

# RESEARCH PROJECTS

STEVEN J. MILLER

ABSTRACT. Here is a collection of research projects, ranging from number theory to probability to statistics to random graphs.... Much of the background material is summarized from [MT-B], though most standard number theory textbooks would have these facts. Each chapter begins with a brief synopsis of the types of problems and background material needed. For more information, see the handouts on-line at [http://www.williams.edu/go/math/sjmillier/public\\_html/projects](http://www.williams.edu/go/math/sjmillier/public_html/projects)

## CONTENTS

1. Irrationality questions	1
1.1. Irrationality of $\sqrt{n}$	2
1.2. Irrationality of $\pi^2$ and the infinitude of primes	4
1.3. Transcendental numbers	5
1.4. Continued Fractions	8
2. Additive and Elementary Number Theory	16
2.1. More sums than differences sets	16
2.2. Structure of MSTD sets	23
2.3. Catalan's conjecture and products of consecutive integers	24
2.4. The $3x + 1$ Problem	28
3. Differential equations	29
4. Probability	30
4.1. Products of Poisson Random Variables	30
4.2. Sabermetrics	31
4.3. Die battles	36
4.4. Beyond the Pidgeonhole Principle	37
4.5. Differentiating identities	37
References	40

## 1. IRRATIONALITY QUESTIONS

The interplay between rational and irrational numbers leads to a lot of fun questions with surprising applications. Frequently the behavior of some system of mathematical or physical interest is wildly different if certain parameters are rational or not. We have ways to measure how irrational a number is (in a natural sense, the golden mean  $(1 + \sqrt{5})/2$  is the most irrational of all irrational numbers), and numbers that are just 'barely' irrational are hard to distinguish on a computer, which since it works only with 0s and 1s obviously can only deal with rational numbers.

We'll describe a variety of projects.

- (1) Irrationality of  $\sqrt{n}$ : Absolutely no background math needed, this project is concerned with the search for elementary and elegant proofs of irrationality.
- (2) Irrationality of  $\pi^2$  and the infinitude of primes: Multivariable calculus, elementary group theory, some combinatorics and some elementary analysis.
- (3) Transcendental numbers: Pidgeon-hole principle, some abstract algebra (minimal polynomials), factorial function and analysis.
- (4) Continued fractions: Lots of numerical investigations here requiring just simple programming (Mathematica has a lot of built in functions for these). Many of the projects require half of a course on continued fractions (I can make notes available if needed). Some of the numerical investigations require basic probability and statistics.

**1.1. Irrationality of  $\sqrt{n}$ .** If  $n$  is not a square, obviously  $\sqrt{n}$  is irrational. The most famous proof is in the special case of  $n = 2$ . Assume not, so  $\sqrt{n} = m/n$  for at least one pair of relatively prime  $m$  and  $n$ . Let  $p$  and  $q$  be such that  $\sqrt{2} = p/q$  and there is no pair with a smaller numerator. (It's a nice exercise to show such a pair exists. One solution is to use a descent argument, which you might have seen in cases of Fermat's last theorem or elliptic curves.) Then

$$\begin{aligned}\sqrt{2} &= \frac{p}{q} \\ 2q^2 &= p^2.\end{aligned}\tag{1.1}$$

We can now conclude that  $2|p$ . If we know unique factorization, the proof is immediate. If not, assume  $p = 2m + 1$  is odd. Then  $p^2 = 4m^2 + 4m + 1$  is odd as well, and hence not divisible by two. (Note: I believe I've heard that the Greeks argued along these lines, which is why their proofs stopped at something like the irrationality of  $\sqrt{17}$ , as they were looking at special cases; it would be interesting to look up how they attacked these problems.) We therefore may write  $p = 2r$  with  $0 < r < p$ . Then

$$2q^2 = p^2 = 4r^2,\tag{1.2}$$

which when we divide by 2 gives

$$q^2 = 2r^2.\tag{1.3}$$

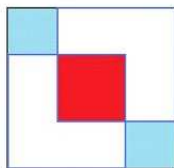
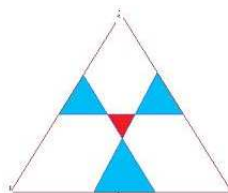
Arguing as before, we find that  $2|q$ , so we may write  $q = 2s$ . We have thus shown that

$$\sqrt{2} = \frac{p}{q} = \frac{2r}{2s} = \frac{r}{s},\tag{1.4}$$

with  $0 < r < p$ . This contradicts the minimality of  $p$ , and therefore  $\sqrt{2}$  is irrational.

On 2/9/09, Margaret Tucker gave a nice colloquium talk at Williams about proofs of the irrationality of  $\sqrt{2}$ . Among the various proofs is an ingenious one due to Conway. Assume  $\sqrt{2}$  is rational. Then there are integers  $m$  and  $n$  such that  $2m^2 = n^2$ . We quickly sketch the proof. As in the first proof, let  $m$  and  $n$  be the smallest such integers where this holds (this implies we have removed all common factors of  $m$  and  $n$ ). Then two squares of side  $m$  have the same area as a square of side  $n$ . This leads to the following picture (Figure 1):

We have placed the two squares of side length  $m$  inside the big square of side length  $n$ ; they overlap in the red region and miss the two blue regions. Thus, as the red region is

FIGURE 1. Conway's proof of the irrationality of  $\sqrt{2}$ FIGURE 2. Miller's proof of the irrationality of  $\sqrt{3}$ , with no attempt made at drawing to scale!

double counted and the area of the two squares of side  $m$  equals that of side  $n$ , we have the area of the red region equals that of the two blue regions. This leads to  $2x^2 = y^2$  for integers  $x$  and  $y$ , with  $x < m$  and  $y < n$ , contradicting the minimality of  $m$  and  $n$ . (One could easily convert this to an infinite descent argument, generating an infinite sequence of rationals.).

