

# DISTRIBUTION OF MISSING SUMS IN SUMSETS

OLEG LAZAREV, STEVEN J. MILLER, AND KEVIN O'BRYANT

**ABSTRACT.** For any finite set of integers  $X$ , define its sumset  $X + X$  to be  $\{x + y : x, y \in X\}$ . In a recent paper, Martin and O'Bryant investigated the distribution of  $|A + A|$  given the uniform distribution on subsets  $A \subseteq \{0, 1, \dots, n - 1\}$ . They also conjectured the existence of a limiting distribution for  $|A + A|$  and showed that the expectation of  $|A + A|$  is  $2n - 11 + O((3/4)^{n/2})$ . Zhao proved that the limits  $m(k) := \lim_{n \rightarrow \infty} \mathbb{P}(2n - 1 - |A + A| = k)$  exist, and that  $\sum_{k \geq 0} m(k) = 1$ .

We continue this program and give exponentially decaying upper and lower bounds on  $m(k)$ , and sharp bounds on  $m(k)$  for small  $k$ . Surprisingly, the distribution is at least bimodal; sumsets have an unexpected bias against missing exactly 7 sums. The proof of the latter is by reduction to questions on the distribution of related random variables, with large scale numerical computations a key ingredient in the analysis. We also derive an explicit formula for the variance of  $|A + A|$  in terms of Fibonacci numbers, finding  $\text{Var}(|A + A|) \approx 35.9658$ . New difficulties arise in the form of weak dependence between events of the form  $\{x \in A + A\}$ ,  $\{y \in A + A\}$ . We surmount these obstructions by translating the problem to graph theory. This approach also yields good bounds on the probability for  $A + A$  missing a consecutive block of length  $k$ .

## CONTENTS

1. Introduction	2
1.1. Terminology and Theorems	3
1.2. Variance and Decay Rates of Missing Sums	4
1.3. Other types of random sets and the divot	6
2. Graph-Theoretic Framework	8
3. Variance of Missing Sums	9
4. Exponential Bounds	17
5. Approximating $\mathbb{P}(k + a_1, k + a_2, \dots, \text{ and } k + a_m \notin A + A)$	20
6. Consecutive Missing Sums	23
7. Bounds on $m(k)$ , $w(k)$ , $y(k)$ , and $z(k)$ for $k < 32$	27
7.1. Making the computation feasible, reliable, and verifiable	31
7.2. Obtaining $y(k)$ , $m(k)$ , and $w(k)$ from $z(k)$	32
8. Conjectures and Future Research	34
Appendix A. Data tables for distributions	36
References	40

---

*Date:* October 20, 2012.

2010 *Mathematics Subject Classification.* 11P99 (primary), 11K99 (secondary).

*Key words and phrases.* sumsets, uniformly random sumsets, Fekete's Lemma.

We thank the participants of the SMALL 2011 REU at Williams College for many enlightening conversations, and the referee for many helpful comments on an earlier draft. The first named author was supported by NSF grants DMS0850577 and Williams College; the second named author was partially supported by NSF grant DMS0970067. This research was supported, in part, under National Science Foundation Grants CNS-0958379 and CNS-0855217 and the City University of New York High Performance Computing Center.

## 1. INTRODUCTION

The central object of additive number theory [N, TV] is the sumset  $X + X$  of a set  $X$  of integers:

$$X + X := \{x_1 + x_2 : x_1, x_2 \in X\}. \quad (1.1)$$

Typically, the theory is concerned with extremal behavior, such as the structure of finite  $X$  when  $|X + X|/|X|$  is nearly minimal (Freiman's Theorem), or the possible densities of  $X$  when  $|X + X|/\binom{|X|+1}{2}$  is maximized (Sidon Sets). See [N, R] for surveys and [F, J] for examples.

Here we focus on *typical* behavior: for a randomly chosen set  $X$  of integers, what is the expected value and variance of  $|X + X|$ ? The answer of course depends on how  $X$  is chosen, and we focus our attention on sets taken uniformly from the  $2^n$  subsets of  $[0, n - 1]$ ; we denote intervals of integers as  $[a, b] := \{x \in \mathbb{Z} : a \leq x \leq b\}$  and such a random set as  $A$ . In §1.3 and §7.2 we discuss some variations on the manner of choosing a random set of natural numbers.

Other authors have considered aspects of typical behavior of sumsets. When Erdős and Rényi [ER] first applied the probabilistic method to number theory, they observed that with probability 1, a uniformly random subset  $C$  of  $\mathbb{N}$  will have  $C + C = \mathbb{N} \setminus F$  for some *finite* set  $F$ , but made no effort to explore  $F$  further. The present work concerns itself with properties of the set

$$F_n := [0, 2n - 2] \setminus (A + A),$$

with  $A$  as above. We prove the existence of

$$\lim_{n \rightarrow \infty} \mathbb{E}[|F_n|^r]$$

for every  $r \geq 1$ , give upper and lower bounds on

$$\mathbb{P}(|F_n| = k)$$

for small  $k$ , large  $n$ , and also as  $k \rightarrow \infty$ , and also bound

$$\mathbb{P}(\{a_1, a_2, \dots, a_k\} \subseteq F_n).$$

Our work is usually quantitatively effective, and we report numerical estimates throughout.

The key obstacle to finding the limiting distribution of  $|F_n|$  is the dependence between different elements occurring or not occurring in  $A + A$ . For example,  $3 \notin A + A$  and  $7 \notin A + A$  are dependent events since both are affected by whether  $2 \in A$ . We develop a graph theoretic framework which makes it much easier to analyze the dependence between such events and to develop bounds that incorporate the dependence. It is possible to avoid this framework, but doing so makes both notation and the underlying issues less clear.

Graph theory has been used in additive number theory before. For example, Plünnecke (see the description in [R]) uses graph theory to estimate the size of  $k$ -fold sumsets in terms of  $|A|$  and  $|A + A|$ , Alon and Erdős [AE] use hypergraphs to study Sidon sets, and Gilbert [G] on the Erdős-Turan conjecture. Our use of graph theory seems to be different from these as we investigate the size of  $A + A$  for typical  $A$ , without reference to the size of  $A$  itself.

The next subsection of this introduction sets up our notation and states our main results. The last two subsections provide more motivation and indicate the nature of our proofs and computations. In §2, we develop a graph theoretic framework for handling the dependencies between events like  $\{a_1 \in F\}$  and  $\{a_2 \in F\}$ . In §3, we find an explicit formula for the limit of the variance of  $|F|$  and prove Theorem 1.5, stated below. In §4, we prove the exponential bounds for Theorem 1.2. In §5, we find the probability of missing certain configurations and prove Theorem 1.6, while in §6 we discuss consecutive missing elements and prove Theorem 1.7 and Theorem 1.8. We return to the

problem of explicit bounds on  $\mathbb{P}(|F_n| = k)$  for small  $k$  and the existence of a limiting distribution for  $|F_n|$  in §7. Finally in §8, we discuss some problems for future research and how the graph theoretic framework may be applied to such problems.

**Remark 1.1.** *Many of the questions in this paper grew out of studying the difference in size between the sumset  $A + A$  and the difference set  $A - A$ . As addition is commutative and subtraction is not, it is natural to expect the difference set of a typical  $A$  drawn uniformly from  $\{0, 1, \dots, n\}$  to be larger than the sumset. Though numerical exploration and heuristics suggested that almost all sets should give rise to more differences, Martin and O’Bryant [MO] proved that a small but positive percentage are sum-dominant. The percentage is quite small, around  $4.5 \cdot 10^{-4}$  [Z]. Understanding the structure of  $A + A$ , in particular when and what sums are missing, has motivated much of the theoretical and numerical work in the field. For other directions, see [HM] for results on non-uniform models or [ILMZ] for multiple comparisons and summands.*

**1.1. Terminology and Theorems.** The main characteristic of  $A + A$  is that it is almost full. Martin and O’Bryant [MO] proved that

$$\mathbb{E}[|A + A|] = 2n - 1 - 10 + O((3/4)^{n/2}). \quad (1.2)$$

Since typical sumsets are almost full, it is more natural to investigate the number of missing sums, which is why we write the above as  $2n - 1$  minus 10. As noted in [MO], sumsets are almost full because middle elements have many representations as a sum of two elements of  $[0, n - 1]$ ; each  $i \in [0, 2n - 2]$  has roughly  $n/4 - |n - i|/4$  representations.

We set

$$\begin{aligned} M_{[0, n-1]} &:= |[0, 2n - 2] \setminus (A + A)| = 2n - 1 - |A + A|, \\ m_n(k) &:= \mathbb{P}(M_{[0, n-1]} = k), \\ m(k) &:= \lim_{n \rightarrow \infty} m_n(k). \end{aligned} \quad (1.3)$$

A special case of Zhao’s theorem [Z] is that  $m(k)$  is well-defined, strictly positive, and that  $\sum_{k=0}^{\infty} m(k) = 1$ , so that we can think of  $m(k)$  as defining a distribution on  $\mathbb{N}$ . Thus, we can speak of “the probability that a large finite set  $X$  has a sumset that misses exactly 17 elements” and mean something sensible. Zhao’s work is numerically impractical and did not give reasonable upper bounds on  $m(k)$ ; we do that in §7, where we also reprove Zhao’s results in this easier setting. See Figure 1 for the experimental estimates and rigorous bounds on  $m(k)$  for  $0 \leq k < 32$ .

The result (1.2) above implies that

$$\lim_{n \rightarrow \infty} \mathbb{E}[M_{[0, n-1]}] = 10.$$

Equivalently, in light of Zhao’s work,  $\sum_{k=0}^{\infty} km(k) = 10$ . To this, we add the following results. Let  $\phi := (1 + \sqrt{5})/2$ , the golden ratio.

**Theorem 1.2.** *Let  $n > 5k$ . Then*

$$2^{-k/2} \ll m_n(k) \ll (\phi/2)^k, \quad (1.4)$$

where the implied constants are independent of  $k$  and  $n$ .

Note that  $2^{-1/2} \approx 0.707$  and  $\phi/2 \approx 0.809$ , so that bounds provided by Theorem 1.2 are reasonably close. We suspect, based on numerical data, that the following conjecture represents the truth of the matter, and perhaps even  $\lambda = \sqrt{\phi - 1}$ .

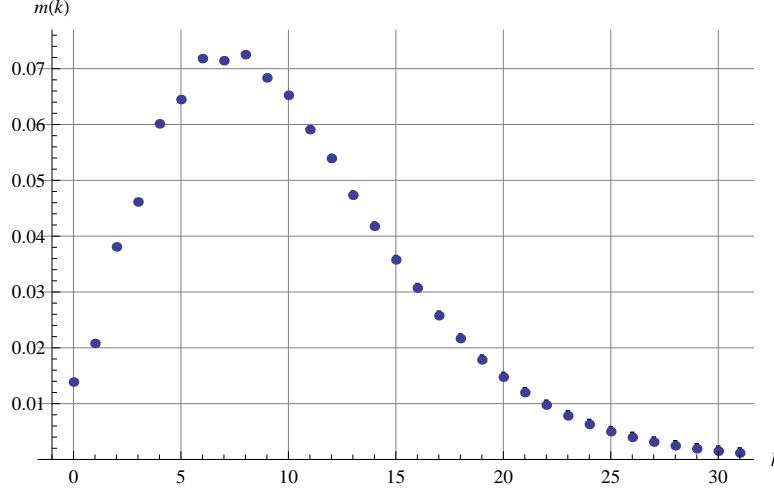


FIGURE 1. Experimental values of  $m(k)$ , with vertical bars depicting the values allowed by our rigorous bounds. In most cases, the allowed interval is smaller than the dot indicating the experimental value. The data comes from generating  $2^{28}$  sets uniformly forced to contain 0 from  $[0, 256)$ ; see §7.2 for details of the calculation.

**Conjecture 1.3.** *There exists  $\lambda$  such that for any  $\epsilon > 0$ ,*

$$(\lambda - \epsilon)^k \ll_{\epsilon} m(k) \ll_{\epsilon} (\lambda + \epsilon)^k. \quad (1.5)$$

*From numerical data,  $\lambda \approx 0.78$ .*

The exponential bounds of Theorem 1.2 already imply that the  $r^{\text{th}}$  moment remains bounded for any  $r \geq 1$ .

**Corollary 1.4.** *The limit of the  $r^{\text{th}}$  moment of  $M_{[0, n-1]}$ ,*

$$\lim_{n \rightarrow \infty} \mathbb{E} [M_{[0, n-1]}^r], \quad (1.6)$$

*exists and is finite.*

**Theorem 1.5.** *The limit*

$$\lim_{n \rightarrow \infty} \text{Var} (M_{[0, n-1]}) \quad (1.7)$$

*exists and is about 35.9658, as these are the first digits of its decimal expansion. This limit can be written as the following convergent series with exponential decay:*

$$\lim_{n \rightarrow \infty} \text{Var} (M_{[0, n-1]}) = 4 \lim_{n \rightarrow \infty} \sum_{i < j < n} \mathbb{P}(i \text{ and } j \notin A + A) - 40. \quad (1.8)$$

Note that “ $i$  and  $j \notin A + A$ ” is meant to be parsed as “ $(i \notin A + A) \text{ AND } (j \notin A + A)$ ”.

**1.2. Variance and Decay Rates of Missing Sums.** The bounds in Theorem 1.2 are due to formulas for probabilities of events such as

$$\mathbb{P}(a_1, a_2, \dots, \text{ and } a_m \notin A + A), \quad (1.9)$$

by which we mean the probability that all of  $a_1, a_2, \dots, a_m$  are in the complement of  $A + A$ . This represents the probability that a particular configuration is not in  $A + A$ . As long as  $n > a_m$ ,

there is no dependence on  $n$  since this probability just depends on  $[0, a_m] \cap A$ . We therefore can assume that  $A \subseteq [0, a_m]$ . Formulas for such probabilities are also important for finding the moments of  $M_{[0, n-1]}$ . For example, to find the expectation of  $|A + A|$ , [MO] find an exact formula for  $\mathbb{P}(k \notin A + A)$ , which is approximately

$$\mathbb{P}(k \notin A + A) = \Theta((3/4)^{k/2}), \quad (1.10)$$

where we say  $g(n) = \Theta(f(n))$  if there exist constants  $C_1, C_2$  such that for all  $n$

$$C_1 f(n) \leq g(n) \leq C_2 f(n). \quad (1.11)$$

Similarly, to find the variance, we can study  $\mathbb{P}(i \text{ and } j \notin A + A)$  as seen from the series expansion in (1.8). In Proposition 3.5, we find an exact formula for this probability and in Corollary 3.6, we show that for fixed  $m$  we have the following approximation:

$$\mathbb{P}(k \text{ and } k + m \notin A + A) = \Theta((\phi/2)^k). \quad (1.12)$$

The implied constants in (1.12) depend significantly on  $m$  and in Corollary 3.6, we also find these constants.

Note that both (1.10) and (1.12) are exponential in  $k$ . In fact, we prove that in general such probabilities are approximately exponential in  $k$ .

**Theorem 1.6.** *For any fixed  $a_1, \dots, a_m$ , there exists  $\lambda_{a_1, \dots, a_m}$  such that*

$$\mathbb{P}(k + a_1, k + a_2, \dots, \text{ and } k + a_m \notin A + A) = \Theta(\lambda_{a_1, \dots, a_m}^k), \quad (1.13)$$

where the implied constants depend on  $a_1, \dots, a_m$  but not  $k$ .

The fact that  $\mathbb{P}(k + a_1, k + a_2, \dots, \text{ and } k + a_m \notin A + A)$  is approximately exponential supports Conjecture 1.3 that the distribution of missing sums is approximately exponential.

For the particular configuration  $a_1 = 1, a_2 = 2, \dots, a_m = m$ , the case of consecutive missing elements, we can approximate  $\lambda_{a_1, \dots, a_m}$  well as seen in the following theorem.

**Theorem 1.7.** *For any  $k, m$*

$$\left(\frac{1}{2}\right)^{(k+m)/2} \ll \mathbb{P}(k + 1, k + 2, \dots, \text{ and } k + m \notin A + A) \ll \left(\frac{1}{2}\right)^{(k+m)/2} (1 + \epsilon_m)^k, \quad (1.14)$$

with  $\epsilon_m \rightarrow 0$  as  $m \rightarrow \infty$ . To be more precise, the exact form of upper bound is  $(1/2)^{(k+m)/2} 2^{k/m}$ . This implies that

$$\lambda_{0, 1, 2, \dots, m} \rightarrow \left(\frac{1}{2}\right)^{1/2} \quad (1.15)$$

as  $m \rightarrow \infty$ .

