)}{\log(1+\delta)}$ , and  $\lambda^* < \frac{1}{2}$  be such that  $H(\lambda^*) = \alpha \log_2(1 + \delta)$ . Then, for all  $\lambda < \lambda^*$ ,

$$\lim_{n \rightarrow \infty} \sum_{w=1}^{\lambda n} \binom{n}{w} \frac{1}{2^m} \left( 1 + \left( 1 - 2 \frac{k \log m}{m} \right)^w \right)^m = 0$$

*Proof.* By Lemma 6, we can select any  $0 < \delta < 1$ ,  $t > 1$  and  $k > -t \frac{\log(\frac{2}{1+\delta}-1)}{\log(1+\delta)}$  so that when  $n$  (or equivalently  $m = \alpha n$ ) is sufficiently large,

$$\begin{aligned} & \sum_{w=1}^{\lambda n} \binom{n}{w} \frac{1}{2^m} \left( 1 + \left( 1 - 2 \frac{k \log m}{m} \right)^w \right)^m \\ & \leq \sum_{w=1}^{\lambda n} \binom{n}{w} (m^{-tw} + (1+\delta)^{-m}) \\ & \leq \sum_{w=1}^{\lambda n} \frac{n^w}{w!} (\alpha n)^{-tw} + \frac{\sum_{w=1}^{\lambda n} \binom{n}{w}}{(1+\delta)^m} \end{aligned}$$

where we used the inequality  $\binom{n}{w} \leq \frac{n^w}{w!}$  for all  $n \in \mathbb{N}^*$  and  $0 \leq w \leq n$ . The first term of the sum can be driven to zero because when we choose any  $t > 1$

$$\begin{aligned} \lim_{n \rightarrow \infty} \sum_{w=1}^{\lambda n} \frac{n^w}{w!} (\alpha n)^{-tw} &= \lim_{n \rightarrow \infty} \sum_{w=1}^{\lambda n} \frac{\alpha^{-tw} n^{(1-t)w}}{w!} \\ &\leq \lim_{n \rightarrow \infty} \sum_{w=0}^{\infty} \frac{\alpha^{-tw} n^{(1-t)w}}{w!} - \frac{\alpha^0 n^0}{0!} \\ &= \lim_{n \rightarrow \infty} e^{\alpha^{-t} n^{1-t}} - 1 = 0 \end{aligned}$$

This requires that there exists  $t > 1$  such that  $k > -t \frac{\log(\frac{2}{1+\delta}-1)}{\log(1+\delta)}$ . For this,  $k > -\frac{\log(\frac{2}{1+\delta}-1)}{\log(1+\delta)}$  suffices.

The second term can be driven to zero when  $H(\lambda) < \alpha \log_2(1+\delta)$  as a direct consequence of Lemma 4.  $\square$

**Lemma 8.** Let  $\alpha \in (0, 1)$ ,  $m = \alpha n$ ,  $c \in \mathbb{N}$ ,  $q = 2^{m+c}$ . Let  $w^*$  and  $\epsilon(n, m, q, f)$  be as in Definition 4. Then, for all  $\gamma > 1$ ,  $k \geq 3.6 - \frac{5}{4} \log_2 \alpha$  and  $f = \frac{k \log m}{m}$ ,  $\exists N_k > 0$ , so that  $\forall n \geq N_k$ , we have

$$\epsilon(n, m, q, f) \leq \gamma \frac{2^c}{q-1}$$

*Proof.* For any  $\delta \in (0, 1)$ , if we choose  $\lambda^* < \frac{1}{2}$ , such that  $H(\lambda^*) = \alpha \log_2(1+\delta)$ , then  $\forall \lambda < \lambda^*$  by Corollary 1, we have for any value of  $n$  sufficiently large,  $\lambda n \leq w^*(n, q)$ .