Professor Morgan commented on the beauty of the proof, but remarked that it is special to proving the irrationality of  $\sqrt{2}$ . The method can be generalized to handle at least one other number:  $\sqrt{3}$ . To see this, note that any equilateral triangle has area proportional to its side length  $s$  (and of course this constant is independent of  $s$ ). Assume  $\sqrt{3}$  is rational, and thus we may write  $3x^2 = y^2$ . Geometrically we may interpret this as the sum of three equilateral triangles of integral side length  $x$  equals an equilateral triangle of integral side length  $y$ . Clearly  $x < y$ , and this leads to the following picture (Figure 2):

Above we have placed the three equilateral triangles of side length  $x$  in the three corners of the equilateral triangle of side length  $y$ . Clearly  $x > y/2$  so there are intersections of these three triangles (if  $x \leq y/2$  then  $3x^2 \leq 3y^2/4 < y^2$ ). Let us color the three equilateral triangles formed where exactly two triangles intersect by blue and the equilateral triangle missed by all by red. (There must be some region missed by all, or the resulting area of the three triangles of side length  $x$  would exceed that of side length  $y$ .) Thus (picture not to scale!) the sum of the three blue triangles equals that of the red triangle. The side length of each blue triangle is  $2x - y$  and that of the red triangle  $x - 2(2x - y) = y - 3x$ , both integers. Thus we have found a smaller pair of integers (say  $a$  and  $b$ ) satisfying  $3a^2 = b^2$ , contradiction.

This leads to the following:

**Project 1.1.** *For what other integers  $k$  can we find some geometric construction along these lines proving  $\sqrt{k}$  is irrational? Or, more generally, for what positive integers  $k$  and  $r$  is  $\sqrt[r]{k}$  irrational?*

**Remark 1.2.** *I have not read Conway's paper, so I do not know what he was able to show.*

**1.2. Irrationality of  $\pi^2$  and the infinitude of primes.** Let  $\pi(x)$  count the number of primes at most  $x$ . The celebrated Prime Number Theorem states that  $\pi(x) \sim x/\log x$  for  $x$  large (even better,  $\pi(x) \sim \text{Li}(x)$ , where  $\text{Li}(x) = \int_2^x dt/\log t$ , which to first order is  $x/\log x$ ). As primes are the building blocks of integers, it is obviously important to know how many we have up to a given height.

There are numerous proofs of the infinitude of primes. Many of the proofs of the infinitude of primes fall naturally into one of two categories. First, there are those proofs which provide a lower bound for  $\pi(x)$ . A classic example of this is Chebyshev's proof that there is a constant  $c$  such that  $cx/\log x \leq \pi(x)$  (many number theory books have this proof; see for example [MT-B]). Another method of proof is to deduce a contradiction from assuming there are only finitely many primes. One of the nicest such arguments is due to Furstenberg (see [AZ]), who gives a topological proof of the infinitude of primes. As is often the case with arguments along these lines, we obtain no information about how rapidly  $\pi(x)$  grows.

Sometimes proofs which at first appear to belong to one category in fact belong to another. For example, Euclid proved there are infinitely many primes by noting the following: if not, and if  $p_1, \dots, p_N$  is a complete enumeration, then either  $p_1 \cdots p_N + 1$  is prime or else it is divisible by a prime not in our list. A little thought shows this proof belongs to the first class, as it yields there are at least  $k$  primes at most  $2^{2^k}$ , that  $\pi(x) \geq \log \log(x)$ .

For the other direction, we examine a standard 'special value' proof; see [MT-B] for proofs of all the claims below. Consider the Riemann zeta function

$$\zeta(s) := \sum_{n=1}^{\infty} \frac{1}{n^s} = \prod_{p \text{ prime}} (1 - p^{-s})^{-1},$$

which converges for  $\Re s > 1$ ; the product representation follows from the unique factorization properties of the integers. One can show  $\zeta(2) = \pi^2/6$ . As  $\pi^2$  is irrational, there must be infinitely many primes; if not, the product over primes at  $s = 2$  would be rational. While at first this argument may appear to belong to the second class (proving  $\pi(x)$  tends to infinity without an estimate of its growth), it turns out that this proof belongs to the first class, and we can obtain an explicit, though *very* weak, lower bound for  $\pi(x)$ . Unfortunately, the argument is a bit circular, for the following reason.

Our lower bounds for  $\pi(x)$  use the fact that the irrationality measure of  $\pi^2/6$  is bounded. An upper bound on the irrationality measure of an irrational  $\alpha$  is a number  $\nu$  such that there are only finitely many pairs  $p$  and  $q$  with

$$\left| \alpha - \frac{p}{q} \right| < \frac{1}{q^\nu}.$$

The irrationality measure  $\mu_{\text{irr}}(\alpha)$  is defined to be the infimum of the bounds and need not itself be a bound. Liouville constructed transcendental numbers by studying numbers with infinite irrationality measure, and Roth proved the irrationality measure of an algebraic number is 2 (see [MT-B]). Currently the best known bound is due to Rhin and Viola [RV2], who give 5.45 as a bound on the irrationality measure of  $\pi^2/6$ . Unfortunately, the published proofs of these bounds use good upper and lower bounds for  $d_n = \text{lcm}(1, \dots, n)$ . These upper and lower bounds are obtained by appealing to the Prime Number Theorem (or Chebyshev type bounds); this is a problem for us, as we are trying to prove a weaker version of the Prime Number Theorem (which we are thus subtly assuming in one of our steps!).