As we will see in the proof of Theorem 1.2, the lower bound  $(1/2)^{(k+m)/2}$  is essentially the probability of missing the first  $k + m$  elements in  $A + A$ . By Theorem 1.7, we have that for large  $m$ ,  $\mathbb{P}(k + 1, k + 2, \dots, \text{ and } k + m \notin A + A)$  is also approximately  $(1/2)^{(k+m)/2}$ . This means that for large  $m$ , essentially the only way to miss  $m$  consecutive elements in  $A + A$  starting at  $k + 1$  is through the trivial way - namely missing all of the first  $k + m$  elements of  $A + A$ .

Theorem 1.7 is in fact a special case of the following inequality.

**Theorem 1.8.** *For  $\lambda_{a_1, \dots, a_m}$  with  $0 \leq a_1 < \dots < a_m$ ,*

$$\lambda_{a_1, \dots, a_m} \leq \mathbb{P}(A, B \subseteq [0, \lfloor a_m/2 \rfloor] \mid a_1, \dots, a_m \notin A + B)^{\frac{1}{a_m+2}}. \quad (1.16)$$

where  $A, B$  are two independently chosen sets.

**1.3. Other types of random sets and the divot.** Figure 1 shows a surprising phenomenon: experimentally,

$$m(7) < m(6) < m(8).$$

That is, a random subset of  $[0, 10^{10}]$  is more likely to have a sumset missing 6 (or 8) elements than one missing 7 elements. That is, the distribution of  $M_{[0, n-1]}$  appears to be bimodal for large  $n$ . We have made a massive computation (details in §7), looping over  $2^{43}$  sets and using only 64-bit integer arithmetic, that lead to the following bounds:

$$0.07177 < m(6) < 0.07202, \quad 0.07138 < m(7) < 0.7170, \quad 0.07243 < m(8) < 0.07282. \quad (1.17)$$

We note that our bounds are actually in the form

$$\frac{107418021089142422011644549535908507304608994344051}{1496577676626844588240573268701473812127674924007424} < m(6) \\ m(6) < \frac{620778536995376440633741122321102716502820362028980739}{8620287417370624828265702027720489157855407562282762240};$$

we hope the reader will excuse our preference for reporting equivalent decimals, rounded in the proper directions to maintain truth.

Closer inspection of Figure 1 also reveals an apparent parity effect:

$$m(2k) + m(2k + 2) > 2m(2k + 1).$$

Here are two plausible explanations for this. The first is that  $M_{[0, n-1]}$  is essentially the sum of two iidrvs: the number of missing sums in  $[0, n - 1]$  and in  $[n, 2n - 2]$ . For any two iidrvs  $X_1, X_2$  taking integer values,  $\mathbb{P}(X_1 + X_2 \text{ even}) \geq \mathbb{P}(X_1 + X_2 \text{ odd})$ , as the calculation comes down to  $x^2 + y^2 \geq 2xy$ . Another parity effect is observed on the ends: as soon as  $0 \notin A$ , then both 0 and 1 are not in  $A + A$ . Thus, on the ends,  $A + A$  always misses an even number of sums.

To compensate for these observations, it is necessary to consider the connections between different ways of selecting a random set. We consider uniformly selecting subsets of  $[0, n - 1]$ , subsets of  $[0, n]$  with diameter  $n$ , subsets of  $\mathbb{N}$ , and subsets of  $\mathbb{N}$  that contain 0. We lay out our notation as follows:

set	setting	condition	missing sums	$\mathbb{P}(\text{missing } k \text{ sums})$
$A$	$[0, n - 1]$	$\emptyset$	$M_{[0, n-1]} := 2n - 1 -  A + A $	$m_n(k)$
$B$	$[0, n]$	$\{0, n\} \subseteq B$	$M_{[0, n] \{0, n\}} := 2n + 1 -  B + B $	$w_n(k)$
$C$	$\mathbb{N}$	$\emptyset$	$M_{\mathbb{N}} :=  \mathbb{N} \setminus (C + C) $	$y(k)$
$D$	$\mathbb{N}$	$0 \in D$	$M_{\mathbb{N} \{0\}} :=  \mathbb{N} \setminus (D + D) $	$z(k)$

Additionally, we set  $m(k) := \lim_{n \rightarrow \infty} m_n(k)$  and  $w(k) := \lim_{n \rightarrow \infty} w_n(k)$ .

Our first parity-effect observation essentially boils down to

$$m_n(k) \rightarrow \sum_{i=0}^k y(i)y(k-i), \quad (1.18)$$

a rigorous exposition of this can be found in [I] and is sketched in §7.2. The second observation and Bayes' Theorem leads us to

$$y(k) = \sum_{i=0}^{\infty} \mathbb{P}(\min C = i) \mathbb{P}(|[2i, \infty) \setminus (C + C)| = k - 2i) = \sum_{i=0}^{\lfloor k/2 \rfloor} 2^{-(i+1)} z(k - 2i). \quad (1.19)$$

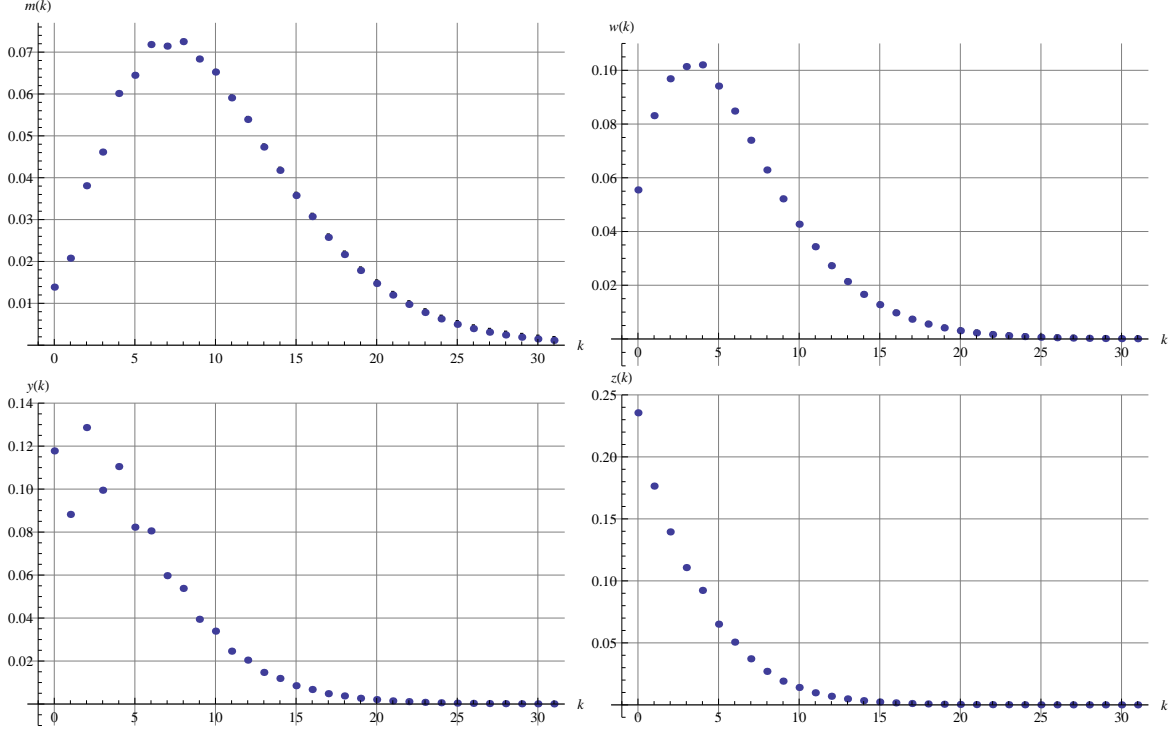


FIGURE 2. Experimental values of  $m(k)$ ,  $w(k)$ ,  $y(k)$ ,  $z(k)$ , with vertical bars depicting the values allowed by our rigorous bounds. See §7 for details.

Similarly to (1.18), one can prove that

$$w_n(k) \rightarrow \sum_{i=0}^k z(i)z(k-i). \quad (1.20)$$

Thus, all four distributions can be understood in terms of  $z(k)$ . Experiments and our bounds (see Figure 2 for small values of  $k$ ) indicate that  $M_{\mathbb{N} \setminus \{0\}}$  has an approximately geometric distribution, and exhibits no obvious parity effect. Computationally, we focus on bounding  $z$  and then allow this to determine bounds on  $m$ ,  $w$  and  $y$ .

We bound  $z$  by conditioning on  $I := D \cap [0, 44)$ , and loop over all  $2^{43}$  possible values of  $I$  (a priori,  $0 \in I$ ). For each  $I \subseteq [0, 44)$ , we explicitly know  $(D+D) \cap [0, 44)$ , we have much information concerning  $(D+D) \cap [44, 88)$ , and theoretically  $(D+D) \cap [88, \infty)$  is  $[88, \infty)$  with high probability. This allows us to give reasonable upper and lower bounds on  $\mathbb{P}(M_{\mathbb{N} \setminus \{0\}} = k \mid D \cap [0, 44) = I)$  for each  $I$ .

If we suppose that  $M_{\mathbb{N} \setminus \{0\}}$  is exactly geometric with parameter  $\lambda$  (i.e., set  $z(k) = (1-\lambda)\lambda^k$ ) and define  $y(k)$  and  $m(k)$  using (1.18) and (1.19), we find that the distribution of  $M_{\mathbb{N} \setminus \{0\}}$  would be bimodal with a divot at  $k = 7$  only for the narrow parameter range  $0.756 < \lambda < 0.771$ . The best-squares fit for  $\lambda$  is 0.765. If we suppose that  $M_{\mathbb{N} \setminus \{0\}}$  has a Poisson distribution, i.e.,  $z(k) = \lambda^k e^{-\lambda} / k!$ , we find that there are no  $\lambda$  whatsoever that give a bimodal distribution with divot at  $k = 7$ .

This implies that the divot's existence relies not only on the above observations but also on the specific values of  $z_k$  for small values. We note that  $z_4$  in particular is larger than the geometric

model predicts; more than half of the least-squares error is from  $z_4$ . The rigorous bounds we give also show this bias towards 4, though we currently have no understanding as to why this is the situation.

**Theorem 1.9.** *The limits defining  $m(k)$  and  $w(k)$  are well-defined, positive, and  $\sum_{k=0}^{\infty} m(k) = \sum_{k=0}^{\infty} w(k) = 1$ . Rigorous bounds on  $m(k)$ ,  $w(k)$ ,  $y(k)$  and  $z(k)$  for  $0 \leq k < 32$  are given in Appendix A. In particular,  $m(7) < m(6) < m(8)$ .*

## 2. GRAPH-THEORETIC FRAMEWORK

We first develop a graph-theoretic framework to study dependent random variables and calculate probabilities like  $\mathbb{P}(a_1, \dots, \text{and } a_m \notin A + A)$ . Note that for odd  $i$

$$\{i \notin A + A\} = \{(0 \notin A \text{ or } i \notin A) \text{ and } \dots \text{ and } ((i-1)/2 \notin A \text{ or } (i+1)/2 \notin A)\}, \quad (2.1)$$

and for even  $i$

$$\{i \notin A + A\} = \{(0 \notin A \text{ or } i \notin A) \text{ and } \dots \text{ and } (i/2 - 1 \notin A \text{ or } i/2 + 1 \notin A) \text{ and } i/2 \notin A\}. \quad (2.2)$$

Therefore for distinct  $i$  the events  $\{i \notin A + A\}$  are dependent as both depend on conditions on  $A$  like  $\{0 \notin A\}$ .

For example, the conditions on  $A$  necessary for  $\{3 \text{ and } 7 \notin A + A\}$  are

$$\begin{array}{ll} i = 3 : & 0 \text{ or } 3 \notin A \\ & \text{and } 1 \text{ or } 2 \notin A \\ j = 7 : & 0 \text{ or } 7 \notin A \\ & \text{and } 1 \text{ or } 6 \notin A \\ & \text{and } 2 \text{ or } 5 \notin A \\ & \text{and } 3 \text{ or } 4 \notin A. \end{array} \quad (2.3)$$

Since the two lists have integers in common, there is dependence between the events  $\{3 \notin A + A\}$  and  $\{7 \notin A + A\}$ .

We construct a graph to represent the dependencies between the random variables. We call this graph the *condition graph* for the probability. We construct the condition graph for  $\mathbb{P}(a_1, \dots, \text{and } a_m \notin A + A)$ , where  $a_1 < \dots < a_m$ , in the following way:

- (1) For every integer in  $[0, a_m]$ , add a vertex labeled with that integer.
- (2) Add an edge between two vertices labeled with  $i$  and  $j$  if  $i + j = a_k$  for some  $1 \leq k \leq m$ .

See Figure 3 for the condition graph for  $\mathbb{P}(3 \text{ and } 7 \notin A + A)$ .

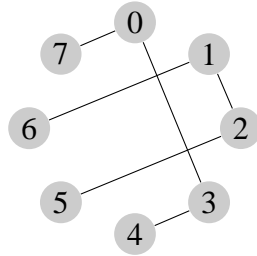


FIGURE 3. Condition Graph for  $\mathbb{P}(3 \text{ and } 7 \notin A + A)$ .

By construction, we have a one-to-one correspondence between edges and conditions and vertices and integers in  $[0, a_m]$ . For example, the edge between vertices labeled with 1 and 6 represents the condition that  $1 \text{ or } 6 \notin A$ , which is one of the conditions necessary for  $7 \notin A + A$  in (2.3). For



each condition, we need to pick at least one element to exclude from  $A$ . Therefore in the condition graph, for each edge we need to pick at least one of its vertices. That is, we need to pick a vertex cover (recall a vertex cover of a graph is a set of vertices such that each edge is incident to at least one vertex in the set). Using this method, we get the following lemma.

**Lemma 2.1.**  $\mathbb{P}(a_1, \dots, \text{ and } a_m \notin A + A)$  equals the probability that we chose a vertex cover for the condition graph.

Note that when we pick vertices in the condition graph for our vertex cover, we are picking to exclude those vertices from  $A$ . For example, note that the vertices 7, 0, 4 and 6, 2 form a vertex cover for the condition graph of  $\mathbb{P}(3 \text{ and } 7 \notin A + A)$  in Figure 3. Then if 7, 0, 4, 6, 2  $\notin A$ , then 3 and 7  $\notin A + A$  since all conditions in (2.3) are met.

Finally, note that when we calculate the probability of choosing a vertex cover for the condition graph, we no longer need to consider a labeled graph. This is because vertices represent elements of  $A$ , and since each element of  $A$  is equally likely to be chosen (as  $A$  is chosen uniformly randomly), we do not need to differentiate between different elements.

### 3. VARIANCE OF MISSING SUMS

We now use the graph-theoretic framework from the previous section to prove Theorem 1.5 and find the variance.