Thus:

$$\begin{aligned} & (q-1) \epsilon(n, m, q, f) \\ &= \sum_{w=1}^{\lambda n} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ & \quad \sum_{w=\lambda n+1}^{w^*} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ & \quad \frac{r}{2^m} (1 + (1-2f)^{w^*+1})^m \\ &\leq \sum_{w=1}^{\lambda n} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ & \quad \sum_{w=1}^{w^*} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^{\lambda n})^m + \\ & \quad \frac{r}{2^m} (1 + (1-2f)^{\lambda n})^m \\ &= \sum_{w=1}^{\lambda n} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ & \quad \frac{q-1}{2^m} (1 + (1-2f)^{\lambda n})^m \\ &= A_n + B_n \end{aligned}$$

where  $A_n = \sum_{w=1}^{\lambda n} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m$  and  $B_n = \frac{q-1}{2^m} (1 + (1-2f)^{\lambda n})^m$ . By our choice of  $\lambda$ , and according to Lemma 7, if we choose any  $k > -\frac{\log(\frac{2}{1+\delta}-1)}{\log(1+\delta)}$  and  $f = \frac{k \log m}{m}$ , we can have

$$\lim_{n \rightarrow \infty} A_n = 0$$

For  $B_n$  we have

$$\begin{aligned} B_n &= \frac{q-1}{2^m} \left( 1 + \left( 1 - 2 \frac{k \log m}{m} \right)^{\lambda n} \right)^m \\ &\leq \frac{q-1}{2^m} \left( 1 + \exp \left( -\frac{2k \log m}{m} \lambda n \right) \right)^m \\ &\leq 2^c \left( 1 + m^{-\frac{2k\lambda}{\alpha}} \right)^m \\ &\leq 2^c \exp \left( m^{1-\frac{2k\lambda}{\alpha}} \right) \end{aligned}$$

where the inequalities follow from  $1+x \leq \exp(x)$ .

If we choose  $k$  such that  $1 - \frac{2k\lambda}{\alpha} < 0$ , or equivalently  $k > \frac{\alpha}{2\lambda}$ , we have

$$\limsup_{n \rightarrow \infty} B_n \leq 2^c$$

If we choose a  $k$  that is sufficiently large to satisfy both  $k > -\frac{\log(\frac{2}{1+\delta}-1)}{\log(1+\delta)}$  and  $k > \frac{\alpha}{2\lambda}$ , we have

$$\limsup_{n \rightarrow \infty} A_n + B_n \leq 2^c$$

which implies that for all  $\gamma > 1$ , and  $n$  sufficiently large,

$$\epsilon(n, m, q, f) \leq \frac{A_n + B_n}{q-1} \leq \gamma \frac{2^c}{q-1}$$

Now we obtain an upper bound on the value of  $k$  (the constraint density  $f$  is proportional to  $k$ , so we'd like this number to be as small as possible). From the derivation above, we can choose any  $0 < \delta < 1$ , and any  $k$  that satisfy the following inequalities:

$$k > -\frac{\log\left(\frac{2}{1+\delta} - 1\right)}{\log(1+\delta)} \quad (19)$$

$$k > \frac{\alpha}{2\lambda} \quad (20)$$

The second inequality also depends on  $\lambda$ , which we are free to choose as long as it satisfies  $\lambda < \lambda^*$ , or

$$H(\lambda) < \log_2(1+\delta)\alpha \equiv \sigma$$

We denote the latter term as  $\sigma$  to lighten the notation. This is satisfied if

$$\lambda < \frac{\sigma}{2\log_2(6/\sigma)} \leq H^{-1}(\sigma)$$

The latter inequality is adapted from Theorem 2.2 in (Calabro 2009)

Therefore, considering that  $\alpha < 1$ , the following condition is tighter than (20)

$$k > \frac{\log_2(6/\sigma)}{\log(1+\delta)}$$

Combining with (19) we have the following condition on  $k$ :

$$k > \max\left(-\frac{\log\left(\frac{2}{1+\delta} - 1\right)}{\log(1+\delta)}, \frac{\log_2(6/\sigma)}{\log_2(1+\delta)}\right)$$