This leads to the following:

**Project 1.3.** *Can we prove that the irrationality measure of  $\pi^2/6$  is finite without appealing to the Prime Number Theorem, Chebyshev's Theorem, or anything along these lines?*

Even if we cannot do this, all hope is not lost in attempting to get a good lower bound on  $\pi(x)$  by studying  $\pi^2/6$ . We can open up the proof of Rhin and Viola [RV2] and see what happens if, infinitely often,  $\pi(x)$  is small. I have some notes to this affect on the webpage (there are some typos there). I think it will be possible to show the following: We say  $f(x) = o(g(x))$  if  $\lim_{x \rightarrow \infty} f(x)/g(x) = 0$ . Let  $f(x)$  be any function satisfying  $f(x) = o(x/\log x)$ . I believe one can show that infinitely often  $\pi(x) > Cf(x)$  for some  $C$ . Thus

**Project 1.4.** *Open up the proof of Rhin and Viola. See where the Prime Number Theorem / Chebyshev's theorem is used to estimate the least common multiple of  $\{1, \dots, n\}$ . Avoid using these results, and instead assume that  $\pi(x) \leq f(x)$  for all  $x$  sufficiently large. Deduce a contradiction. It is essential that their argument can be split into two parts, one part needed the least common multiple and one part independent.*

(Note: if interested, I have a copy of Rhin and Viola's paper.)

**1.3. Transcendental numbers.** While it is easy to construct irrational numbers, it is much harder to prove that a given irrational number is transcendental (even though, in a certain sense, almost every irrational number is transcendental!). Recall the following definitions:

**Definition 1.5** (Algebraic Number). *An  $\alpha \in \mathbb{C}$  is an algebraic number if it is a root of a polynomial with finite degree and integer coefficients.*

**Definition 1.6** (Transcendental Number). *An  $\alpha \in \mathbb{C}$  is a transcendental number if it is not algebraic.*

It has been known for a long time that numbers such as  $e$  and  $\pi$  are transcendental, though it is an open question as to whether or not  $e + \pi$  or  $e\pi$  is transcendental (we can show at least one is, and we expect both are). Certain numbers are readily shown to be transcendental. These special numbers are called Liouville numbers. We'll describe their form below, and why they are transcendental.

We need a definition first; though this was defined in a previous subsection, to make this part self-contained we repeat the preliminaries. Let  $\alpha$  be a real number. We desire

a rational number  $\frac{p}{q}$  such that  $\left| \alpha - \frac{p}{q} \right|$  is small. Some explanation is needed. In some sense, the size of the denominator  $q$  measures the “cost” of approximating  $\alpha$ , and we want an error that is small relative to  $q$ . For example, we could approximate  $\pi$  by  $314159/100000$ , which is accurate to 5 decimal places (about the size of  $q$ ), or we could use  $103993/33102$ , which uses a smaller denominator and is accurate to 9 decimal places (about twice the size of  $q$ )!

**Definition 1.7** (Approximation Exponent). *The real number  $\xi$  has approximation order (or exponent)  $\tau(\xi)$  if  $\tau(\xi)$  is the smallest number such that for all  $e > \tau(\xi)$  the inequality*

$$\left| \xi - \frac{p}{q} \right| < \frac{1}{q^e} \quad (1.5)$$

*has only finitely many solutions.*

Good exercises are to show that rationals have approximation exponent of 1 and irrationals have irrationality exponent at least 2 (the standard proof uses Dirichlet’s pigeon-hole principle). Another good exercise is

**Exercise 1.8** (Approximation Exponent). *Show  $\xi$  has approximation exponent  $\tau(\xi)$  if and only if for any fixed  $C > 0$  and  $e > \tau(\xi)$  the inequality*

$$\left| \xi - \frac{p}{q} \right| < \frac{C}{q^e} \quad (1.6)$$

*has only finitely many solutions with  $p, q$  relatively prime.*

**Theorem 1.9** (Liouville’s Theorem). *Let  $\alpha$  be a real algebraic number of degree  $d$ . Then  $\alpha$  is approximated by rationals to order at most  $d$ .*

*Proof.* Let

$$f(x) = a_d x^d + \cdots + a_1 x + a_0 \quad (1.7)$$

be the polynomial with relatively prime integer coefficients of smallest degree (called the **minimal polynomial** such that  $f(\alpha) = 0$ ). The condition of minimality implies that  $f(x)$  is irreducible over  $\mathbb{Z}$ . (It is a good exercise to prove this.)