We first note that the result of [MO] in (1.2) is really that

$$\mathbb{E} [M_{[0, n-1]}(A)] = \sum_{0 \leq i \leq 2n-2} \mathbb{P}(i \notin A + A) = 10 + O((3/4)^{n/2}). \quad (3.1)$$

Since

$$\text{Var} (M_{[0, n-1]}(A)) = \mathbb{E} [M_{[0, n-1]}(A)^2] - (\mathbb{E} [M_{[0, n-1]}(A)])^2 \quad (3.2)$$

and we know  $\mathbb{E} [M_{[0, n-1]}(A)]$  from (3.1), to find the variance we just need to determine  $\mathbb{E} [M_{[0, n-1]}(A)^2]$ , which equals the following:

$$\begin{aligned} \mathbb{E} [M_{[0, n-1]}(A)^2] &= \frac{1}{2^n} \sum_{A \subseteq [0, n-1]} |\{\text{missing sums in } A + A\}|^2 \\ &= \frac{1}{2^n} \sum_{A \subseteq [0, n-1]} \sum_{\substack{0 \leq i, j \leq 2n-2 \\ i, j \notin A + A}} 1 \\ &= \frac{1}{2^n} \sum_{0 \leq i, j \leq 2n-2} \sum_{\substack{A \subseteq [0, n-1] \\ i, j \notin A + A}} 1 \\ &= \sum_{0 \leq i, j \leq 2n-2} \mathbb{P}(A \subseteq [0, n-1] \mid i \text{ and } j \notin A + A) \\ &= 2 \sum_{0 \leq i < j \leq 2n-2} \mathbb{P}(i \text{ and } j \notin A + A) + \sum_{0 \leq i \leq 2n-2} \mathbb{P}(i \notin A + A). \quad (3.3) \end{aligned}$$

Combining (3.2), (3.1), and (3.3), we get

$$\text{Var} (M_{[0, n-1]}(A)) = 2 \sum_{0 \leq i < j \leq 2n-2} \mathbb{P}(i \text{ and } j \notin A + A) - 90 + O((3/4)^{n/2}). \quad (3.4)$$

We first simplify the sum over  $i, j$ . Note that if  $i, j < n$ , then

$$\mathbb{P}(i \text{ and } j \notin A + A) = \mathbb{P}(2n - 2 - i \text{ and } 2n - 2 - j \notin A + A), \quad (3.5)$$

and so

$$\sum_{0 \leq i < j < n} \mathbb{P}(i \text{ and } j \notin A + A) = \sum_{n \leq i < j \leq 2n-2} \mathbb{P}(i \text{ and } j \notin A + A). \quad (3.6)$$

Also, note that if  $i < n/2$  and  $j > 3n/2$ , then  $\{i \notin A + A\}$  and  $\{j \notin A + A\}$  are independent. This is because  $\{i \notin A + A\}$  depends only on  $[0, i] \cap A$  and  $\{j \notin A + A\}$  depends only on  $[j - n + 1, n - 1] \cap A$  and if  $i < n/2$  and  $j > 3n/2$ , these sets are disjoint. Therefore for such  $i, j$ , we have

$$\mathbb{P}(i \text{ and } j \notin A + A) = \mathbb{P}(i \notin A + A) \mathbb{P}(j \notin A + A). \quad (3.7)$$

Finally note that if  $n/2 \leq i < n$  or  $n \leq j \leq 3n/2$ , then

$$\mathbb{P}(i \text{ and } j \notin A + A) = O((3/4)^{n/4}) \quad (3.8)$$

by (1.10). Therefore

$$\begin{aligned} & \sum_{i < n, n \leq j} \mathbb{P}(i \text{ and } j \notin A + A) \\ &= \sum_{i < n/2 \text{ and } 3n/2 < j} \mathbb{P}(i \text{ and } j \notin A + A) + \sum_{n/2 \leq i < n \text{ or } n \leq j \leq 3n/2} \mathbb{P}(i \text{ and } j \notin A + A) \\ &= \sum_{i < n/2, 3n/2 < j} \mathbb{P}(i \text{ and } j \notin A + A) + O(n^2(3/4)^{n/4}) \\ &= \left( \sum_{i < n/2} \mathbb{P}(i \notin A + A) \right) \cdot \left( \sum_{3n/2 < j \leq 2n-2} \mathbb{P}(j \notin A + A) \right) + O(n^2(3/4)^{n/4}) \\ &= (5 + O((3/4)^{n/4})) \cdot (5 + O((3/4)^{n/4})) + O(n^2(3/4)^{n/4}) \\ &= 25 + O(n^2(3/4)^{n/4}), \end{aligned} \quad (3.9)$$

where we use (3.1) and (3.5) to get the second to last equality. Combining (3.6) and (3.9), we have

$$\begin{aligned} & \sum_{0 \leq i < j \leq 2n-2} \mathbb{P}(i \text{ and } j \notin A + A) \\ &= \sum_{0 \leq i < j < n} \mathbb{P}(i \text{ and } j \notin A + A) + \sum_{n \leq i < j \leq 2n-2} \mathbb{P}(i \text{ and } j \notin A + A) + \sum_{i < n, n \leq j} \mathbb{P}(i \text{ and } j \notin A + A) \\ &= 2 \sum_{0 \leq i < j \leq n-1} \mathbb{P}(i \text{ and } j \notin A + A) + 25 + O(n^2(3/4)^{n/4}), \end{aligned} \quad (3.10)$$

and so by (3.4)

$$\text{Var} \left( M_{[0, n-1]}(A) \right) = 4 \sum_{0 \leq i < j \leq n-1} \mathbb{P}(i \text{ and } j \notin A + A) - 40 + O(n^2(3/4)^{n/4}). \quad (3.11)$$

Therefore to find the variance, we just need to study  $\mathbb{P}(i \text{ and } j \notin A + A)$  for  $i < j < n$ .

Since the other cases are handled similarly, we only present the details for the case when  $i$  and  $j$  are both odd. By Lemma 2.1, we just need to study the condition graph for  $\mathbb{P}(i \text{ and } j \notin A + A)$ .

Recall that we already found the condition graph for  $\mathbb{P}(3 \text{ and } 7 \notin A + A)$  in Figure 3. After untangling this graph, we see that it really consists of two components, as seen in Figure 4.



FIGURE 4. Untangled condition graph for  $\mathbb{P}(3 \text{ and } 7 \notin A + A)$ .

Also note that each component is a *segment graph*, a graph that consists of a sequence of vertices such that each vertex is connected only to the vertices to its immediate left and right. A similar situation holds in general, as seen by the following proposition.

**Proposition 3.1.** *The condition graph for  $\mathbb{P}(i \text{ and } j \notin A + A)$  has components that are segment graphs.*

*Proof.* The condition graph for  $\mathbb{P}(i \text{ and } j \notin A + A)$  has vertices with degree less than or equal to 2; if the vertex is labeled with  $\ell$ , it can only be connected to vertices labeled  $i - \ell$  or  $j - \ell$  (if such vertices exist).

Furthermore, there are no cycles in the condition graph. Suppose there is a cycle in the condition graph. Consider the vertex in the cycle with the maximum label  $\ell$  and consider the vertices around this vertex. Each of these vertices must have exactly two neighbors and so we have the following situation as seen in Figure 5.



FIGURE 5. Vertices around a labeled vertex  $\ell$ .

Notice that  $\ell + j - i > \ell$  since  $j > i$ . Therefore,  $\ell$  is not the maximum label, which is a contradiction and proves that we cannot have a cycle. Thus all components are trees with all vertices having degree less than or equal to 2, implying that all are segment graphs.  $\square$

Since labels in different components are distinct and there are no edges between different components, each component is independent. That is, the probability of getting a vertex cover for the entire graph is the product of the probability of getting vertex covers for each component. In this way, we just need to find the probability of getting a vertex cover for each component. To do this, we find the number of vertex covers for an arbitrary segment graph, which we do in the following proposition.

**Proposition 3.2.** *The number of vertex covers  $g(n)$  for a segment graph with  $n$  vertices satisfies  $g(n) = F_{n+2}$ , where  $F_k$  is the  $k^{\text{th}}$  Fibonacci number.*

*Proof.* There are two cases: the first vertex of the segment graph is in the vertex cover, or it is not. If the first vertex is in the cover, then the first edge already has one of its vertices picked. Therefore we just need a vertex cover for the subgraph with  $n - 1$  vertices that follows the first edge, and by definition there are  $g(n - 1)$  such covers. If the first vertex is not in the cover, then the second vertex must be the cover since the first edge must have one of its vertices chosen. Since the second vertex is now in the cover, then the second edge automatically has one of its vertices in the cover. Therefore we just need a vertex cover for the subgraph with  $n - 2$  vertices that follows the second

edge, and by definition there are  $g(n-2)$  such vertex covers. Therefore, we have the Fibonacci recursive relationship  $g(n) = g(n-1) + g(n-2)$ . As  $g(2) = 3 = F_4$  and  $g(3) = 5 = F_5$ , these initial conditions and the recurrence imply  $g(n) = F_{n+2}$ , completing the proof.  $\square$

Therefore, we have

$$\mathbb{P}(\text{chose a vertex cover for a segment graph with } n \text{ vertices}) = \frac{F_{n+2}}{2^n}. \quad (3.12)$$

Returning to our example with 3 and 7, we note that since the condition graph in this case consists of two segment graph components each of length 4, we have

$$\mathbb{P}(3 \text{ and } 7 \notin A + A) = \frac{F_{4+2}}{2^4} \cdot \frac{F_{4+2}}{2^4} = \frac{1}{4}, \quad (3.13)$$

where we can multiply the probabilities by the independence of the components.

In general, as the condition graph may have many components we must find how many segment graph components there are in the entire graph for  $\mathbb{P}(i \text{ and } j \notin A + A)$ .

**Proposition 3.3.** *There are  $(j-i)/2$  segment graph components for the graph of  $\mathbb{P}(i \text{ and } j \notin A + A)$ .*

*Proof.* Note that in total  $j+1$  vertices are used in the graph of  $\mathbb{P}(i \text{ and } j \notin A + A)$ ; since  $\{i \text{ and } j \notin A + A\}$  depends just on  $A \cap [0, j]$ , the graph uses exactly the integers in  $[0, j]$ . Also note that each component must end with a vertex labeled by an integer greater than  $i$ . If a component ends with a vertex labeled by  $\ell \leq i$ , then it can be connected to two other vertices  $i - \ell$  and  $j - \ell$ . Remember that we are assuming  $i$  and  $j$  are odd (the other cases are similar). As they are odd,  $i - \ell \neq \ell$  and  $j - \ell \neq \ell$  and so  $i - \ell, j - \ell, \ell$  are all distinct. Since  $\ell$  is connected to two other vertices, it cannot be an end vertex. Therefore, each end vertex is labeled by some integer in  $[i+1, j]$ . Also note that each of these integers must be end vertex since it cannot be used to add up to  $i$ . Therefore, the set  $[i+1, j]$  coincides with the set of end vertices and since each component has two end vertices with distinct labels, there are  $(j-i)/2$  components.  $\square$

We also need to find the length of each component. Fortunately, there are only two possible component lengths for the graph of  $\mathbb{P}(i \text{ and } j \notin A + A)$ , as seen by the following lemma.

**Proposition 3.4.** *The length of each segment graph component for the graph of  $\mathbb{P}(i \text{ and } j \notin A + A)$  is always either*

$$2 \left\lceil \frac{i+1}{j-i} \right\rceil \quad \text{or} \quad 2 \left\lceil \frac{i+1}{j-i} \right\rceil + 2. \quad (3.14)$$

*Proof.* First note that the difference between a given vertex and another vertex that is two edges away is  $j-i$ . This is because the sum of the vertices that share an edge alternates between  $i$  and  $j$ , so that we have segments of the form given in Figure 6. The difference between  $j-x$  and  $i-x$  is  $j-i$  as needed.

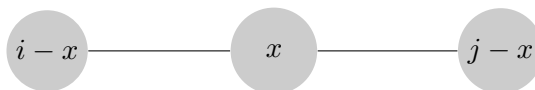


FIGURE 6. Difference between every other vertex.

Now note that these differences can be used to determine the size of each component. Suppose the end vertex of a segment graph is  $m$ . Since we decrease by  $j - i$  for every two vertices and since we only use non-negative integers, there can only be

$$\left\lfloor \frac{m}{j-i} \right\rfloor + 1 = \left\lceil \frac{m+1}{j-i} \right\rceil \quad (3.15)$$

decreases. Since we decrease once for every two vertices, we have that the length is twice the number of decreases. Therefore the length is

$$2 \left\lceil \frac{m+1}{j-i} \right\rceil. \quad (3.16)$$

From Proposition 3.3, we also know that the end vertex  $m$  of each segment graph satisfies  $i < m \leq j$ . Therefore, the length of each segment graph is always

$$2 \left\lceil \frac{i+1}{j-i} \right\rceil \quad \text{or} \quad 2 \left\lceil \frac{i+1}{j-i} \right\rceil + 2, \quad (3.17)$$

as desired.  $\square$

For simplicity, we denote the first of the two values in (3.14) by  $q$  and the second by  $q + 2$ . We must find the number of components with size  $q$  and  $q + 2$ . Suppose there are  $r$  components of size  $q$  and  $r'$  components of size  $q + 2$ . Then conditions on the number of components from Proposition 3.3 and the length of each component from Proposition 3.4 gives us the following two equations:

$$\begin{aligned} qr + (q+2)r' &= j+1 \\ r + r' &= \frac{j-i}{2}. \end{aligned} \quad (3.18)$$

Solving these equations for  $r, r'$  in terms of  $q$  gives

$$\begin{aligned} r &= \frac{1}{2} \left( \frac{j-i}{2} q - (i+1) \right) = \frac{1}{2} \left( (j-i) \left\lceil \frac{i+1}{j-i} \right\rceil - (i+1) \right) \\ r' &= \frac{1}{2} \left( j+1 - \frac{j-i}{2} q \right) = \frac{1}{2} \left( j+1 - (j-i) \left\lceil \frac{i+1}{j-i} \right\rceil \right). \end{aligned} \quad (3.19)$$

Therefore, again by independence of components, we have for odd  $i, j$  that

$$\mathbb{P}(i \text{ and } j \notin A + A) = \frac{1}{2^{j+1}} F_{q+2}^r F_{q+4}^{r'} \quad (3.20)$$

with  $q, r, r'$  as given in (3.17) and (3.19). Arguing similarly leads to formulas for the other three cases, which we state below.

**Proposition 3.5.** *Consider  $i, j$  such that  $i < j$ .*

*For  $i, j$  both odd:*

$$\mathbb{P}(i \text{ and } j \notin A + A) = \frac{1}{2^{j+1}} F_{q+2}^r F_{q+4}^{r'} \quad (3.21)$$

where

$$\begin{aligned}
q &= 2 \left\lceil \frac{i+1}{j-i} \right\rceil \\
r &= \frac{1}{2} \left( (j-i) \left\lceil \frac{i+1}{j-i} \right\rceil - (i+1) \right) \\
r' &= \frac{1}{2} \left( j+1 - (j-i) \left\lceil \frac{i+1}{j-i} \right\rceil \right).
\end{aligned} \tag{3.22}$$

For  $i$  even,  $j$  odd:

$$\mathbb{P}(i \text{ and } j \notin A + A) = \frac{1}{2^{j+1}} F_o F_{q+2}^r F_{q+4}^{r'} \tag{3.23}$$

where

$$\begin{aligned}
o &= 2 \left\lceil \frac{i/2+1}{j-i} \right\rceil - 1 \\
q &= 2 \left\lceil \frac{i+1}{j-i} \right\rceil \\
r &= \frac{1}{2} \left( (j-i-1) \left\lceil \frac{i+1}{j-i} \right\rceil - (i+1) + o \right) \\
r' &= \frac{1}{2} \left( j - (j-i-1) \left\lceil \frac{i+1}{j-i} \right\rceil - o \right).
\end{aligned} \tag{3.24}$$

For  $i$  odd,  $j$  even:

$$\mathbb{P}(i \text{ and } j \notin A + A) = \frac{1}{2^{j+1}} F_{o'+2} F_{q+2}^r F_{q+4}^{r'} \tag{3.25}$$

where

$$\begin{aligned}
o' &= 2 \left\lceil \frac{j/2+1}{j-i} \right\rceil - 2 \\
q &= 2 \left\lceil \frac{i+1}{j-i} \right\rceil \\
r &= \frac{1}{2} \left( (j-i-1) \left\lceil \frac{i+1}{j-i} \right\rceil - (i+1) + o' \right) \\
r' &= \frac{1}{2} \left( j - (j-i-1) \left\lceil \frac{i+1}{j-i} \right\rceil - o' \right).
\end{aligned} \tag{3.26}$$

For  $i, j$  both even:

$$\mathbb{P}(i \text{ and } j \notin A + A) = \frac{1}{2^{j+1}} F_o F_{o'} F_{q+2}^r F_{q+4}^{r'} \tag{3.27}$$

where

$$\begin{aligned}
o &= 2 \left\lceil \frac{i/2 + 1}{j - i} \right\rceil - 1 \\
o' &= 2 \left\lceil \frac{j/2 + 1}{j - i} \right\rceil - 2 \\
q &= 2 \left\lceil \frac{i + 1}{j - i} \right\rceil \\
r &= \frac{1}{2} \left( (j - i - 2) \left\lceil \frac{i + 1}{j - i} \right\rceil - (i + 1) + o + o' \right) \\
r' &= \frac{1}{2} \left( j - 1 - (j - i - 2) \left\lceil \frac{i + 1}{j - i} \right\rceil - o - o' \right).
\end{aligned} \tag{3.28}$$

We conclude this section with some bounds on  $\mathbb{P}(i \text{ and } j \notin A + A)$ . We have (Binet's formula)

$$F_n = \frac{1}{\sqrt{5}}(\phi^n - (-1/\phi)^n), \tag{3.29}$$

where  $\phi = (1 + \sqrt{5})/2$  is the golden ratio. Therefore, for even  $n$  we have

$$F_n \leq \frac{1}{\sqrt{5}}\phi^n. \tag{3.30}$$

Since  $q + 2$  and  $q + 4$  are always even, then for any  $i, j$  both odd, we have

$$\begin{aligned}
\mathbb{P}(i \text{ and } j \notin A + A) &= \frac{1}{2^{j+1}} F_{q+2}^r F_{q+4}^{r'} \\
&\leq \frac{1}{2^{j+1}} \left( \frac{\phi^{q+2}}{\sqrt{5}} \right)^r \left( \frac{\phi^{q+4}}{\sqrt{5}} \right)^{r'} \\
&= \frac{1}{2^{j+1}} \frac{\phi^{(qr + (q+2)r') + (2r + 2r')}}{5^{(r+r')/2}} \\
&= \frac{1}{2^{j+1}} \frac{\phi^{j+1+j-i}}{5^{(j-i)/4}} \\
&= \frac{\phi^{2j+1}}{2^{j+1} 5^{j/4}} \frac{5^{i/4}}{\phi^i},
\end{aligned} \tag{3.31}$$

where the second to last equality comes from (3.18). In fact, we can use Proposition 3.5 to show that (3.31) holds for all  $i, j$  (slightly better constants hold for the other  $i, j$ ).