We are allowed to choose  $\delta$  to give us the best bound. This choice is asymptotically insignificant, so we choose an arbitrary but empirically well performing  $\delta = 3/4$ , and derive

$$k > \max\left(-\frac{\log\left(\frac{2}{7/4} - 1\right)}{\log(7/4)}, \frac{\log_2\frac{6}{\log_2 7/4} - \log_2 \alpha}{\log_2(7/4)}\right)$$

which is approximately

$$k > \max(3.47, 3.58 - 1.23 \log_2 \alpha)$$

which is implied by

$$k \geq 3.6 - \frac{5}{4} \log_2 \alpha \quad (21)$$

as desired.  $\square$

The bounds in Lemma 8 are graphically shown in Figure 3. When  $m = \alpha n$ , the plot on the left shows the minimum  $f^*$  so that  $\epsilon(n, m, q, f^*)$  defined in Definition 4 is less than  $2/2^m$ . The plot on the right shows the empirical  $k$  so that the  $f^* \equiv k \frac{\log m}{m}$ . We also show the proved asymptotic bounds  $k = 3.6 - \frac{5}{4} \log_2 \alpha$  for comparison. As expected, the value of  $k$  found empirically does not exceed the bound (21).

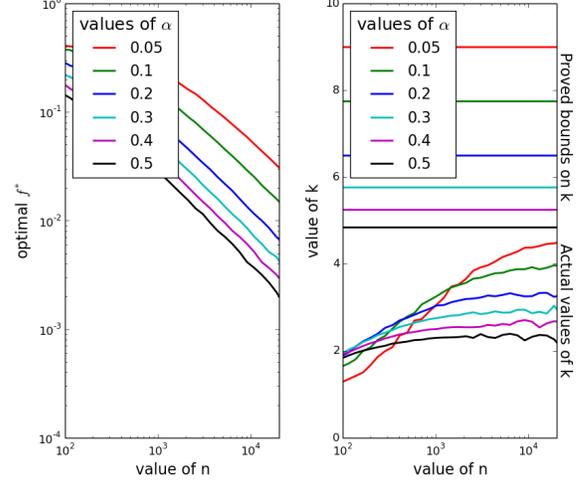


Figure 3: Empirical vs. proved bound on the value of  $k$

**Proof of Theorem 2 (part 2).** By Corollary 1 and Theorem 2 of Ermon et al. (2014), for set  $S$  with size  $|S| = q = 2^{m+c}$  and  $h \in \mathcal{H}_{m \times n}^f$ , a sufficient condition for ensuring that  $S$  is  $\epsilon$ -shattered, i.e.,  $\Pr[S(h) \geq 1] \geq \frac{1}{2} + \epsilon$ , is the “weak-concentration” condition given by:<sup>3</sup>

$$\epsilon(n, m, |S|, f) \leq \frac{\mu/(\frac{1}{2} + \epsilon) - 1}{|S| - 1} = \frac{2^c/(\frac{1}{2} + \epsilon) - 1}{q - 1} \quad (22)$$

By Lemma 8, when  $\gamma > 1$ ,  $f > (3.6 - \frac{5}{4} \log_2 \alpha) \frac{\log m}{m}$ ,  $c \geq 2$ , and  $m$  is sufficiently large:

$$\epsilon(n, m, |S|, f) \leq \gamma \frac{2^c}{q - 1}$$

Hence, to satisfy requirement (22), it suffices to have:

$$\gamma 2^c \leq \frac{2^c}{\frac{1}{2} + \epsilon} - 1$$

that is,  $\gamma \leq 1/(\frac{1}{2} + \epsilon) - 2^{-c}$ . We can therefore choose a  $\gamma > 1$  whenever  $1/(\frac{1}{2} + \epsilon) > 1 + 2^{-c}$ . Rearranging terms, this yields  $\epsilon < \frac{2^c - 1}{2(2^c + 1)}$ . Hence, for  $c \geq 2$ , it suffices to have  $\epsilon < 3/10$ . This completes the proof of part two of Theorem 2.  $\square$

### Part III: Upper Bound when $m = \Theta(n^\beta)$

Similar to Part II, we will first establish a few lemmas. We will assume for the rest of the reasoning that  $m = \alpha n^\beta$  for some constant  $\alpha, \beta \in (0, 1)$ .