In particular, as  $f(x)$  is irreducible over  $\mathbb{Q}$ ,  $f(x)$  does not have any rational roots. If it did then  $f(x)$  would be divisible by a linear polynomial  $(x - \frac{a}{b})$ . Therefore  $f$  is non-zero at every rational. Our plan is to show the existence of a rational number  $\frac{p}{q}$  such that  $f(\frac{p}{q}) = 0$ . Let  $\frac{p}{q}$  be such a candidate. Substituting gives

$$f\left(\frac{p}{q}\right) = \frac{N}{q^d}, \quad N \in \mathbb{Z}. \quad (1.8)$$

Note the integer  $N$  depends on  $p, q$  and the  $a_i$ ’s. To emphasize this dependence we write  $N(p, q; \alpha)$ . As usual, the proof proceeds by showing  $|N(p, q; \alpha)| < 1$ , which then forces  $N(p, q; \alpha)$  to be zero; this contradicts  $f$  is irreducible over  $\mathbb{Q}$ .

We find an upper bound for  $N(p, q; \alpha)$  by considering the Taylor expansion of  $f$  about  $x = \alpha$ . As  $f(\alpha) = 0$ , there is no constant term in the Taylor expansion. We may assume  $\frac{p}{q}$  satisfies  $|\alpha - \frac{p}{q}| < 1$ . Then

$$f(x) = \sum_{i=1}^d \frac{1}{i!} \frac{d^i f}{dx^i}(\alpha) \cdot (x - \alpha)^i. \quad (1.9)$$

Consequently

$$\begin{aligned} \left| f\left(\frac{p}{q}\right) \right| &= \left| \frac{N(p, q; \alpha)}{q^d} \right| \leq \left| \frac{p}{q} - \alpha \right| \cdot \sum_{i=1}^d \left| \frac{1}{i!} \frac{d^i f}{dx^i}(\alpha) \right| \cdot \left| \frac{p}{q} - \alpha \right|^{i-1} \\ &\leq \left| \frac{p}{q} - \alpha \right| \cdot d \cdot \max_i \left| \frac{1}{i!} \frac{d^i f}{dx^i}(\alpha) \right| \cdot 1^{i-1} \\ &\leq \left| \frac{p}{q} - \alpha \right| \cdot A(\alpha), \end{aligned} \quad (1.10)$$

where  $A(\alpha) = d \cdot \max_i \left| \frac{1}{i!} \frac{d^i f}{dx^i}(\alpha) \right|$ . If  $\alpha$  were approximated by rationals to order greater than  $d$ , then (Exercise 1.8) for some  $\epsilon > 0$  there would exist a constant  $B(\alpha)$  and infinitely many  $\frac{p}{q}$  such that

$$\left| \frac{p}{q} - \alpha \right| \leq \frac{B(\alpha)}{q^{d+\epsilon}}. \quad (1.11)$$

Combining yields

$$\left| f\left(\frac{p}{q}\right) \right| \leq \frac{A(\alpha)B(\alpha)}{q^{d+\epsilon}}. \quad (1.12)$$

Therefore

$$|N(p, q; \alpha)| \leq \frac{A(\alpha)B(\alpha)}{q^\epsilon}. \quad (1.13)$$

For  $q$  sufficiently large,  $A(\alpha)B(\alpha) < q^\epsilon$ . As we may take  $q$  arbitrarily large, for sufficiently large  $q$  we have  $|N(p, q; \alpha)| < 1$ . As the only non-negative integer less than 1 is 0, we find for  $q$  large that  $f\left(\frac{p}{q}\right) = 0$ , contradicting  $f$  is irreducible over  $\mathbb{Q}$ .  $\square$

We may use the above to construct transcendental numbers; see [MT-B] (among numerous other sources!) for a proof.

**Theorem 1.10** (Liouville). *The number*

$$\alpha = \sum_{m=1}^{\infty} \frac{1}{10^{m!}} \quad (1.14)$$

*is transcendental.*

This gives us one transcendental number. Can we get more?

**Project 1.11.** *Consider the binary expansion for  $x \in [0, 1)$ , namely*

$$x = \sum_{n=1}^{\infty} \frac{b_n(x)}{2^n}, \quad b_n(x) \in \{0, 1\}. \quad (1.15)$$

*For irrational  $x$  this expansion is unique. Consider the function*

$$M(x) = \sum_{n=0}^{\infty} 10^{-(b_n(x)+1)n!}. \quad (1.16)$$

*Prove for irrational  $x$  that  $M(x)$  is transcendental. Thus the above is an explicit construction for uncountably many transcendentals! Investigate the properties of this function. Is it continuous or differentiable (everywhere or at some points)? What is the*

measure of these numbers? These are “special” transcendental numbers; do they have any interesting properties?

#### 1.4. Continued Fractions.

1.4.1. *Introduction.* For many problems (such as approximations by rationals and algebraicity), the continued fraction expansion of a number provides information that is hidden in the binary or decimal expansion. There are many applications of this knowledge, ranging from digit bias in data to the behavior of the fractional parts of  $n^k\alpha$  (which arises in certain physical systems).

There are many ways to represent numbers. A common way is to use decimal or base 10 expansions. For a positive real number  $x$ ,

$$\begin{aligned} x &= x_n 10^n + x_{n-1} 10^{n-1} + \cdots + x_1 10^1 + x_0 + x_{-1} 10^{-1} + x_{-2} 10^{-2} + \cdots \\ x_i &\in \{0, 1, \dots, 9\}. \end{aligned} \quad (1.17)$$

We can obviously generalize this to an arbitrary base.

Unfortunately the decimal expansion is not ‘natural’; the universe almost surely does not care that we have 10 fingers on our hand! Thus, we want an expansion that is base-independent, and hopefully this will highlight key properties of our number.