If  $i = k$  and  $j = k + m$ , where  $m$  is fixed and  $k$  goes to infinity, a lower bound similar to (3.31) also holds. First note that for even  $n$

$$\begin{aligned}
F_n^r &= \frac{1}{5^{r/2}} (\phi^n - \phi^{-n})^r \\
&= \frac{1}{5^{r/2}} \phi^{nr} (1 - \phi^{-2n})^r \\
&= \frac{1}{5^{r/2}} \phi^{nr} (1 - r(1 - c)^{r-1} \phi^{-2n})
\end{aligned} \tag{3.32}$$

for some  $c$  such that  $0 < c < 1/\phi^{2n}$  by Taylor expansion. Therefore for odd  $i, j$ , we have

$$\begin{aligned}
\mathbb{P}(i \text{ and } j \notin A + A) &= \frac{1}{2^{j+1}} F_{q+2}^r F_{q+4}^{r'} \\
&\geq \frac{1}{2^{j+1}} \frac{1}{5^{(q+2)/2}} \phi^{(q+2)r} (1 - r\phi^{-2(q+2)}) \frac{1}{5^{(q+4)/2}} \phi^{(q+4)r'} (1 - r'\phi^{-2(q+4)}) \\
&= \frac{\phi^{2j+1}}{2^{j+1} 5^{j/4}} \frac{5^{i/4}}{\phi^i} (1 - r\phi^{-2(q+2)}) (1 - r'\phi^{-2(q+4)}) \\
&\geq \frac{\phi^{2j+1}}{2^{j+1} 5^{j/4}} \frac{5^{i/4}}{\phi^i} (1 - (r + r')\phi^{-2(q+2)}) \\
&\geq \frac{\phi^{2j+1}}{2^{j+1} 5^{j/4}} \frac{5^{i/4}}{\phi^i} (1 - (j - i)\phi^{-i/(j-i)}), \tag{3.33}
\end{aligned}$$

and similar formulas hold for the other parity cases. If  $j/i \rightarrow 1$  not too slowly, then the remainder term on the right-hand-side of (3.33) goes to 1. For example, if  $i = k$  and  $j = k + m$ , then we have the following corollary by combining (3.31) and (3.33).

**Corollary 3.6.** *For any fixed  $m$ ,*

$$\mathbb{P}(k \text{ and } k + m \notin A + A) \sim \frac{\phi^{2(k+m)+1}}{2^{(k+m)+1} 5^{(k+m)/4}} \frac{5^{k/4}}{\phi^k} = \frac{\phi^{k+1}}{2^{k+1}} \frac{\phi^{2m}}{2^m 5^{m/4}}, \tag{3.34}$$

as  $k$  goes to infinity with  $k, k + m$  are both odd. Similar asymptotics hold for general  $k, k + m$ . If we ignore the constants related to  $m$ , we have

$$\mathbb{P}(k \text{ and } k + m \notin A + A) = \Theta((\phi/2)^k) \tag{3.35}$$

as  $k$  goes to infinity with any  $k, k + m$ .

Note that since  $\mathbb{P}(i \text{ and } j \notin A + A)$  has exponential decay in  $i, j$  as seen in (3.31), then (3.11) converges as  $n \rightarrow \infty$ ; that is

$$\lim_{n \rightarrow \infty} \text{Var} (M_{[0, n-1]}(A)) = 4 \sum_{i < j} \mathbb{P}(i \text{ and } j \notin A + A) - 40 \tag{3.36}$$

exists and is finite. In particular, we know that the limit is an infinite sum of Fibonacci products. However, we could not find a closed form for this sum. Nonetheless, because of the exponential decay in the terms in the sum, we can approximate the variance well. In particular, note that the tail of the sum has exponential decay:

$$\begin{aligned}
\sum_{n \leq i < j} P(i \text{ and } j \notin A + A) &\leq \sum_{n \leq i < j} \frac{\phi}{2} \left( \frac{\phi^2}{2 \cdot 5^{1/4}} \right)^j \left( \frac{5^{1/4}}{\phi} \right)^i \\
&\leq \frac{\phi}{2} \left( \sum_{n \leq j} \left( \frac{\phi^2}{2 \cdot 5^{1/4}} \right)^j \right) \left( \sum_{n \leq i} \left( \frac{5^{1/4}}{\phi} \right)^i \right) \\
&\leq \frac{\phi}{2} \left( \frac{1}{1 - \phi^2/2 \cdot 5^{1/4}} \right) \left( \frac{1}{1 - 5^{1/4}/\phi} \right) \left( \frac{\phi^2}{2 \cdot 5^{1/4}} \right)^n \left( \frac{5^{1/4}}{\phi} \right)^n \\
&\leq 87 \left( \frac{\phi}{2} \right)^n \leq 87(0.81)^n. \tag{3.37}
\end{aligned}$$



Here we use that (3.31) holds for all  $i, j$ . Using Mathematica to sum the first 300 terms of (3.36), whose exact form is given in Proposition 3.5, we get the following approximation for the variance:

$$\lim_{n \rightarrow \infty} \text{Var} (M_{[0, n-1]}(A)) = 35.9658 + E, \quad (3.38)$$

where  $|E| < 10^{-4}$ . The error term  $E$  comes mostly from truncating the computation of the 300-term series given by Mathematica. By (3.37), the error term from truncating the series at  $n = 300$  is less than  $87(0.81)^{300} \sim 3 \cdot 10^{-28}$ , which is much less than the Mathematica error term. This proves Theorem 1.5.

#### 4. EXPONENTIAL BOUNDS

We now prove Theorem 1.2 and find exponential bounds for the distribution of  $M_{[0, n-1]}(A)$ .

*Proof of Theorem 1.2.* For the lower bound, we construct many  $A$  such that  $A + A$  is missing  $k$  elements. First suppose that  $k$  is even. Let the first  $k/2$  non-negative integers not be in  $A$ . Then let the rest of the elements of  $A$  be any subset  $A'$  that fills in (so  $A' + A'$  has no missing elements between its largest and smallest elements); that is  $M_{n-k/2}(A') = 0$ . By [MO, Proposition 8], we can show that

$$\mathbb{P}(M_{[0, n-1]}(A') = 0) > 1/2^{10} \quad (4.1)$$

independent of  $n$ . If  $L \subseteq [0, \ell - 1]$  and  $U \subseteq [n - u, n - 1]$  are fixed, then their proposition says that

$$\mathbb{P}([2\ell - 1, 2n - 2u - 1] \subseteq A' + A' \mid A' \cap [0, \ell - 1] = L, A' \cap [n - u, n - 1] = U) > 1 - 6(2^{-|L|} + 2^{-|U|}), \quad (4.2)$$

independent of  $n$ . Therefore,

$$\begin{aligned} \mathbb{P}([2\ell - 1, 2n - 2u - 1] \subseteq A' + A' \text{ and } A' \cap [0, \ell - 1] = L, A' \cap [n - u, n - 1] = U) \\ > (1 - 6(2^{-|L|} + 2^{-|U|}))2^{-\ell}2^{-u}. \end{aligned} \quad (4.3)$$

By letting  $L = [0, \ell - 1], U = [n - u, n - 1]$  so the ends fill in, we get that

$$\mathbb{P}(A' + A' = [0, 2n - 2]) > (1 - 6(2^{-\ell} + 2^{-u}))2^{-\ell}2^{-u}. \quad (4.4)$$

Letting  $\ell = u = 4$  so that the first term in the product is positive, we get that

$$\mathbb{P}(A' + A' = [0, 2n - 2]) > (1 - 6(2^{-4} + 2^{-4}))2^{-4}2^{-4} = 1/2^{10}, \quad (4.5)$$

independent of  $n$ , which gives us (4.1).

As  $A = k/2 + A'$ , we have  $A + A = k + A' + A' = [k, 2n - 2]$  and so  $M_{[0, n-1]}(A) = k$  as seen by Figure 7.

Therefore we have

$$\begin{aligned} \mathbb{P}(M_{[0, n-1]}(A) = k) &\geq \mathbb{P}(A = k/2 + A' \text{ and } M_{n-k/2}(A') = 0) \\ &= \left(\frac{1}{2}\right)^{k/2} \mathbb{P}(M_{n-k/2}(A') = 0) \\ &\gg \left(\frac{1}{2}\right)^{k/2} \geq (0.70)^k, \end{aligned} \quad (4.6)$$

where the implied constants are independent of  $n$  by (4.1). This proves the lower bound in Theorem 1.2 when  $k$  is even.

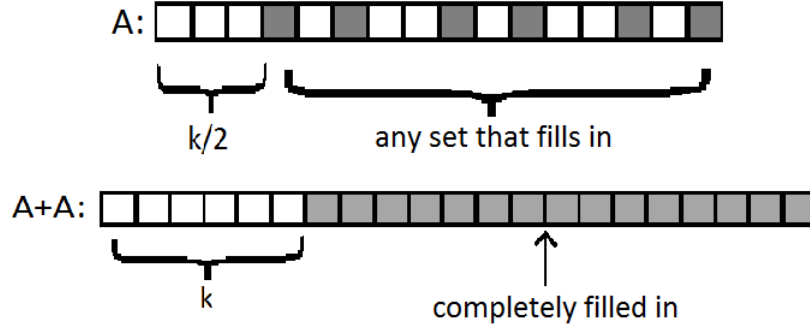


FIGURE 7.  $A$  and  $A + A$  for lower bound.

If  $k$  is odd, then we can let  $L = [0, \ell - 1] \setminus \{2, 3\}$  and  $U = [n - u, n - 1]$  so that only the element 3 is missing from  $A' + A'$ . Then we get a bound for  $\mathbb{P}(M_{[0, n-1]}(A') = 1)$  as in (4.1). Letting  $A = (k - 1)/2 + A'$ , we get the desired lower bound in Theorem 1.2 for when  $k$  is odd.

For the upper bound, we can use bounds like

$$\mathbb{P}(k \notin A + A) \leq \left(\frac{3}{4}\right)^{k/2} \quad (4.7)$$

from [MO]. Again, first suppose that  $k$  is even. Note that if  $A + A$  is missing  $k$  elements, then one of these missing elements must be at least  $k/2$  elements away from the ends of  $[0, 2n - 2]$ . That is, we have the following situation (see Figure 8).

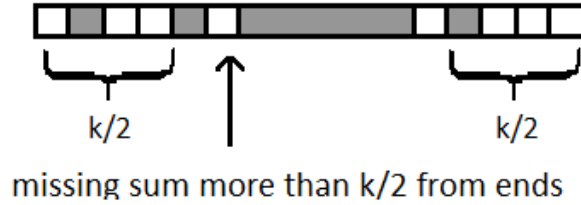


FIGURE 8. Upper bound for  $\mathbb{P}(M_{[0, n-1]}(A) = k)$

Therefore

$$\begin{aligned} \mathbb{P}(M_{[0, n-1]}(A) = k) &\leq \mathbb{P}(A + A \text{ missing element at least } k/2 \text{ away from edges}) \\ &= \mathbb{P}(j \notin A + A, j \in [k/2, 2n - k/2]) \\ &\leq 2 \sum_{j \geq k/2} \left(\frac{3}{4}\right)^{j/2} \\ &\ll \left(\frac{3}{4}\right)^{k/4} \approx (0.93)^k. \end{aligned} \quad (4.8)$$

Note that this bound does not use the fact that there may be missing elements on both ends at the same time. By focusing on one particular side, we can get a stronger result. For example, we

have the following inequality for the probability of missing  $k$  elements in  $[0, n/2]$ :

$$\begin{aligned}
\mathbb{P}(|[0, n/2] \setminus (A + A)| = k) &\leq \mathbb{P}(j \notin A + A, j \in [k, n/2]) \\
&\leq 2 \sum_{j \geq k} \left(\frac{3}{4}\right)^{j/2} \\
&\ll \left(\frac{3}{4}\right)^{k/2} \approx (0.87)^k
\end{aligned} \tag{4.9}$$

and similarly for  $\mathbb{P}(|[3n/2, 2n] \setminus (A + A)| = k)$ . Furthermore, (7.27) from Section 7 connects the probability of missing  $k$  elements to the probability of missing elements in  $[0, n/2]$  and  $[3n/2, 2n]$ :

$$\mathbb{P}(M_{[0, n-1]}(A) = k) = \sum_{i+j=k} \mathbb{P}(|[0, n/2] \setminus (A+A)| = i) \mathbb{P}(|[3n/2, 2n] \setminus (A+A)| = j) + O\left(\left(\frac{3}{4}\right)^{n/4}\right). \tag{4.10}$$

Combining (4.9) and (4.10), we get

$$\begin{aligned}
&\mathbb{P}(M_{[0, n-1]}(A) = k) \\
&= \sum_{i+j=k} \mathbb{P}(|[0, n/2] \setminus (A + A)| = i) \mathbb{P}(|[3n/2, 2n] \setminus (A + A)| = j) + O\left(\left(\frac{3}{4}\right)^{n/4}\right) \\
&\ll \sum_{i+j=k} \left(\frac{3}{4}\right)^{i/2} \left(\frac{3}{4}\right)^{j/2} + \left(\frac{3}{4}\right)^{n/4} \\
&\ll k \left(\frac{3}{4}\right)^{k/2} + \left(\frac{3}{4}\right)^{n/4}.
\end{aligned} \tag{4.11}$$

Therefore if  $k/2 < n/4$ , then we get the desired bound

$$\mathbb{P}(M_{[0, n-1]}(A) = k) \ll k \left(\frac{3}{4}\right)^{k/2} \approx (0.87)^k. \tag{4.12}$$

Note that the bound in (4.12) for the distribution is exactly the same as the bound in (4.9) for missing elements on a single side. Since all our bounds are exponential and (4.10) multiplies  $\mathbb{P}(|[0, n/2] \setminus (A + A)| = i)$  with  $\mathbb{P}(|[3n/2, 2n] \setminus (A + A)| = j)$ , we can always use this approach to transform bounds on the probability of missing elements in  $(A + A) \cap [0, n/2]$  to equally good bounds on number of missing elements in all of  $A + A$ . So it is sufficient to just develop bounds on missing elements on one side of  $A + A$ . In particular, we can use this approach to transform the

bounds in Corollary 3.6 to improve the bounds in (4.12). By Corollary 3.6, we have

$$\begin{aligned}
\mathbb{P}(|[0, n/2] \setminus (A + A)| = k) &\leq \mathbb{P}(A + A \text{ misses } 2 \text{ elements greater than } k - 3) \\
&= \mathbb{P}(i, j \notin A + A, i, j \in [k - 3, n/2]) \\
&= \sum_{k-3 < i < j} \mathbb{P}(i \text{ and } j \notin A + A) \\
&\ll \sum_{k-3 < i < j} \frac{\phi^{2j+1}}{2^{j+1} 5^{j/4}} \frac{5^{i/4}}{\phi^i} \\
&\ll \frac{\phi^{2k+1}}{2^{k+1} 5^{k/4}} \frac{5^{k/4}}{\phi^k} = \left(\frac{\phi}{2}\right)^k \approx (0.81)^k.
\end{aligned} \tag{4.13}$$

Then using the previous approach, we get a similar bound on the total number of missing sums:

$$\mathbb{P}(M_{[0, n-1]}(A) = k) \ll \left(\frac{\phi}{2}\right)^k \approx (0.81)^k. \tag{4.14}$$

Note that as in (4.10), we always have an extra  $(3/4)^{n/4}$  term. To make this term negligible, we need to have  $(3/4)^{n/4} < (0.81)^k$ , which means  $n > k \cdot 4 \log(0.81)/\log(3/4) \sim 2.92k$  or that  $k < 0.34n$ . This condition is sufficient in this case where we have the bound  $(\phi/2)^k$ . However in general, we know that we have a lower bound of  $(1/2)^{k/2}$  for the distribution. Therefore, to make the  $(3/4)^{n/4}$  term always negligible, we can have  $(3/4)^{n/4} < (1/2)^{k/2}$ , which means  $n > k \cdot 2 \log(1/2)/\log(3/4) \sim 5k$ , as in the statement of Theorem 1.2. Note that then the implied constants are independent of  $n$ . Combining (4.6) and (4.14), we get Theorem 1.2.  $\square$

## 5. APPROXIMATING $\mathbb{P}(k + a_1, k + a_2, \dots, \text{ AND } k + a_m \notin A + A)$

In this section, we prove Theorem 1.6 which says that for any fixed  $a_1, \dots, a_m$ , there exists  $\lambda_{a_1, \dots, a_m}$  such that

$$\mathbb{P}(k + a_1, k + a_2, \dots, \text{ and } k + a_m \notin A + A) = \Theta(\lambda_{a_1, \dots, a_m}^k), \tag{5.1}$$

where the implied constants depend on  $a_1, \dots, a_m$  but not  $k$ . Therefore, the probability is approximately exponential.