**Lemma 9.** Let  $\alpha, \beta \in (0, 1), \gamma > 0, m = \alpha n^\beta$ , and  $\lambda^* = \frac{\gamma}{1-\beta}$ . Then for all  $\lambda < \lambda^*$ ,

$$\lim_{n \rightarrow \infty} \frac{\sum_{j=1}^{\lambda m / \log n} \binom{n}{j}}{2^{\gamma m}} = 0$$

<sup>3</sup>Note that the notation used for  $1/2 + \epsilon$  (for  $\epsilon > 0$ ) by Ermon et al. (2014) is  $1 - 1/\delta$  (for  $\delta > 2$ ).

*Proof.* For any  $1 \leq w \leq \frac{n}{2}$ ,

$$\begin{aligned} \log \left( \sum_{j=1}^w \binom{n}{j} \right) &\leq \log \left( w \binom{n}{w} \right) \leq \log \left( w \left( \frac{ne}{w} \right)^w \right) \\ &\leq \log w + w \log \frac{ne}{w} \leq w \log \frac{2ne}{w} \end{aligned}$$

When  $n$  is sufficiently large,  $1 \leq \lambda \frac{m}{\log n} \leq \frac{n}{2}$ . Let  $\epsilon > 0$  be any constant. Substituting  $w = \lambda m = \lambda \frac{\alpha n^\beta}{\log n}$ :

$$\begin{aligned} \log \left( \sum_{j=1}^{\lambda \alpha n^\beta / \log n} \binom{n}{j} \right) &\leq \frac{\lambda \alpha n^\beta}{\log n} \log \left( \frac{2e n^{1-\beta} \log n}{\lambda \alpha} \right) \\ &\leq (1 - \beta + \epsilon) \lambda \alpha n^\beta \end{aligned}$$

for large enough  $n$ . We thus have,

$$\lim_{n \rightarrow \infty} \frac{\sum_{j=1}^{\lambda m / \log n} \binom{n}{j}}{2^{\gamma m}} \leq \lim_{n \rightarrow \infty} 2^{((1-\beta+\epsilon)\lambda-\gamma)\alpha n^\beta}$$

Let  $\lambda^* = \gamma / (1 - \beta)$ . It follows that for any  $\lambda < \lambda^*$ ,  $\exists \epsilon > 0$ , such that

$$(1 - \beta + \epsilon) \lambda < \frac{1 - \beta + \epsilon}{1 - \beta} \gamma < \gamma$$

and the above limit can be driven to zero.  $\square$

**Corollary 2.** Let  $\alpha, \beta \in (0, 1)$ ,  $m = \alpha n^\beta$ ,  $c \geq 2$ ,  $w^*$  as in Definition 4, and  $\lambda^* = 1 / (1 - \beta)$ . Then, for all  $\lambda < \lambda^*$ ,  $\exists N > 0$  such that  $\forall n > N$ , we have

$$\frac{\lambda m}{\log n} \leq w^*(n, 2^{m+c})$$

**Lemma 10.** Let  $t > 0$ ,  $\delta \in (0, 1)$ ,  $w \geq 0$ , and  $k > 0$ . Then for all values of  $m$  sufficiently large,

$$\frac{\left( \frac{1}{2} + \frac{1}{2} \left( 1 - \frac{k \log^2 m}{m} \right)^w \right)^m}{m^{-tw} + (1 + \delta)^{-m}} < 1$$

*Proof.* Similar to Lemma 6, we will simplify notation by defining

$$\zeta(w) = \left( \frac{1}{2} + \frac{1}{2} \left( 1 - \frac{k \log^2 m}{m} \right)^w \right)^m$$

and assume that  $1 - \frac{k \log^2 m}{m} > 0$ . Consider the bound at

$$w_0 = \frac{m \log(1 + \delta)}{t \log m}$$

where  $m^{-tw} = (1 + \delta)^{-m}$ , then separately consider the cases of  $w$  being smaller and larger than  $w_0$ . When  $w = w_0$ ,

we have:

$$\begin{aligned} &\frac{\zeta(w_0)}{(1 + \delta)^{-m}} \\ &\leq \frac{(1 + \delta)^m}{2^m} \left( 1 + \exp \left( \frac{-w_0 k \log^2 m}{m} \right) \right)^m \\ &= \frac{(1 + \delta)^m}{2^m} \left( 1 + \exp \left( \frac{-k \log(1 + \delta) \log m}{t} \right) \right)^m \\ &= \frac{(1 + \delta)^m}{2^m} \left( 1 + (1 + \delta)^{-\frac{k}{t} \log m} \right)^m \\ &= \left( \frac{(1 + \delta)(1 + (1 + \delta)^{-\frac{k}{t} \log m})}{2} \right)^m \end{aligned}$$

It is easy to see that for any  $k > 0$ ,  $t > 0$ , when  $m$  is sufficiently large the base of this exponential quantity, and hence the quantity itself, is smaller than 1.

Now consider the general case of  $w > w_0$ . Because  $m > 0$  and  $\zeta(w)$  is monotonically non-increasing in  $w$ , we have:

$$\frac{\zeta(w)}{m^{-tw} + (1 + \delta)^{-m}} < \frac{\zeta(w_0)}{(1 + \delta)^{-m}}$$

which, by the above argument, is smaller than 1, as desired.

The remaining case of  $w \leq w_0$  is proved similar to the proof of Lemma 6. Due to the convexity of  $\log \zeta(w)$  with respect to  $w$ , combined with the fact that for any  $m > 0$ :

$$\log \zeta(0) = \log [m^{-t0}]$$

and for  $m$  sufficiently large:

$$\log \zeta(w_0) \leq \log [m^{-tw_0}]$$

we have that  $\exists M$  such that  $\forall m > M$ ,  $0 \leq w \leq w_0$

$$\log \zeta(w) \leq \log [m^{-tw}]$$

which implies

$$\zeta(w) < m^{-tw} + (1 + \delta)^{-m}$$

when  $w \leq w_0$ . Combined with the earlier similar result for  $w > w_0$ , this finishes the proof.  $\square$

**Lemma 11.** Let  $\alpha, \beta, \delta \in (0, 1)$ ,  $m = \alpha n^\beta$ ,  $\lambda^* = \frac{\log_2(1+\delta)}{1-\beta}$ ,  $\lambda < \lambda^*$ , and  $k > 0$ . Then

$$\lim_{n \rightarrow \infty} \sum_{w=1}^{\lambda m / \log n} \binom{n}{w} \frac{1}{2^m} \left( 1 + \left( 1 - 2 \frac{k \log^2 m}{m} \right)^w \right)^m = 0$$

*Proof.* By Lemma 10, for any  $t > 0$  and large enough  $n$  (and thus  $m$ ), the desired expression is at most:

$$\begin{aligned} &\sum_{w=1}^{\lambda m / \log n} \binom{n}{w} (m^{-tw} + (1 + \delta)^{-m}) \\ &\leq \sum_{w=1}^{\lambda m / \log n} \frac{n^w m^{-tw}}{w!} + \frac{\sum_{w=1}^{\lambda m / \log n} \binom{n}{w}}{(1 + \delta)^m} \end{aligned}$$

The second term here converges to zero as  $n \rightarrow \infty$  by applying Lemma 9 with  $\gamma$  set to  $\log_2(1 + \delta)$ . For the first term, we get:

$$\begin{aligned} \sum_{w=1}^{\lambda m / \log n} \frac{\alpha^{-tw} n^{(1-\beta)t w}}{w!} &\leq \sum_{w=0}^{\infty} \frac{\alpha^{-tw} n^{(1-\beta)t w}}{w!} - \frac{\alpha^0 n^0}{0!} \\ &= \exp(\alpha^{-t} n^{1-\beta t}) - 1 \end{aligned}$$

Choose any  $t > 1/\beta$ . Then the second term converges to zero as well.  $\square$