A **Finite Continued Fraction** is a number of the form

$$a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{\ddots + \frac{1}{a_n}}}}, \quad a_i \in \mathbb{R}. \quad (1.18)$$

As  $n$  is finite, the above expression makes sense provided we never divide by 0. Since this notation is cumbersome to write, we introduce the following shorthand notations. The first is

$$a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \cdots + \frac{1}{a_n}}}. \quad (1.19)$$

A more common notation, which we often use, is

$$[a_0, a_1, \dots, a_n]. \quad (1.20)$$

We state a few standard definitions.

**Definition 1.12** (Positive Continued Fraction). *A continued fraction  $[a_0, \dots, a_n]$  is positive if each  $a_i > 0$  for  $i \geq 1$ .*

**Definition 1.13** (Digits). *If  $\alpha = [a_0, \dots, a_n]$  we call the  $a_i$  the digits of the continued fraction. Note some books call  $a_i$  the  $i^{\text{th}}$  partial quotient of  $\alpha$ .*

**Definition 1.14** (Simple Continued Fraction). *A continued fraction is simple if for each  $i \geq 1$ ,  $a_i$  is a positive integer.*

Below we mostly concern ourselves with simple continued fractions; however, in truncating infinite simple continued fractions we encounter expansions which are simple except for the last digit.



**Definition 1.15** (Convergents). Let  $x = [a_0, a_1, \dots, a_n]$ . For  $m \leq n$ , set  $x_m = [a_0, \dots, a_m]$ . Then  $x_m$  can be written as  $\frac{p_m}{q_m}$ , where  $p_m$  and  $q_m$  are polynomials in  $a_0, a_1, \dots, a_m$ . The fraction  $x_m = \frac{p_m}{q_m}$  is the  $m^{\text{th}}$  convergent of  $x$ .

There turns out to be a very simple algorithm to compute continued fraction expansions; in fact, it's basically just the famous Euclidean algorithm! We want to find integers  $a_i$  (all positive except possibly for  $a_0$ ) such that

$$x = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \dots}}. \quad (1.21)$$

Obviously  $a_0 = [x]$ , the greatest integer at most  $x$ . Then

$$x - [x] = \frac{1}{a_1 + \frac{1}{a_2 + \dots}}, \quad (1.22)$$

and the inverse is

$$x_1 = \frac{1}{x - [x]} = a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \dots}}. \quad (1.23)$$

Therefore the next digit of the continued fraction expansion is  $[x_1] = a_1$ . Then  $x_2 = \frac{1}{x_1 - [x_1]}$ , and  $[x_2] = a_2$ , and so on.

**Project 1.16.** Let  $p/q \in (0, 2]$  be a rational number. Prove it may be written as a sum of distinct rationals of the form  $1/n$  (for example,  $31/30 = 1/2 + 1/3 + 1/5$ ). Hard: is the claim still true if  $p/q > 2$ ? (I forget if this is known!)

1.4.2. *Quadratic Irrationals.* An  $x \in \mathbb{R}$  is rational if and only if  $x$  has a finite continued fraction. This is a little different than decimal expansions, as there are some infinite decimal expansions that correspond to rational numbers. Things get interesting when we look at irrational numbers.

First, some notation. By a periodic continued fraction we mean a continued fraction of the form

$$[a_0, a_1, \dots, a_k, \dots, a_{k+m}, a_k, \dots, a_{k+m}, a_k, \dots, a_{k+m}, \dots]. \quad (1.24)$$

For example,

$$[1, 2, 3, 4, 5, 6, 7, 8, 9, 7, 8, 9, 7, 8, 9, 7, 8, 9, \dots]. \quad (1.25)$$

The following theorem is one of the most important in the subject; see [MT-B] for a proof.

**Theorem 1.17** (Lagrange). A number  $x \in \mathbb{R}$  has a periodic continued fraction if and only if it satisfies an irreducible quadratic equation; i.e., there exist  $A, B, C \in \mathbb{Z}$  such that  $Ax^2 + Bx + C = 0$ ,  $A \neq 0$ , and  $x$  does not satisfy a linear equation with integer coefficients.

**Project 1.18.** Give an explicit upper bound for the constant  $M$  that arises in the proof of the above theorem in [MT-B]; the bound should be a function of the coefficients of the quadratic polynomial. Use this bound to determine an  $N$  such that we can find

three numbers  $a_{n_1}, a_{n_2}, a_{n_3}$  as in the proof with  $n_i \leq N$ . Deduce a bound for where the periodicity must begin. Similarly, deduce a bound for the length of the period. Note: I am not sure how much is known here, but it is an interesting problem seeing how the period varies with  $A, B$  and  $C$ .

We have shown that  $x$  is a quadratic irrational if and only if its continued fraction is periodic from some point onward. Thus, given any repeating block we can find a quadratic irrational. In some sense this means we completely understand these numbers; however, depending on how we traverse countable sets we can see greatly different behavior.

For example, consider the following ordered subsets of  $\mathbb{N}$ :

$$\begin{aligned} S_1 &= \{1, \mathbf{2}, 3, \mathbf{4}, 5, \mathbf{6}, 7, \mathbf{8}, 9, \mathbf{10}, 11, \mathbf{12}, \dots\} \\ S_2 &= \{1, 3, \mathbf{2}, 5, 7, \mathbf{4}, 9, 11, \mathbf{6}, 13, 15, \mathbf{8}, \dots\}. \end{aligned} \quad (1.26)$$

For  $N$  large, in the first set the even numbers make up about half of the first  $N$  numbers, while in the second set, only one-third. Simply by reordering the terms, we can adjust certain types of behavior. What this means is that, depending on how we transverse a set, we can see different limiting behaviors.