To prove this theorem, we use a version of Fekete's Lemma, which says that sub-additive sequences are approximately linear. From [S] we have the following version in which the sequence is both sub-additive and super-additive.

**Lemma 5.1.** *If  $b_n$  is a sequence such that*

$$b_n + b_m - 1 \leq b_{n+m} \leq b_n + b_m + 1 \tag{5.2}$$

*for all  $n, m$ , then  $\lambda = \inf b_n/n$  exists and for all  $n$ ,*

$$\left| \frac{b_n}{n} - \lambda \right| < \frac{1}{n}. \tag{5.3}$$

**Remark 5.2.** *The proof of this Lemma can be easily modified to get that if*

$$b_n + b_m - c \leq b_{n+m} \leq b_n + b_m + c \tag{5.4}$$

for some constant  $c > 0$ , then

$$\left| \frac{b_n}{n} - \lambda \right| < \frac{c}{n}. \quad (5.5)$$

Suppose that  $a_n$  is approximately multiplicative rather than approximately additive so that for some constant  $c > 1$

$$c^{-1} \cdot a_m a_n \leq a_{m+n} \leq c \cdot a_m a_n \quad (5.6)$$

for all  $m, n$ . As  $b_n = \log a_n$  satisfies the properties of Lemma 5.1, for  $\lambda = \inf \frac{\log a_n}{n}$  we have

$$\left| \frac{\log a_n}{n} - \lambda \right| < \frac{\log c}{n} \quad (5.7)$$

for all  $n$ . That is,

$$c^{-1} \lambda^n \leq a_n \leq c \lambda^n \quad (5.8)$$

for all  $n$ , implying

$$a_n = \Theta(\lambda^n). \quad (5.9)$$

Therefore we just need to relate  $\mathbb{P}(k + a_1, k + a_2, \dots, \text{ and } k + a_m \notin A + A)$  as a function of  $k$  to some approximately multiplicative function satisfying (5.6).

For example, consider  $\mathbb{P}(18, 19, \text{ and } 21 \notin A + A)$ , whose condition graph is in Figure 9. Note that this graph has a loop from vertex 9 to itself since  $9 + 9 = 18$ . We can symmetrize this graph by removing this loop and also removing the edge between vertices 8 and 10 and the edge between vertices 9 and 10, resulting in the modified condition graph in Figure 10.

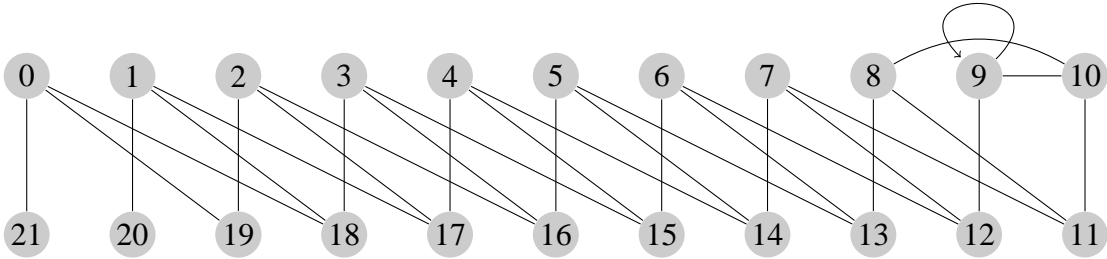


FIGURE 9. Condition graph for  $\mathbb{P}(18, 19, 21 \notin A + A)$ .

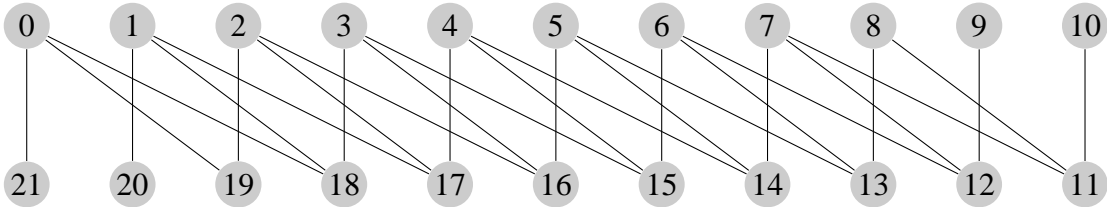


FIGURE 10. Modified condition graph for  $\mathbb{P}(18, 19, 21 \notin A + A)$ .

Denote the probability of getting a vertex cover for graphs like the one in Figure 10 of length  $n$  by  $f(n)$ ; so the probability of getting a vertex cover in Figure 10 is  $f(11)$ .

Note that  $f(11)$  is an upper bound for the probability in the original condition graph in Figure 9 since we have removed some edges. On the other hand, we have the following lower bound:

$$\begin{aligned}
& \mathbb{P}(18, 19, \text{ and } 21 \notin A + A) \\
& \geq \mathbb{P}(18, 19, 21 \notin A + A \text{ and } 9, 10, 11, 12 \notin A) \\
& = \mathbb{P}(18, 19, 21 \notin A + A \mid 9, 10, 11, 12 \notin A) \mathbb{P}(9, 10, 11, 12 \notin A). \tag{5.10}
\end{aligned}$$

Note that the condition graph for  $\mathbb{P}(18, 19, 21 \notin A + A \mid 9, 10, 11, 12 \notin A)$  is the original condition graph in Figure 9 with all edges incident on vertices 9, 10, 11 or 12 removed, as depicted in Figure 11.

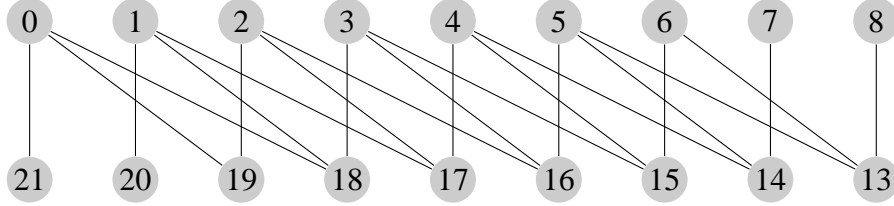


FIGURE 11. Condition graph for  $\mathbb{P}(18, 19, \text{ and } 21 \notin A + A \mid 9, 10, 11, 12 \notin A)$ .

Note that in Figure 11 we have removed vertices 9, 10, 11 and 12 completely since there are no longer any conditions on them in  $\mathbb{P}(18, 19, \text{ and } 21 \notin A + A \mid 9, 10, 11, 12 \notin A)$ . Finally, note that the probability of getting a vertex cover in the graph in Figure 11 is just  $f(9)$ . Therefore, by (5.10), we have

$$(1/2)^4 f(9) \leq \mathbb{P}(18, 19, \text{ and } 21 \notin A + A) \leq f(11), \tag{5.11}$$

where we use that  $\mathbb{P}(9, 10, 11, 12 \notin A) = (1/2)^4$ .

Since the condition graph for  $\mathbb{P}(k, k+1, \text{ and } k+3 \notin A + A)$  is just a longer version of the condition graph for  $\mathbb{P}(18, 19, \text{ and } 21 \notin A + A)$ , we can apply the same method as before to get that

$$(1/2)^4 f(k/2) \leq \mathbb{P}(k, k+1, \text{ and } k+3 \notin A + A) \leq f((k+4)/2) \tag{5.12}$$

for even  $k$ , with a similar formula holding for odd  $k$ . Therefore we are reduced to studying  $f(n)$ , which is easier to investigate since the condition graph is more symmetric. We will show that  $f(n)$  satisfies (5.6), implying it is approximately exponential.

For example, to see that  $f(11) \leq f(4)f(7)$ , we can separate the graph in Figure 10 at the 4<sup>th</sup> vertex and remove edges that cross this gap, resulting in the graph in Figure 12.

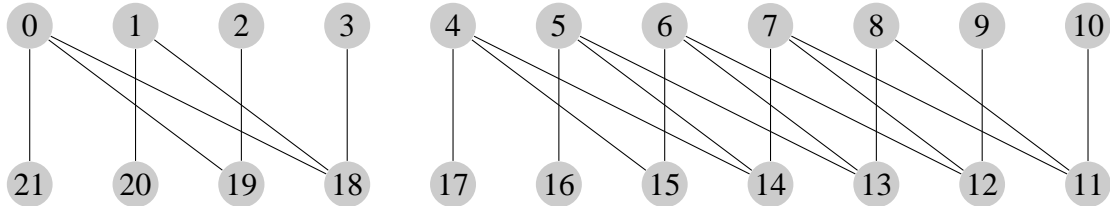


FIGURE 12. Upper Bound for  $f(11)$ .

Since the components are independent smaller copies of the original, the probability of getting a vertex cover for the graph in Figure 12 is  $f(4)f(7)$ . We can do this for any integer less than

11, defining  $f(n)$  for small integers by truncating at the  $n^{\text{th}}$  vertex. Since we have removed some edges to get the graph in Figure 12, we have

$$f(11) \leq f(4)f(7) \quad (5.13)$$

as desired.

To get a lower bound for  $f(11)$ , we use that

$$f(11) \geq f(11 \mid 4, 5, 6, 15, 16, 17 \text{ chosen})\mathbb{P}(4, 5, 6, 15, 16, 17 \text{ chosen}), \quad (5.14)$$

where  $f(11 \mid 4, 5, 6, 15, 16, 17 \text{ chosen})$  denotes the probability of getting a vertex cover for the graph in Figure 12 given that the vertices 4, 5, 6, 15, 16, 17 are chosen. The graph for  $f(11 \mid 4, 5, 6, 15, 16, 17 \text{ chosen})$  is depicted in Figure 13. The probability of getting vertex covers for the

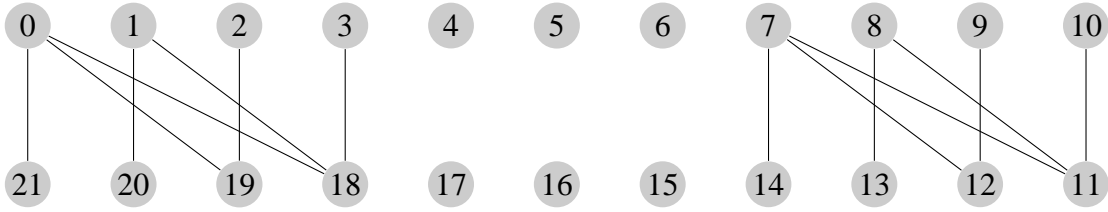


FIGURE 13. Lower Bound for  $f(11)$ .

two independent components is  $f(4)f(4)$ . Therefore from (5.14), we get that

$$f(11) \geq (1/2)^6 f(4)f(4) \geq (1/2)^6 f(4)f(7), \quad (5.15)$$

with the last inequality since  $f(n)$  is decreasing. Therefore, in general we have

$$(1/2)^6 f(m)f(n) \leq f(m+n) \leq f(m)f(n), \quad (5.16)$$

and so  $f(n)$  satisfies the conditions of (5.6). By the modified version of Fekete's Lemma, we have

$$f(n) = \Theta(\lambda^n) \quad (5.17)$$

for some  $\lambda$ . Therefore by (5.12), we have

$$\mathbb{P}(k, k+1, \text{ and } k+3 \notin A+A) = \Theta(\lambda^{k/2}), \quad (5.18)$$

which proves Theorem 1.6 for the case  $a_1 = 0, a_2 = 1, a_3 = 3$ .

The general situation follows in exactly the same way: by first making the configuration graph of  $\mathbb{P}(k+a_1, \dots, \text{ and } k+a_m \notin A+A)$  look more symmetric and then using the modified Fekete's Lemma.

## 6. CONSECUTIVE MISSING SUMS

In this section, we prove Theorem 1.7 and its generalization Theorem 1.8. We begin by proving Theorem 1.7, which says that

$$\left(\frac{1}{2}\right)^{(k+m)/2} \ll \mathbb{P}(k+1, \dots, \text{ and } k+m \notin A+A) \ll \left(\frac{1}{2}\right)^{(k+m)/2} (1+\epsilon_m)^k. \quad (6.1)$$

The lower bound comes from the construction in Figure 7 by letting the first  $\lfloor (k+m)/2 \rfloor$  elements of  $A$  be missing, which forces the first  $k+m$  elements of  $A+A$  to be missing as well. That is,

$$\begin{aligned} & \mathbb{P}(0, 1, \dots, k+m-1, \text{ and } k+m \notin A+A) \\ &= \mathbb{P}(0, 1, \dots, \text{ and } \lfloor (k+m)/2 \rfloor \notin A) \\ &= (1/2)^{\lfloor (k+m)/2 \rfloor + 1}. \end{aligned} \tag{6.2}$$

Therefore, we only need to prove the upper bound.

Before giving the proof, we consider an example with condition graphs which illustrates the idea. Consider  $\mathbb{P}(16, 17, 18, 19, 20 \notin A+A)$ . The condition graph here is given in Figure 14.

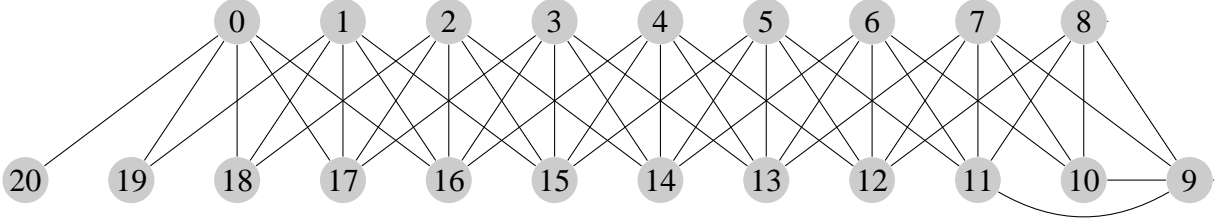


FIGURE 14. Condition graph for  $\mathbb{P}(16, 17, 18, 19, 20 \notin A+A)$ .

We need to find the probability of getting a vertex cover for this graph. If we remove some edges, the probability of getting a vertex cover for the resulting graph is an upper bound for the probability of getting a vertex cover for the original graph. We can remove some edges to get the graph of Figure 15.