**Lemma 12.** Let  $\alpha, \beta \in (0, 1)$ ,  $m = \alpha n^\beta$ ,  $c \in \mathbb{Z}$ , and  $q = 2^{m+c}$ . Let  $w^*$  and  $\epsilon(n, m, q, f)$  be as in Definition 4. Then, for all  $\gamma > 1$ ,  $k > \frac{1-\beta}{2\beta}$  and  $f = \frac{k \log^2 m}{m}$ , and values of  $n$  greater than some  $N > 0$ ,

$$\epsilon(n, m, q, f) \leq \gamma \frac{2^c}{q-1}$$

*Proof.* For any  $0 < \delta < 1$ , let  $\lambda^* = \frac{\log_2(1+\delta)}{1-\beta}$ , then  $\forall \lambda < \lambda^*$ , by Lemma 2, for all values of  $n$  sufficiently large,

$$\frac{\lambda m}{\log n} \leq w^*(n) = w^*$$

We can write:

$$\begin{aligned} &(q-1)\epsilon(n, m, q, f) \\ &= \sum_{w=1}^{\lambda m / \log n} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ &\quad \sum_{w=1+\lambda m / \log n}^{w^*} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ &\quad \frac{r}{2^m} (1 + (1-2f)^{w^*+1})^m \\ &\leq \sum_{w=1}^{\lambda m / \log n} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ &\quad \sum_{w=1}^{w^*} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^{\lambda m / \log n})^m + \\ &\quad \frac{r}{2^m} (1 + (1-2f)^{\lambda m / \log n})^m \\ &= \sum_{w=1}^{\lambda m / \log n} \binom{n}{w} \frac{1}{2^m} (1 + (1-2f)^w)^m + \\ &\quad \frac{q-1}{2^m} (1 + (1-2f)^{\lambda m / \log n})^m \\ &= A_n + B_n \end{aligned}$$

By Lemma 11, from our choice of  $\lambda$ , for any  $k > 0$ ,  $f = \frac{k \log^2 m}{m}$ , we have:

$$\lim_{n \rightarrow \infty} A_n = 0$$

For  $B_n$ , when  $n$  is sufficiently large, we have:

$$\begin{aligned} B_n &= \frac{q-1}{2^m} \left( 1 + \left( 1 - 2 \frac{k \log^2 m}{m} \right)^{\lambda m / \log n} \right)^m \\ &\leq 2^c \left( 1 + \exp \left( -2 \frac{k \log^2 m}{m} \frac{\lambda m}{\frac{1}{\beta} (\log m - \log \alpha)} \right) \right)^m \\ &\leq 2^c \left( 1 + \exp \left( -2 \frac{k \log m}{m} \lambda \beta m \right) \right)^m \\ &= 2^c (1 + m^{-2k\lambda\beta})^m \\ &\leq 2^c \exp(m^{-2k\lambda\beta} m) = 2^c \exp(m^{1-2k\lambda\beta}) \end{aligned}$$

for all  $\lambda > 1$ , if we choose  $k$  such that  $k > 1/(2\lambda\beta)$ , then  $1 - 2k\lambda\beta < 0$  and we have

$$\limsup_{n \rightarrow \infty} B_n \leq 2^c$$

Combining the two results, as long as we choose  $k > 1/(2\lambda\beta)$ , we have,

$$\limsup_{n \rightarrow \infty} A_n + B_n \leq 2^c$$

which implies that for all  $\gamma > 1$ , for sufficiently large  $n$ :

$$\epsilon(n, m, q, f) \leq \frac{A_n + B_n}{q-1} \leq \gamma \frac{2^c}{q-1}$$

Since  $\lambda < \log_2(1+\delta)/(1-\beta)$ , we need  $k$  to be larger than  $(1-\beta)/(2\beta \log_2(1+\delta))$ . Since  $\delta$  can be chosen arbitrarily from  $(0, 1)$ , it suffices to have:

$$k > \frac{1-\beta}{2\beta}$$

as for any such  $k$ , we can always find a  $\delta$  close enough to 1 such that the above condition is satisfied.  $\square$

**Proof of Theorem 2 (part 3).** This proof is almost exactly the same as that of Theorem 2 (part 2), following as a direct consequence of Lemma 12 along with Corollary 1 and Theorem 2 of Ermon et al. (2014).  $\square$