**Exercise 1.19** (Rearrangement Theorem). *Consider a sequence of real numbers  $a_n$  that is conditionally convergent but not absolutely convergent:  $\sum_{n=1}^{\infty} a_n$  exists and is finite, but  $\sum_{n=1}^{\infty} |a_n| = \infty$ ; for example,  $a_n = \frac{(-1)^n}{n}$ . Prove by re-arranging the order of the  $a_n$ 's one can obtain a new series which converges to any desired real number! Moreover, one can design a new sequence that oscillated between any two real numbers.*

Therefore, when we decide to investigate quadratic irrationals, we need to specify how the set is ordered. This is similar to our use of height functions to investigate rational numbers. One interesting set is  $\mathcal{F}_N = \{\sqrt{n} : n \leq N\}$ ; another is  $\mathcal{G}_N = \{x : ax^2 + bx + c = 0, |a|, |b|, |c| \leq N\}$ . We could fix a quadratic irrational  $x$  and study its powers  $\mathcal{H}_N = \{x^n : 0 < |n| \leq N\}$  or its multiples  $\mathcal{I}_N = \{nx : 0 < |n| \leq N\}$  or ratios  $\mathcal{J}_N = \{\frac{x}{n} : 0 < |n| \leq N\}$ .

**Remark 1.20** (Dyadic intervals). *In many applications, instead of considering  $0 < n \leq N$  one investigates  $N \leq n \leq 2N$ . There are many advantages to such studies. For  $N$  large, all elements are of a comparable magnitude. Additionally, often there are low number phenomena which do not persist at larger values: by starting the count at 1, these low values could pollute the conclusions. For example, looking at*

$$\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}, \quad (1.27)$$

*we conclude 40% of numbers are prime, and 50% of primes  $p$  also have  $p+2$  prime (i.e., start a twin prime pair); further, these percentages hold if we extend to  $\{1, \dots, 20\}$ ! Both these conclusions are false. The Prime Number Theorem states that the proportion of numbers less than  $x$  that are prime is like  $\frac{1}{\log x}$ , and heuristics (using the Circle Method) indicate the proportion that are twin primes is  $\frac{2C_2}{\log^2 x}$ , where  $C_2 \approx .66016$  is the Hardy-Littlewood twin prime constant. See [So] for further details.*

One must be very careful about extrapolations from data. A terrific example is Skewes' number. Let  $\pi(x)$  equal the number of primes at most  $x$ . A good approximation to  $\pi(x)$  is  $\text{Li}(x) = \int_2^x \frac{dt}{\log t}$ ; note to first order, this integral is  $\frac{x}{\log x}$ . By studying

tables of primes, mathematicians were led to the conjecture that  $\pi(x) < \text{Li}(x)$ . While simulations supported this claim, Littlewood proved that which of the two functions is larger changes infinitely often; his student Skewes [Sk] proved the first change occurs by  $x = 10^{10^{10^3}}$ . This bound has been significantly improved; however, one expects the first change to occur around  $10^{250}$ . See [Rie] for investigations of  $\pi(x) - \text{Li}(x)$ . Numbers this large are beyond the realm of experimentation. The moral is: for phenomena whose natural scale is logarithmic (or log-logarithmic, and so on), numerics can be very misleading.

**Project 1.21.** *Determine if possible simple closed formulas for the sets  $\mathcal{H}_N, \mathcal{I}_N$  and  $\mathcal{J}_N$  arising from  $\phi$  (the golden mean) and  $x_{m,+}$ . In particular, what can one say about  $nx_{m,+}^k$  or  $x_{m,+}^k/n$ ? How are the lengths of the periods related to  $(m, k, n)$ , and what digits occur in these sets (say for fixed  $m$  and  $k$ ,  $0 < n \leq N$ )? If  $m = 1$ ,  $x_{1,+} = \phi$ , the golden mean. For  $\frac{p}{q} \in \mathbb{Q}$ , note  $\frac{p}{q}x^k$  can be written as  $\frac{p_1}{q_1}\sqrt{5} + \frac{p_2}{q_2}$ . Thus, in some sense, it is sufficient to study  $\frac{p_1}{q_1}\sqrt{5}$ .*

**Remark 1.22** (Important). *Many of the formulas for the continued fraction expansions were first seen in numerical experiments. The insights that can be gained by investigating cases on the computer cannot be underestimated. Often seeing the continued fraction for some cases leads to an idea of how to do the calculations, and at least as importantly what calculations may be interesting and worthwhile. For example,*

$$\begin{aligned}
\frac{\sqrt{5}}{4} &= [0, 1, \overline{1, 3, 1, 2}] \\
\frac{\sqrt{5}}{8} &= [0, 3, \overline{1, 1, 2, 1, 2, 1, 1, 6}] \\
\frac{\sqrt{5}}{16} &= [0, 7, \overline{6, 2, 3, 3, 3, 2, 6, 14, 6, 2, 3, 3, 3, 2, 6, 14}] \\
\frac{\sqrt{5}}{10} &= [0, 4, \overline{2, 8}] \\
\frac{\sqrt{5}}{6} &= [0, 2, \overline{1, 2, 6, 2, 1, 4}] \\
\frac{\sqrt{5}}{12} &= [0, 5, \overline{2, 1, 2, 1, 2, 10}] \\
\frac{\sqrt{5}}{14} &= [0, 6, \overline{3, 1, 4, 1, 14, 1, 4, 1, 3, 12}] \\
\frac{\sqrt{5}}{28} &= [0, 12, \overline{1, 1, 10, 1, 6, 1, 10, 1, 1, 24}] \\
\frac{\sqrt{5}}{42} &= [0, 18, \overline{1, 3, 1, 1, 1, 1, 4, 1, 1, 1, 1, 3, 1, 36}]. \tag{1.28}
\end{aligned}$$