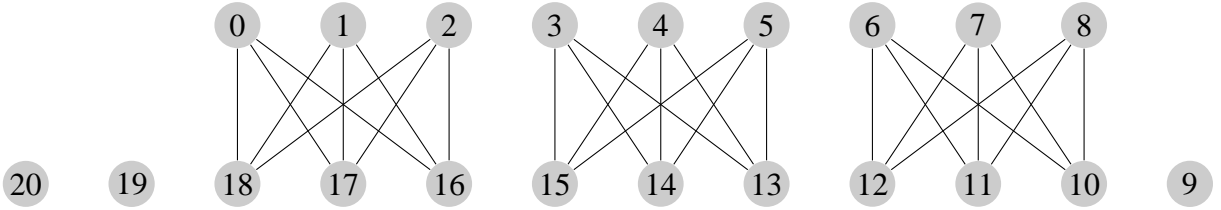


FIGURE 15. Graph after removing some edges.

The resulting graph has  $3 \sim 20/6$  components that are all complete bipartite graphs with 6 vertices. These are easier to handle since the only way to get a vertex cover for such graphs is to have all vertices on one side be chosen. So the probability of getting a vertex cover for one of these complete bipartite components is less than  $(1/2)^3 + (1/2)^3 = 2/2^3$ . Since the components are also independent, we have

$$\mathbb{P}(16, 17, 18, 19, 20 \notin A+A) \leq \left(\frac{2}{2^3}\right)^3 \sim \left(\frac{1}{4}\right)^{20/6}. \tag{6.3}$$

and in general we get that

$$\mathbb{P}(k, k+1, k+2, k+3, k+4 \notin A+A) \leq \left(\frac{2}{2^3}\right)^{(k+4)/6} = \left(\frac{2^{1/3}}{2}\right)^{(k+4)/2}. \tag{6.4}$$



We use this approach in the general proof. Notice that as  $m \rightarrow \infty$ , the size of the complete bipartite graphs grows, and so we will be taking out relatively fewer and fewer constraints. Therefore, this approach gets us closer to the correct answer.

Now we give a formal proof of Theorem 1.7 that does not rely on the condition graphs.

*Proof.* We first do the proof for  $\mathbb{P}(k, k+1, \dots, \text{ and } k+2m-1 \notin A+A)$  with  $2m-1$  instead of  $m$ . Note that since the probability depends only on  $[0, k+2m-1] \cap A$ , we can assume that  $A \subseteq [0, k+2m-1]$ . We will also assume that  $m$  divides  $k$  and that

$$k = qm \quad (6.5)$$

with  $q$  even.

We begin by writing  $A$  as the following disjoint union:

$$A = A_0 \cup A_1 \cup \dots \cup A_q \cup A_{q+1}, \quad (6.6)$$

where

$$A_j = A \cap [jm, (j+1)m - 1]. \quad (6.7)$$

Then if  $[k, k+2m-1] \cap (A+A) = \emptyset$ , then  $[k, k+2m-1] \cap (A_j + A_{q-j}) = \emptyset$  for all  $j$ . Note that

$$A_j + A_{q-j} \subseteq [k, k+2m-2]. \quad (6.8)$$

Therefore,  $[k, k+2m-1] \cap (A_j + A_{q-j}) = \emptyset$  implies  $A_j + A_{q-j} = \emptyset$ . If  $q$  is even, we have

$$\begin{aligned} \mathbb{P}(k, k+1, \dots, \text{ and } k+2m-1 \notin A+A) &< \mathbb{P}([k, k+2m] \cap (A_j + A_{q-j}) = \emptyset \text{ for all } j \leq q/2) \\ &= \mathbb{P}(A_j + A_{q-j} = \emptyset \text{ for all } j \leq q/2) \\ &= \mathbb{P}(A_j = \emptyset \text{ or } A_{q-j} = \emptyset \text{ for all } j \leq q/2). \end{aligned} \quad (6.9)$$

For different  $j$ , the pairs of sets  $A_j, A_{q-j}$  are disjoint. Therefore, we have independence:

$$\mathbb{P}(A_j = \emptyset \text{ or } A_{q-j} = \emptyset \text{ for all } j \leq q/2) = \mathbb{P}(A_{q/2} = \emptyset) \prod_{j=0}^{q/2-1} \mathbb{P}(A_j = \emptyset \text{ or } A_{q-j} = \emptyset). \quad (6.10)$$

Finally, note that

$$\mathbb{P}(A_j = \emptyset \text{ or } A_{q-j} = \emptyset) \leq \mathbb{P}(A_j = \emptyset) + \mathbb{P}(A_{q-j} = \emptyset) = \frac{2}{2^m}. \quad (6.11)$$

Combining (6.9), (6.10), and (6.11), we find

$$\begin{aligned} \mathbb{P}(k, k+1, \dots, \text{ and } k+2m-1 \notin A+A) &\leq \frac{1}{2^m} \prod_{j=0}^{q/2-1} \frac{2}{2^m} \\ &= 2^{q/2} \left( \frac{1}{2^m} \right)^{q/2+1} \\ &= 2^{k/2m} \left( \frac{1}{2} \right)^{(k+2m)/2}. \end{aligned} \quad (6.12)$$

This inequality is true for all  $m, k$  such that  $q = k/m$  is an even integer.

Changing  $m$  to  $m/2$ , we get that

$$\mathbb{P}(k, k+1, \dots, \text{ and } k+m-1 \notin A+A) \leq 2^{k/m} \left( \frac{1}{2} \right)^{(k+m)/2} \quad (6.13)$$

for even  $m$  and  $q = k/m$  still an even integer. Note that (6.13) is similar to the bound we get in (6.4) using the condition graph approach.

For odd  $m$ , we just need to use (6.13), noting that

$$\mathbb{P}(k, k+1, \dots, k+m-1, \text{ and } k+m \notin A+A) \leq \mathbb{P}(k, k+1, \dots, \text{ and } k+m-1 \notin A+A). \quad (6.14)$$

For odd  $q$ , we need to partition  $A$  such that there is a block in the very middle of  $A$ . This ensures that this middle block is matched with itself (just like  $A_{q/2}$  was matched with itself when  $q$  was even). This gives us the extra  $1/2^m$  that is needed in order to achieve the bound. For non-integer  $q$ , we need to repartition  $A$  in a similar way. Therefore the bound in (6.13) holds in general, up to a constant.

Finally, note that as  $m \rightarrow \infty$ , we have  $2^{1/m} \rightarrow 1$ . Writing  $2^{1/m} = 1 + \epsilon_m$ , we have

$$\mathbb{P}(k, \dots, \text{ and } k+m-1 \notin A+A) < \left(\frac{1}{2}\right)^{(k+m)/2} (1 + \epsilon_m)^k, \quad (6.15)$$

where  $\epsilon_m \rightarrow 0$  as  $m \rightarrow \infty$ . By raising  $2^{1/m} = 1 + \epsilon_m$  to the  $m^{\text{th}}$  power, we see that

$$\epsilon_m < \frac{1}{m}. \quad (6.16)$$

Therefore a weakened version of the inequality says that

$$\left(\frac{1}{2}\right)^{(k+m)/2} \ll \mathbb{P}(k+1, \dots, \text{ and } k+m \notin A+A) \ll \left(\frac{1}{2}\right)^{(k+m)/2} (1 + \epsilon_m)^k, \quad (6.17)$$

where the implied constants are independent of  $m$  and  $k$ .

This bound is interesting since it means that the trivial lower bound is almost the right answer for the exact bound. The trivial lower bound makes us miss all of  $[0, k+m]$  in  $A+A$  as seen in (6.2) but we only need  $[k+1, k+m]$  to be missing. In this sense, we see that essentially the only way to miss  $m$  consecutive elements at  $k+1$  for large  $m$  is to miss all the previous elements as well.

Also, note that (6.17) implies that

$$\lambda_{0,1,\dots,m} \rightarrow \left(\frac{1}{2}\right)^{1/2} \quad (6.18)$$

as  $m \rightarrow \infty$  by definition of  $\lambda_{0,1,\dots,m}$ . □

Now we will prove Theorem 1.8, which says that

$$\lambda_{a_1,\dots,a_m} \leq \mathbb{P}(A, B \subseteq [0, \lfloor a_m/2 \rfloor] \mid a_1, \dots, a_m \notin A+B)^{1/(a_m+2)}. \quad (6.19)$$

Note that Theorem 1.7 is indeed a special case of this theorem since we have the following upper bound

$$\begin{aligned} \lambda_{0,1,\dots,m} &\leq \mathbb{P}(A, B \subseteq [0, \lfloor m/2 \rfloor] \mid 0, \dots, m \notin A+B)^{1/(m+2)} \\ &= \mathbb{P}(A, B \subseteq [0, \lfloor m/2 \rfloor] \mid A = \emptyset \text{ or } B = \emptyset)^{1/(m+2)} \\ &\leq \left(2 \left(\frac{1}{2}\right)^{\lfloor m/2 \rfloor + 1}\right)^{1/(m+2)}, \end{aligned} \quad (6.20)$$

which converges to  $\sqrt{1/2}$ .

The proof of Theorem 1.8 will be almost exactly the same as the proof of Theorem 1.7.

*Proof.* We will first show that for  $0 \leq a_1 < \dots < a_m$ ,

$$\begin{aligned} & \mathbb{P}(A \subseteq [0, k + a_m] \mid k + a_1, \dots, k + a_m \notin A + A) \\ & \leq \mathbb{P}(A, B \subseteq [0, a_m/2] \mid a_1, \dots, a_m \notin A + B)^{k/(a_m+2)} \end{aligned} \quad (6.21)$$

for all  $k, a_m$  such that  $a_m$  is even and  $a_m + 2$  divides  $k$ . Similar results hold in the other cases of  $k, m$ . Furthermore, we first assume that  $a_m = 2r - 2$ . Note that since the probability depends only on  $[0, k + 2r - 2] \cap A$ , we can take  $A \subseteq [0, k + 2r - 2]$ . Again, we first assume that  $r$  divides  $k$  and that  $k = qr$ . Then as before,

$$\begin{aligned} & \mathbb{P}(k + a_1, \dots, \text{ and } k + a_m \notin A + A) \\ & \leq \mathbb{P}(k + a_1, \dots, \text{ and } k + a_m \notin A_j + A_{q-j} \text{ for all } j \leq \lfloor q/2 \rfloor) \\ & = \prod_{j=0}^{\lfloor q/2 \rfloor} \mathbb{P}(k + a_1, \dots, \text{ and } k + a_m \notin A_j + A_{q-j}). \end{aligned} \quad (6.22)$$

The key fact is that if  $j \neq q - j$ , the sets  $A_j, A_{q-j}$  are independent and

$$\mathbb{P}(k + a_1, \dots, \text{ and } k + a_m \notin A_j + A_{q-j}) = \mathbb{P}(A, B \subseteq [0, r - 1] \mid a_1, \dots, a_m \notin A + B) \quad (6.23)$$

for all  $j$ . Therefore, if  $q$  is odd

$$\begin{aligned} & \mathbb{P}(k + a_1, \dots, \text{ and } k + a_m \notin A + A) \\ & \leq \mathbb{P}(A, B \subseteq [0, r - 1] \mid a_1, \dots, a_m \notin A + B)^{\lfloor q/2 \rfloor + 1} \\ & = \mathbb{P}(A, B \subseteq [0, r - 1] \mid a_1, \dots, a_m \notin A + B)^{k/2r + 1/2} \end{aligned} \quad (6.24)$$

and if  $q$  is even,

$$\begin{aligned} & \mathbb{P}(k + a_1, \dots, \text{ and } k + a_m \notin A + A) \\ & \leq \mathbb{P}(A \subseteq [0, r - 1] \mid a_1, \dots, a_m \notin A + A) \\ & \quad \times \mathbb{P}(A, B \subseteq [0, r - 1] \mid a_1, \dots, a_m \notin A + B)^{k/2r}. \end{aligned} \quad (6.25)$$

If we drop the terms that do not depend on  $k$ , we have for all even  $a_m$  and all  $k$  divisible by  $a_m + 2$

$$\begin{aligned} & \mathbb{P}(A \subseteq [0, k + a_m] \mid k + a_1, \dots, k + a_m \notin A + A) \\ & \leq \mathbb{P}(A, B \subseteq [0, a_m/2] \mid a_1, \dots, a_m \notin A + B)^{k/(a_m+2)}, \end{aligned} \quad (6.26)$$

which is (6.21). Note that if  $k$  is not divisible by  $a_m + 2$  or if  $a_m$  is not even, we have

$$\begin{aligned} & \mathbb{P}(A \subseteq [0, k + a_m] \mid k + a_1, \dots, k + a_m \notin A + A) \\ & \leq \mathbb{P}(A, B \subseteq [0, \lfloor a_m/2 \rfloor] \mid a_1, \dots, a_m \notin A + B)^{\lfloor k/(a_m+2) \rfloor}, \end{aligned} \quad (6.27)$$

which proves that (6.19). □

## 7. BOUNDS ON $m(k)$ , $w(k)$ , $y(k)$ , AND $z(k)$ FOR $k < 32$

As mentioned in §1.3 and covered in more detail in §7.2, it suffices to bound  $z(k)$ . Our strategy is this: if  $D + D$  (where  $D$  is a uniformly chosen subset of  $\mathbb{N}$  that contains 0) is missing exactly 7 elements, then it is very likely that those 7 missing sums are all smaller than 88 and typically even all smaller than 44. If we loop over all  $2^{43}$  possibilities  $\beta$  for  $D \cap [0, 44)$ , for each possibility we can compute  $(D + D) \cap [0, 44) = (\beta + \beta) \cap [0, 44)$  and a subset of  $(D + D) \cap [44, 48) \supseteq (\beta + \beta) \cap [44, 88)$ . From this (with a little theory to handle the tail of the sumset) we can bound the likelihood of

missing exactly 7 sums, given  $D \cap [0, 44)$ . By combining these estimates, we acquire bounds on  $z(7)$ .

Let  $n \geq 2$  be a natural number parameter (the computations reported here use  $n = 44$ , although  $n = 43$  is already enough to show  $m(7) < m(6) < m(8)$ ), and set

$$z(k \mid \beta) := \mathbb{P}(|\mathbb{N} \setminus (D + D)| = k \mid D \cap [0, n) = \beta). \quad (7.1)$$

We have

$$z(k) = \sum_{0 \in \beta \subseteq [0, n)} z(k \mid \beta) \mathbb{P}(D \cap [0, n) = \beta) = 2^{-(n-1)} \sum_{0 \in \beta \subseteq [0, n)} z(k \mid \beta), \quad (7.2)$$

so that it suffices to bound  $z(k \mid \beta)$  above and below for all  $0 \leq k < 32$  (our arbitrary notion of “small  $k$ ” is  $0 \leq k < 32$ ) and all  $0 \in \beta \subseteq [0, n)$ .

Further, set

$$\begin{aligned} \mathcal{B} &:= D \cap [0, n) \\ \mathcal{D} &:= [0, n) \setminus (\beta + \beta) \\ \mathcal{L} &:= [n, 2n) \setminus (\beta + \beta) \\ m &:= \min \mathcal{L} \\ \mathcal{T} &:= [2n, \infty) \\ \eta &:= \mathbb{E}[|[n, \infty) \setminus (D + D)| \mid \mathcal{B} = \beta] \\ \mu &:= 2^{-|\beta \cap [0, m-n]|}. \end{aligned} \quad (7.3)$$

If we condition on  $\mathcal{B} = \beta$ , then the elements of  $\mathcal{D}$  are *Definitely* missing from  $D + D$ , the elements of  $\mathcal{L}$  are *Likely* but not certain to be missing, and the elements of  $\mathcal{T}$ , the *Tail* of the natural numbers, are very likely to be missing. Note that  $2n - 1 \in \mathcal{L}$ , so  $\mathcal{L}$  is nonempty and  $m$  is well-defined.