**Project 1.23.** *The data in Remark 1.22 seem to indicate a pattern between the length of the repeating block and the factorization of the denominator, as well as what the largest digit is. Discover and prove interesting relations. How are the digits distributed (i.e., how many are 1's, 2's, 3's and so on. Also, the periodic expansions*

are almost symmetric (if one removes the final digit, the remaining piece is of the form  $abc\dots xyzyx\dots cba$ ). Is this always true? What happens if we divide by other  $n$ , say odd  $n$ ?

**Project 1.24.** How are the continued fractions of  $n$ -equivalent numbers related? We have seen quadratic irrationals have periodic continued fractions. Consider the following generalization. Fix functions  $f_1, \dots, f_k$ , and study numbers of the form

$$[f_1(1), \dots, f_k(1), f_1(2), \dots, f_k(2), f_1(3), \dots, f_k(3), f_1(4), \dots]. \quad (1.29)$$

Which numbers have such expansions (say if the  $f_i$ 's are linear)? See [Di] for some results. For results on multiplying continued fractions by rationals see [vdP1], and see [PS1, PS2, vdP3] for connections between power series and continued fractions.

**Project 1.25.** For more on the lengths of the period of  $\sqrt{n}$  or  $\sqrt{p}$ , as well as additional topics to investigate, see [Bec, GI]. For a generalization to what has been called "linearly periodic" expansions, see [Di].

1.4.3. *More on digits of continued fractions.* We start with an easily stated but I believe still wide open problem:

**Project 1.26** (Davenport). Determine whether the digits of the continued fraction expansion of  $\sqrt[3]{2} = [1, 3, 1, 5, 1, 1, 4, 1, \dots]$  are bounded or not. This problem appears on page 107 of [Da1].

Given  $\alpha \in \mathbb{R}$ , we can calculate its continued fraction expansion and investigate the distribution of its digits. Without loss of generality we assume  $\alpha \in (0, 1)$ , as this shift changes only the zeroth digit. Thus

$$\alpha = [0, a_1, a_2, a_3, a_4, \dots]. \quad (1.30)$$

Given any sequence of positive integers  $a_i$ , we can construct a number  $\alpha$  with these as its digits. However, for a generic  $\alpha$  chosen uniformly in  $(0, 1)$ , how often do we expect to observe the  $n^{\text{th}}$  digit in the continued fraction expansion equal to 1? To 2? To 3? And so on.

If  $\alpha \in \mathbb{Q}$  then it has a finite continued fraction expansion; if  $\alpha$  is a quadratic irrational then its continued fraction expansion is periodic. In both of these cases there are really only finitely many digits; however, if we stay away from rationals and quadratic irrationals, then  $\alpha$  will have a bona fide infinite continued fraction expansion, and it makes sense to ask the above questions.

For the decimal expansion of a generic  $\alpha \in (0, 1)$ , we expect each digit to take the values 0 through 9 with equal probability; as there are infinitely many values for the digits of a continued fraction, each value cannot be equally likely. We will see, however, that as  $n \rightarrow \infty$  the probability of the  $n^{\text{th}}$  digit equalling  $k$  converges to  $\log_2 \left(1 + \frac{1}{k(k+2)}\right)$ . An excellent source is [Kh].

For notational convenience, we adopt the following convention. Let  $A_{1, \dots, n}(a_1, \dots, a_n)$  be the event that  $\alpha \in [0, 1)$  has its continued fraction expansion  $\alpha = [0, a_1, \dots, a_n, \dots]$ . Similarly  $A_{n_1, \dots, n_k}(a_{n_1}, \dots, a_{n_k})$  is the event where the zeroth digit is 0, digit  $n_1$  is  $a_{n_1}$ ,  $\dots$ , and digit  $n_k$  is  $a_{n_k}$ , and  $A_n(k)$  is the event that the zeroth digit is 0 and the  $n^{\text{th}}$  digit is  $k$ .

Gauss conjectured that as  $n \rightarrow \infty$  the probability that the  $n^{\text{th}}$  digit equals  $k$  converges to  $\log_2 \left( 1 + \frac{1}{k(k+2)} \right)$ . In 1928, Kuzmin proved Gauss' conjecture, with an explicit error term:

**Theorem 1.27** (Gauss-Kuzmin). *There exist positive constants  $A$  and  $B$  such that*

$$\left| A_n(k) - \log_2 \left( 1 + \frac{1}{k(k+2)} \right) \right| \leq \frac{A}{k(k+1)} e^{-B\sqrt{n-1}}. \quad (1.31)$$

This is clearly compatible with Gauss' conjecture, as for  $B > 0$  the expression  $e^{-B\sqrt{n-1}}$  tends to zero when  $n$  approaches  $+\infty$ . The error term has been improved by Lévy [Le] to  $Ae^{-Cn}$ , and then further by Wirsing [Wir].