**Lemma 7.1.** *For all  $k < |\mathcal{D}|$ , we have  $z(k \mid \beta) = 0$ .*

*Proof.* Conditioning on  $\mathcal{B} = \beta$ , we have  $\mathcal{D} \subseteq \mathbb{N} \setminus (D + D)$ . In fact,  $\mathcal{D} = [0, n) \setminus (D + D)$ .  $\square$

**Lemma 7.2.** *We have  $\eta = 5 \cdot 2^{-|\mathcal{B}|} + \sum_{\ell \in \mathcal{L}} 2^{-|\mathcal{B} \cap [0, \ell-n]|}$ .*

*Proof.* By linearity of expectation

$$\eta := \mathbb{E}[|[n, \infty) \setminus (D + D)|] = \mathbb{E}[|[n, 2n) \setminus (D + D)|] + \mathbb{E}[|\mathcal{T} \setminus (D + D)|]. \quad (7.4)$$

Again using linearity of expectation, we have

$$\mathbb{E}[|[n, 2n) \setminus (D + D)|] = \sum_{\ell \in \mathcal{L}} \mathbb{P}(\ell \notin D + D) \quad (7.5)$$

Since  $\ell \notin D + D$  is the same as (for  $n \leq \ell < 2n$ )

$$\begin{aligned} \ell \notin D + D &= \bigwedge_{i=0}^{\ell/2} (i \notin D \vee \ell - i \notin D) \\ &= \bigwedge_{\substack{b \in \beta \\ b \leq \ell - n}} \ell - b \notin D. \end{aligned} \quad (7.6)$$

Thus

$$\mathbb{P}(\ell \notin D + D) = 2^{-|\beta \cap [0, \ell - n]|}, \quad (7.7)$$

and so

$$\sum_{\ell \in \mathcal{L}} \mathbb{P}(\ell \notin D + D) = \sum_{\ell \in \mathcal{L}} 2^{-|\beta \cap [0, \ell - n]|}. \quad (7.8)$$

That

$$\mathbb{E}[|\mathcal{T} \setminus (D + D)|] = 5 \cdot 2^{-|\beta|} \quad (7.9)$$

is essentially in [MO], but we derive it here for the reader's convenience. Using linearity of expectation,

$$\mathbb{E}[|\mathcal{T} \setminus (D + D)|] = \sum_{t=2n}^{\infty} \mathbb{P}(t \notin D + D) \quad (7.10)$$

and

$$\mathbb{P}(t \notin D + D) = \mathbb{P}\left(\left(\bigwedge_{b \in \beta} t - b \notin D\right) \wedge \left(\bigwedge_{i=n}^{t/2} i \notin D \vee t - i \notin D\right)\right). \quad (7.11)$$

Now this has two cases leading to

$$\mathbb{P}(t \notin D + D) = \begin{cases} 2^{-|\beta|} (3/4)^{(t-2n+1)/2} & t \text{ is odd,} \\ 2^{-|\beta|} (1/2) (3/4)^{(t-2n)/2} & t \text{ is even.} \end{cases} \quad (7.12)$$

The infinite sum (7.10) now simplifies  $5 \cdot 2^{-|\beta|}$ . □

**Lemma 7.3.** *We have  $\max\{0, 1 - \eta\} \leq z(|\mathcal{D}| \mid \beta) \leq 1 - \mu$ .*

*Proof.* Trivially  $z(|\mathcal{D}| \mid \beta) \geq 0$ . Since

$$\begin{aligned} \eta &= \mathbb{E}[|[n, \infty) \setminus (D + D)| \mid \mathcal{B} = \beta] \\ &= \sum_{i=0}^{\infty} z(|\mathcal{D}| + i \mid \beta) \cdot i \\ &\geq \sum_{i=1}^{\infty} z(|\mathcal{D}| + i \mid \beta) \\ &= 1 - z(|\mathcal{D}| \mid \beta), \end{aligned} \quad (7.13)$$

we also have  $z(|\mathcal{D}| \mid \beta) \geq 1 - \eta$ .

Observe that the event  $|\mathbb{N} \setminus (D + D)| > |\mathcal{D}|$  contains the event  $\{m \notin D + D\}$ , and so

$$\mathbb{P}(|\mathbb{N} \setminus (D + D)| = |\mathcal{D}|) \leq 1 - \mathbb{P}(m \notin D + D) = 1 - \mu, \quad (7.14)$$

concluding the proof of this lemma. □

**Lemma 7.4.** *We have  $\max\{0, 2\mu - \eta\} \leq z(|\mathcal{D}| + 1 \mid \beta) \leq \min\{1, \eta\}$ .*

*Proof.* Trivially  $0 \leq z(|\mathcal{D}| + 1 \mid \beta) \leq 1$ . We have

$$\eta = \sum_{k=0}^{\infty} k \cdot z(|\mathcal{D}| + k \mid \beta) \geq z(|\mathcal{D}| + 1 \mid \beta), \quad (7.15)$$

which leaves only the bound  $2\mu - \eta \leq z(|\mathcal{D}| + 1 \mid \beta)$  to prove.

The idea here is that if exactly  $|\mathcal{D}| + 1$  sums are missing, they are very likely to be the  $|\mathcal{D}|$  elements of  $\mathcal{D}$ , and  $m$ . Formally,

$$\begin{aligned} \{|\mathbb{N} \setminus (D + D)| = |\mathcal{D}| + 1\} &\supseteq \{m \notin D + D\} \cap \bigcap_{\substack{\ell \in \mathcal{L} \\ \ell > m}} \{\ell \in D + D\} \cap \bigcap_{t \in \mathcal{T}} \{t \in D + D\} \\ &\supseteq \{m \notin D + D\} \setminus \left( \bigcup_{\substack{\ell \in \mathcal{L} \\ \ell > m}} \{\ell \notin D + D\} \cup \bigcup_{t \in \mathcal{T}} \{t \notin D + D\} \right) \end{aligned}$$

and so

$$\begin{aligned} z(|\mathcal{D}| + 1 \mid \beta) &\geq \mathbb{P}(m \notin D + D) - \sum_{\substack{\ell \in \mathcal{L} \\ \ell > m}} \mathbb{P}(\ell \notin D + D) - \sum_{t \in \mathcal{T}} \mathbb{P}(t \notin D + D) \\ &= 2\mathbb{P}(m \notin D + D) - \sum_{i \in \mathcal{L} \cup \mathcal{T}} \mathbb{P}(i \notin D + D) \\ &= 2\mu - \eta. \end{aligned} \tag{7.16}$$

□

**Lemma 7.5.** For  $k \geq 2$ ,  $0 \leq z(|\mathcal{D}| + k \mid \beta) \leq \frac{1}{k} \min\{\eta, 2\eta - 2\mu\}$ .

We note that sometimes this bound is weaker than  $z(|\mathcal{D}| + k \mid \beta) \leq 1$ . This happens for few enough  $\beta$  that, from a computational vantage point, it is not worth checking for.

*Proof.* Trivially,  $0 \leq z(|\mathcal{D}| + k \mid \beta)$ . We have

$$\eta = \sum_{i=0}^{\infty} i \cdot z(|\mathcal{D}| + i) \geq kz(|\mathcal{D}| + k), \tag{7.17}$$

whence  $z(|\mathcal{D}| + k) \leq \eta/k$ . But also,

$$\begin{aligned} \eta &= \sum_{i=0}^{\infty} i \cdot z(|\mathcal{D}| + i) \\ &= z(|\mathcal{D}| + 1) + \sum_{i=2}^{\infty} i \cdot z(|\mathcal{D}| + i) \\ &\geq 2\mu - \eta + kz(|\mathcal{D}| + k), \end{aligned} \tag{7.18}$$

and so  $z(|\mathcal{D}| + k) \geq (2\eta - 2\mu)/k$ .

□

**7.1. Making the computation feasible, reliable, and verifiable.** A massive computation has been performed, so some words are necessary as to how this is feasible. Set

$$\begin{aligned} \text{LOWER}(k \mid \beta) &:= \begin{cases} \max\{0, 2^n - 2^n \eta\}, & k = |\mathcal{D}| \\ \max\{0, 2 \cdot 2^n \mu - 2^n \eta\}, & k = |\mathcal{D}| + 1 \\ 0, & \text{otherwise} \end{cases} \\ \text{UPPER}(k \mid \beta) &:= \begin{cases} 2^n - 2^n \mu, & k = |\mathcal{D}| \\ \min\{2^n, 2^n \eta\}, & k = |\mathcal{D}| + 1 \\ 0, & \text{otherwise} \end{cases} \\ \text{OVERHANG}(k \mid \beta) &:= \begin{cases} \min\{2^n \eta, 2 \cdot 2^n \eta - 2 \cdot 2^n \mu\}, & k = |\mathcal{D}| \\ 0, & \text{otherwise.} \end{cases} \end{aligned} \quad (7.19)$$

The lemmas above imply that the vector

$$2^{2n-1} \langle z(0), z(1), \dots, z(31) \rangle = \sum_{0 \in \beta \subseteq [0, n)} 2^n \langle z(0 \mid \beta), z(1 \mid \beta), \dots, z(31 \mid \beta) \rangle \quad (7.20)$$

is bounded below componentwise by

$$\sum_{0 \in \beta \subseteq [0, n)} \langle \text{LOWER}(0 \mid \beta), \text{LOWER}(1 \mid \beta), \dots, \text{LOWER}(31 \mid \beta) \rangle \quad (7.21)$$

and is bounded above componentwise by

$$\begin{aligned} \sum_{0 \in \beta \subseteq [0, n)} \bigg( &\langle \text{UPPER}(0 \mid \beta), \text{UPPER}(1 \mid \beta), \dots, \text{UPPER}(31 \mid \beta) \rangle + \\ &\langle \text{OVERHANG}(0 \mid \beta), \text{OVERHANG}(1 \mid \beta), \dots, \text{OVERHANG}(31 \mid \beta) \rangle \cdot M \bigg), \end{aligned}$$

where  $M$  is the  $32 \times 32$  matrix whose  $(i, j)^{\text{th}}$  entry (running the indices from 0 to 31) is  $\frac{1}{j-i}$  if  $j \geq i + 2$ , and is 0 otherwise. This allows us to compute an upper bound on  $z(0), \dots, z(31)$  from

$$\sum_{0 \in \beta \subseteq [0, n)} \langle \text{UPPER}(0 \mid \beta), \text{UPPER}(1 \mid \beta), \dots, \text{UPPER}(31 \mid \beta) \rangle \quad (7.22)$$

and

$$\sum_{0 \in \beta \subseteq [0, n)} \langle \text{OVERHANG}(0 \mid \beta), \text{OVERHANG}(1 \mid \beta), \dots, \text{OVERHANG}(31 \mid \beta) \rangle. \quad (7.23)$$

Observe that LOWER, UPPER and Overhang are always integral, as  $2^n \mu$  and  $2^n \eta$  are both integers; this means that we can compute (7.21), (7.22) and (7.23) using only integer arithmetic.

We need to compute  $\beta + \beta$  and  $\beta \cap [0, k]$  (for various  $k$ ) for each  $\beta$ . This work can be tremendously reduced by using a Gray code. That is, the subsets of  $[1, n)$  can be enumerated in such a way that each set differs from its predecessor in only one element (either put in or taken out). By storing the representation function for  $\beta + \beta$  (that is, the number of times each sum can be written as a sum of two elements of  $\beta$ ), we can simply update the necessary computations instead of re-computing.

Unfortunately, the size of the computation requires us to use  $2n + 1$ -bit integers, and this is not a supported data type in most languages for  $n \geq 32$ . The options of using C with GMP, Mathematica, or some other route to arbitrary size integers is prohibited by the size of our computation and the modesty of our actual needs (we add, but never multiply, and know a priori the number of bits we will need). Therefore, we choose to represent our numbers as arrays of 64-bit integers in C++ (each element of the array represents a separate digit of the binary expansion of the number, but the digits aren't restricted to  $\{0, 1\}$ ). To further extend our reach, we ran the code on the parallel computing cluster at the High Performance Computing Cluster at the City University of New York. To facilitate parallelization, we break  $\beta$  into  $\beta_1 = \beta \cap [0, n_1)$  and  $\beta_2 = \beta \cap [n_1, n)$ . This makes the algorithm “embarrassingly parallel”, and allows us to store intermediate calculations both to recover from any system or power failings, and to allow for spot checking of results.

To ensure correctness of the results, we have written the code in Mathematica using the simplest algorithms conceivable. Such code becomes intractably slow around  $n \approx 25$ , but this provides a sequence of values against which we can test our progressively more subtly written code, both in Mathematica and in C++. Our most sophisticated code is in C++.

Finally, we have the bounds on  $\mathbb{P}(|\mathbb{N} \setminus (D + D)| = k \mid D \cap [0, 2^{10}) = \beta_1)$  for all  $\beta_1$  in a publicly available file, together with our source code. We invite the reader to spot check our implementation.

**7.2. Obtaining  $y(k)$ ,  $m(k)$ , and  $w(k)$  from  $z(k)$ .** While it is clear that  $z(k)$  is defined, that is, the event “ $|\mathbb{N} \setminus (D + D)| = k$ ” is measurable, it is less clear that  $z(\infty) = 0$ . This, and that  $y(\infty) = 0$ , follows from the Borel-Cantelli lemma and bounds such as (1.10). We can define  $D$  (a uniformly chosen subset of  $\mathbb{N}$  containing 0) as  $C - \min C$  (where  $C$  is a uniformly chosen subset of  $\mathbb{N}$ ), and so

$$\begin{aligned}
y(k) &= \sum_{i=0}^{\infty} \mathbb{P}(\min C = i \text{ AND } |\mathbb{N} \setminus (C + C)| = k) \\
&= \sum_{i=0}^{\infty} \mathbb{P}(\min C = i \text{ AND } |\mathbb{N} \setminus ((C - \min C) + (C - \min C))| = k - 2i) \\
&= \sum_{i=0}^{\lfloor k/2 \rfloor} \mathbb{P}(\min C = i) \mathbb{P}(|\mathbb{N} \setminus (D + D)| = k - 2i) \\
&= \sum_{i=0}^{\lfloor k/2 \rfloor} \frac{1}{2^{i+1}} z(k - 2i).
\end{aligned} \tag{7.24}$$

To obtain the formulas

$$m(k) = \sum_{i=0}^k y(i)y(k-i), \quad w(k) = \sum_{i=0}^k z(i)z(k-i) \tag{7.25}$$



we refer the reader to [I]. The gist of the argument is that

$$\begin{aligned}
m(k) &:= \mathbb{P}(|[0, 2n-2] \setminus (A+A)| = k) \\
&= \sum_{i=0}^k \mathbb{P}(|[0, n/2] \setminus (A+A)| = i \text{ AND } |(3n/2, 2n-2] \setminus A+A| = k-i) \\
&\quad + \mathbb{P}(A+A \text{ misses an element in } [n/2, 3n/2]) \\
&= \sum_{i=0}^k \mathbb{P}(|[0, n/2] \setminus (A+A)| = i \text{ AND } |(3n/2, 2n-2] \setminus A+A| = k-i) + O\left(\left(\frac{3}{4}\right)^{n/4}\right).
\end{aligned} \tag{7.26}$$

Since  $A+A \cap [0, n/2]$  is only affected by  $A \cap [0, n/2]$  and  $A+A \cap (3n/2, 2n-2]$  is only affected by  $A \cap (n/2, n)$ , we can use independence to write

$$m(k) = \sum_{i=0}^k \mathbb{P}(|[0, n/2] \setminus (A+A)| = i) \mathbb{P}(|(3n/2, 2n-2] \setminus A+A| = k-i) + O\left(\left(\frac{3}{4}\right)^{n/4}\right). \tag{7.27}$$

so that

$$m(k) \sim \sum_{i=0}^k \mathbb{P}(|[0, n/2] \setminus (A+A)| = i) \mathbb{P}(|(3n/2, 2n-2] \setminus A+A| = k-i). \tag{7.28}$$

As  $n \rightarrow \infty$ , the set  $[0, n/2] \setminus (A+A)$  looks more and more like  $\mathbb{N} \setminus (C+C)$ , so that

$$\mathbb{P}(|[0, n/2] \setminus (A+A)| = i) \rightarrow y(i), \tag{7.29}$$

and similarly (after replacing  $A$  with  $n-1-A$ ) for  $\mathbb{P}(|(3n/2, 2n-2] \setminus A+A| = k-i)$ . The argument for  $w(k)$  is identical, but with “ $D$ ” in place “ $C$ ”.

Let  $Z_1, Z_2$  be independent random variables with the same distribution as  $M_{\mathbb{N} \setminus \{0\}}$ , and set  $W := Z_1 + Z_2$ . Then  $\mathbb{P}(W = k) = \sum_{i=0}^k z(i)z(k-i) = w(k)$ , whence  $\sum_{i=0}^{\infty} w(i) = 1$ , and similarly  $\sum_{i=0}^{\infty} m(i) = 1$ .

Since  $y(k)$  is a linear combination of  $z(0), \dots, z(k)$  with *positive* coefficients, the lower bounds on  $z(0), \dots, z(k)$  immediately give a lower bound on  $y(k)$ , and likewise upper bounds on  $z(0), \dots, z(k)$  yield an upper bound on  $y(k)$ . The situation is the same between  $y$  and  $m$  and between  $z$  and  $w$ , even though the combination is not linear!

To experimentally estimate  $z(k)$ , we hypothesized that  $\mathbb{P}(\mathbb{N} \setminus (D+D) \not\subseteq [0, 256])$  is sufficiently small as to be ignored. Then, using Mathematica 8, we generated  $2^{28}$  pseudorandom subsets  $E$  of  $[0, 256)$ , forced each to contain 0, and then computed  $k := |[0, 256) \setminus (E+E)|$  and kept a running tally of the number of times each value of  $k$  arose. This estimates (with an enormous sample size)

$$\mathbb{P}(|\mathbb{N} \setminus (D+D)| = k \mid \mathbb{N} \setminus (D+D) \subseteq [0, 256]) \approx z(k). \tag{7.30}$$

The estimates  $\widehat{z(k)}$ , along with conservative 99.9% confidence intervals, are given in Table 16 and shown in Figure 2. The implied bounds on  $w$ ,  $m$ , and  $y$  are given in Tables 17, 18, and 19 respectively, and shown in Figure 2.

## 8. CONJECTURES AND FUTURE RESEARCH

We end with some conjectures that are supported by numerical data. Our main conjecture remains Conjecture 1.3, which says that the distribution of missing sums is approximately exponential. One possible method of studying this distribution is finding where the first present sum in  $A + A$  occurs, given that  $A + A$  has  $k$  missing elements. Recall that the lower bound in §4 was proven by constructing  $A$  such that  $M_{[0, n-1]}(A) = k$  by letting the first  $k/2$  elements of  $A$  be missing. In this case, the index of the first present sum in  $A + A$  occurs at index  $k$ . But from numerical data, the index of the first present element will not be  $k$  for typical  $A + A$  that is missing  $k$  elements. This also suggests that this trivial construction does not account for the real ‘random’ way of constructing  $A$  such that  $A + A$  is missing  $k$  elements, which is consistent with the fact that the conjectured decay constant for the distribution is approximately 0.78 but the lower bound gives only the decay constant approximately 0.70. Even though the index of the first present element is not  $k$ , from numerical data, the index seems to be linear in  $k$ .

To be precise, let

$$X_n(A) = \max\{m : \text{if } \ell < m \text{ then } \ell \notin A + A\}$$

be the index of the first present sum of  $A + A$ . Then we have the following conjecture.

**Conjecture 8.1.** *For large  $k$ ,*

$$\lim_{n \rightarrow \infty} \mathbb{E}(X_n(A) \mid M_{[0, n-1]}(A) = k) \quad (8.1)$$

*is asymptotically linear in  $k$ .*

Similarly, we can investigate how far we must move to the right of zero and to the left of the maximum possible sum,  $2n - 2$ , so that there are no missing sums of  $A + A$  in this interval. Given  $A \in [0, n - 1]$  missing exactly  $k$  sums, as  $n \rightarrow \infty$  each of the  $k$  missing elements of  $A + A$  are either near 0 or near  $2n - 2$ . Thus all of the action is happening near the two fringes, and we want to understand what is happening there. This suggests studying

$$\max\{Y_n(A) - W_n(A) : [W_n(A), Y_n(A)] \subset A + A\}.$$

**Conjecture 8.2.** *With  $W_n(A)$  as above*

$$\lim_{n \rightarrow \infty} \mathbb{E}(W_n(A) \mid M_{[0, n-1]}(A) = k) \quad (8.2)$$

*is asymptotically linear in  $k$ .*

Note a similar conjecture should hold for  $2n - 2 - Y_n(A)$ .

Another direction is to improve the exponential bounds for  $\mathbb{P}(M_{[0, n-1]}(A) = k)$ . One approach to do this is to find upper bounds on probabilities like  $\mathbb{P}(a_1, \dots, a_m \notin A + A)$  for arbitrary  $a_1, a_2, \dots, a_m$  around  $k$ .

Recall that in §4 we first used  $\mathbb{P}(i \notin A + A)$  to get an upper bound for  $\mathbb{P}(M_{[0, n-1]}(A) = k)$  of  $\Theta((3/4)^{k/2})$  and then used  $\mathbb{P}(i, j \notin A + A)$  to get a bound of  $\Theta((\phi/2)^k)$ , an improvement. Knowing  $\mathbb{P}(a_1, \dots, a_m \notin A + A)$  would result in similar improvement. Using the current approach, this would require studying the number of vertex covers for graphs that have vertices with degree  $m$  instead of 2.

Finally, note that it is possible to use the graph-theoretic approach to study higher moments of  $M_{[0, n-1]}$ . Recall that the variance was calculated by finding explicit formulas for  $\mathbb{P}(i \text{ and } j \notin A + A)$ . Similarly, the  $m^{\text{th}}$  moment can be found by finding explicit formulas for  $\mathbb{P}(a_1, \dots, a_m \notin A + A)$

for arbitrary  $a_1, \dots, a_m$ , which requires finding the number of vertex covers in certain graphs that have vertices with degree  $m$ . Note that we again need to study  $\mathbb{P}(a_1, \dots, a_m \notin A + A)$ , as we do when we try to improve the bounds for  $\mathbb{P}(M_{[0, n-1]}(A) = k)$ ; however now we need an exact formula for  $\mathbb{P}(a_1, \dots, a_m \notin A + A)$ , whereas before we just needed an upper bound.

APPENDIX A. DATA TABLES FOR DISTRIBUTIONS

$k$	rigorous lower	lower CI	$10^5 \widehat{z(k)}$	upper CI	rigorous upper
0	23532	23543	23554	23566	23535
1	17651	17634	17644	17655	17662
2	13955	13941	13950	13960	13975
3	11074	11065	11073	11082	11101
4	9233	9225	9233	9241	9266
5	6502	6502	6509	6516	6540
6	5049	5055	5061	5067	5090
7	3700	3710	3716	3721	3745
8	2687	2698	2703	2708	2733
9	1898	1910	1914	1918	1945
10	1384	1400	1404	1407	1433
11	958	973	976	979	1006
12	677	691	694	697	725
13	467	480	483	485	515
14	323	337	339	341	370
15	219	231	233	235	266
16	149	161	162	164	195
17	100	110	111	112	145
18	66	75	76	77	110
19	43	51	52	53	86
20	28	35	36	37	70
21	18	23	24	25	58
22	11	16	16	17	51
23	7	11	11	12	45
24	4	7	8	8	42
25	2	4	5	6	39
26	1	3	4	4	37
27	0	2	2	3	36
28	0	1	2	2	35
29	0	1	1	2	35
30	0	0	1	1	34
31	0	0	1	1	34

FIGURE 16. The first and last columns give our rigorous lower and upper bounds on  $10^5 z(k)$ . The second and fourth columns give the bounds of a conservative 99.9% confidence interval for  $10^5 \widehat{z(k)}$ . The middle column gives our best guess for the integer closest to  $10^5 z(k)$ , which we denote  $10^5 \widehat{z(k)}$ .

$k$	rigorous lower	lower CI	$10^5 \widehat{y(k)}$	upper CI	rigorous upper
0	11766	11771	11777	11783	11768
1	8825	8817	8822	8828	8831
2	12860	12856	12864	12871	12872
3	9950	9941	9948	9955	9966
4	11047	11041	11048	11056	11069
5	8226	8221	8228	8235	8253
6	8048	8048	8055	8062	8079
7	5963	5966	5972	5978	5999
8	5367	5373	5379	5385	5406
9	3931	3938	3943	3948	3972
10	3376	3387	3391	3396	3419
11	2444	2455	2459	2463	2489
12	2026	2039	2043	2046	2072
13	1456	1468	1471	1474	1502
14	1174	1188	1191	1193	1221
15	837	850	852	855	884
16	662	674	676	679	708
17	468	480	482	483	514
18	364	375	376	378	409
19	256	265	267	268	300
20	196	205	206	207	240
21	137	144	146	147	179
22	103	110	111	112	145
23	72	77	78	79	112
24	54	59	59	60	93
25	37	41	42	43	76
26	27	31	32	32	65
27	19	21	22	23	56
28	14	16	17	17	50
29	9	11	12	12	45
30	7	8	9	9	42
31	4	5	6	7	40

FIGURE 17. The first and last columns give our rigorous lower and upper bounds on  $10^5 y(k)$ . The second and fourth columns give the bounds of a conservative 99.9% confidence interval for  $10^5 \widehat{y(k)}$ . The middle column gives our best guess for the integer closest to  $10^5 y(k)$ , which we denote  $10^5 \widehat{y(k)}$ .

$k$	rigorous lower	lower CI	$10^5 \widehat{m(k)}$	upper CI	rigorous upper
0	1384	1385	1387	1389	1385
1	2076	2075	2078	2081	2079
2	3805	3804	3808	3813	3810
3	4611	4607	4613	4618	4619
4	6010	6005	6012	6020	6022
5	6445	6439	6447	6455	6463
6	7177	7172	7181	7191	7202
7	7138	7133	7143	7153	7170
8	7243	7240	7251	7261	7282
9	6825	6824	6835	6846	6871
10	6510	6513	6523	6534	6563
11	5892	5897	5907	5918	5951
12	5374	5382	5392	5402	5439
13	4712	4724	4733	4742	4783
14	4153	4168	4176	4185	4228
15	3551	3567	3575	3583	3629
16	3046	3064	3071	3079	3127
17	2550	2569	2576	2582	2633
18	2139	2159	2165	2172	2225
19	1759	1780	1785	1790	1846
20	1449	1469	1474	1479	1536
21	1173	1193	1198	1202	1260
22	951	970	974	978	1038
23	760	778	782	785	846
24	608	625	628	631	693
25	480	496	498	501	564
26	379	394	396	398	462
27	296	309	311	313	378
28	232	243	245	247	311
29	179	189	191	193	258
30	139	148	149	150	216
31	106	114	115	117	182

FIGURE 18. The first and last columns give our rigorous lower and upper bounds on  $10^5 m(k)$ . The second and fourth columns give the bounds of a conservative 99.9% confidence interval for  $10^5 \widehat{m(k)}$ . The middle column gives our best guess for the integer closest to  $10^5 m(k)$ , which we denote  $10^5 \widehat{m(k)}$ .

$k$	rigorous lower	lower CI	$10^5 \widehat{w(k)}$	upper CI	rigorous upper
0	5537	5543	5548	5554	5539
1	8307	8303	8312	8321	8314
2	9684	9674	9685	9696	9698
3	10138	10127	10139	10152	10162
4	10202	10190	10203	10217	10236
5	9411	9401	9414	9427	9454
6	8475	8470	8483	8497	8528
7	7384	7385	7397	7410	7445
8	6273	6279	6291	6302	6342
9	5194	5204	5215	5226	5269
10	4247	4262	4272	4282	4327
11	3405	3424	3433	3441	3490
12	2696	2718	2726	2733	2784
13	2107	2130	2137	2144	2197
14	1629	1654	1660	1666	1720
15	1245	1270	1275	1281	1337
16	943	968	973	977	1035
17	708	732	736	740	800
18	527	549	553	556	617
19	389	410	412	415	478
20	285	304	306	309	372
21	207	224	226	228	293
22	149	164	166	168	233
23	106	120	121	123	189
24	75	87	88	90	156
25	53	63	64	65	132
26	37	45	46	48	114
27	25	32	33	34	101
28	17	23	24	25	91
29	12	16	17	18	84
30	8	11	12	13	79
31	5	8	9	10	76

FIGURE 19. The first and last columns give our rigorous lower and upper bounds on  $10^5 w(k)$ . The second and fourth columns give the bounds of a conservative 99.9% confidence interval for  $10^5 \widehat{w(k)}$ . The middle column gives our best guess for the integer closest to  $10^5 w(k)$ , which we denote  $10^5 \widehat{w(k)}$ .

## REFERENCES

- [AE] Noga Alon and Paul Erdős, *An application of graph theory to additive number theory*, European J. Combin. **6** (1985), no. 3, 201–203. MR818591 (87d:11015)
- [ER] Paul Erdős and Alfréd Rényi, *Additive properties of random sequences of positive integers*, Acta Arith. **6** (1960), 83–110. MR0120213 (22 #10970)
- [F] Gregory A. Freĭman, *On the addition of finite sets*, Dokl. Akad. Nauk SSSR **158** (1964), 1038–1041 (Russian). MR0168529 (29 #5791)
- [HM] Peter Hegarty and Steven J. Miller, *When almost all sets are difference dominated*, Random Structures Algorithms **35** (2009), no. 1, 118–136, DOI 10.1002/rsa.20268. MR2532877 (2010f:11016)
- [I] Tiffany C. Inglis, *Distributions of missing sums and differences* (2007), available at [arXiv:1204.4938v1](https://arxiv.org/abs/1204.4938v1). NSERC USRA Report.
- [ILMZ] Geoffrey Iyer, Oleg Lazarev, Steven J. Miller, and Liyang Zhang, *Generalized More Sums Than Differences Sets*, Journal of Number Theory **132** (2012), no. 5, 1054–1073.
- [J] Renling Jin, *Applications of nonstandard analysis in additive number theory*, Bull. Symbolic Logic **6** (2000), no. 3, 331–341, DOI 10.2307/421059. MR1803637 (2001k:11262)
- [G] Frédéric Gilbert, *A finite problem related to the Erdős-Turan conjecture on additive bases*, preprint (2012).
- [MO] Greg Martin and Kevin O’Byrant, *Many sets have more sums than differences*, Additive combinatorics, CRM Proc. Lecture Notes, vol. 43, Amer. Math. Soc., Providence, RI, 2007, pp. 287–305. MR2359479 (2008i:11038)
- [MS] Steven J. Miller and Daniel Scheinerman, *Explicit constructions of infinite families of MSTD sets*, with Appendix 2 by Steven J. Miller and Peter Hegarty, Additive number theory, Springer, New York, 2010, pp. 229–248. MR2744760 (2012b:11041)
- [N] Melvyn B. Nathanson, *Additive number theory*, Graduate Texts in Mathematics, vol. 165, Springer-Verlag, New York, 1996. Inverse problems and the geometry of sumsets. MR1477155 (98f:11011)
- [R] Imre Z. Ruzsa, *Sumsets and structure*, Combinatorial number theory and additive group theory, Adv. Courses Math. CRM Barcelona, Birkhäuser Verlag, Basel, 2009, pp. 87–210. MR2522038 (2010m:11013)
- [S] J. Michael Steele, *Probability theory and combinatorial optimization*, CBMS-NSF Regional Conference Series in Applied Mathematics, vol. 69, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997. MR1422018 (99d:60002)
- [TV] Terence Tao and Van H. Vu, *Additive combinatorics*, Cambridge Studies in Advanced Mathematics, vol. 105, Cambridge University Press, Cambridge, 2010. Paperback edition [of MR2289012]. MR2573797
- [Z] Yufei Zhao, *Sets characterized by missing sums and differences*, J. Number Theory **131** (2011), no. 11, 2107–2134, DOI 10.1016/j.jnt.2011.05.003. MR2825117

*E-mail address:* olazarev@Princeton.edu

DEPARTMENT OF MATHEMATICS, PRINCETON UNIVERSITY, PRINCETON, NJ 08544

*E-mail address:* sjml@williams.edu, Steven.Miller.MC.96@aya.yale.edu

DEPARTMENT OF MATHEMATICS AND STATISTICS, WILLIAMS COLLEGE, WILLIAMSTOWN, MA 01267

*E-mail address:* kevin@member.ams.org

DEPARTMENT OF MATHEMATICS, CUNY, THE COLLEGE OF STATEN ISLAND AND THE GRADUATE CENTER, STATEN ISLAND, NY 10314