See [Kh, MT-B] for a proof. It is important to note that the digits are *not* independent; the probability of observing a 1 followed by a 2 is *not* the product of the two probabilities! See [MT-B] for this calculation.

**Project 1.28.** *Assign explicit values to the constants  $A$  and  $B$  in the Gauss-Kuzmin Theorem, or find  $A_0, B_0, N_0$  such that for all  $n \geq N_0$ , one may take  $A = A_0$  and  $B = B_0$ . Note: I'm not sure if this has been done, but it would be nice to have explicit constants.*

There are many open questions concerning the digits of a generic continued fraction expansion. We know the digits in the continued fraction expansions of rationals and quadratic irrationals do not satisfy the Gauss-Kuzmin densities in the limits; in the first case there are only finitely many digits, while in the second the expansion is periodic. What can one say about the structure of the set of  $\alpha \in [0, 1)$  whose distribution of digits satisfy the Gauss-Kuzmin probabilities? We know such a set has measure 1, but what numbers are in this set?

The set of algebraic numbers is countable, hence of measure zero. Thus it is possible for the digits of every algebraic number to violate the Gauss-Kuzmin law. Computer experimentation, however, indicates that the digits of algebraic numbers do seem to follow the Gauss-Kuzmin probabilities (except for quadratic irrationals, of course). The following subsets of real algebraic numbers were extensively tested by students at Princeton (where the number of digits with given values were compared with the predictions from the Gauss-Kuzmin Theorem, and in some cases pairs and triples were also compared) and shown to have excellent agreement with predictions:  $\sqrt[n]{p}$  for  $p$  prime and  $n \leq 5$  ([Ka, Law1, Mic1]) and roots of polynomials with different Galois groups ([AB]). To analyze the data from such experiments, one should perform basic hypothesis testing. For some results on numbers whose digits violate the Gauss-Kuzmin Law, see [Mic2].

**Project 1.29.** *Investigate the digits of other families of algebraic numbers, for example, the positive real roots of  $x^n - p = 0$  (see the mentioned student reports for more details and suggestions). Alternatively, for a fixed real algebraic number  $\alpha$ , one can investigate its powers or rational multiples. There are two different types of experiments one can perform. First, one can fix a digit, say the millionth digit, and examine its value as we vary the algebraic number. Second, one could look at the same large block of digits for an algebraic number, and then vary the algebraic number.*

While similar, there are different features in the two experiments. In the first we are checking digit by digit. For a fixed number, its  $n^{\text{th}}$  digit is either  $k$  or not; thus, the only probabilities we see are 0 or 1. To have a chance of observing the Gauss-Kuzmin probabilities, we need to perform some averaging (which is accomplished by looking at roots of many different polynomials).

For the second, since we are looking at a large block of digits there is already a chance of observing probabilities close to the Gauss-Kuzmin predictions. For each root and each value (or pairs of values and so on), we obtain a probability in  $[0, 1]$ . One possibility is to perform a second level of averaging by averaging these numbers over roots of different polynomials. Another possibility is to construct a histogram plot of the probabilities for each value. This allows us to investigate more refined questions. For example, are the probabilities as likely to undershoot the predicted values as overshoot? How does that depend on the value? How are the observed probabilities for the different values for each root distributed about the predictions: does it look like a uniform distribution or a normal distribution?

**Remark 1.30.** If one studies say  $x^3 - p = 0$ , as we vary  $p$  the first few digits of the continued fraction expansions of  $\sqrt[3]{p}$  are often similar. For example,

$$\begin{aligned}\sqrt[3]{10000000087} &= [1000, 34482, 1, 3, 6, 4, \dots] \\ \sqrt[3]{10000000093} &= [1000, 32258, 15, 3, 1, 3, 1, \dots] \\ \sqrt[3]{10000000097} &= [1000, 30927, 1, 5, 10, 19, \dots].\end{aligned}\tag{1.32}$$

The zeroth digit is 1000, which isn't surprising as these cube roots are all approximately  $10^3$ . Note the first digit in the continued fraction expansions is about 30000 for each. Hence if we know the continued fraction expansion for  $\sqrt[3]{p}$  for one prime  $p$  around  $10^9$ , then we have some idea of the first few digits of  $\sqrt[3]{q}$  for primes  $q$  near  $p$ . Thus if we were to look at the first digit of the cube roots of ten thousand consecutive primes near  $10^9$ , we would not expect to see the Gauss-Kuzmin probabilities.

Consider a large number  $n_0$ . Primes near it can be written as  $n_0 + x$  for  $x$  small. Then

$$\begin{aligned}(n_0 + x)^{1/3} &= n_0^{1/3} \cdot \left(1 + \frac{x}{n_0}\right)^{1/3} \\ &\approx n_0^{1/3} \cdot \left(1 + \frac{1}{3} \frac{x}{n_0}\right) \\ &= n_0^{1/3} + \frac{x}{3n_0^{2/3}}.\end{aligned}\tag{1.33}$$

If  $n_0$  is a perfect cube, then for small  $x$  relative to  $n_0$  we see these numbers should have a large first digit. Thus, if we want investigate cube roots of lots of primes  $p$  that are approximately the same size, the first few digits are not independent as we vary  $p$ . In many of the experiments digits 50,000 to 1,000,000 were investigated: for cube roots of primes of size  $10^9$ , this was sufficient to see independent behavior (though ideally one should look at autocorrelations to verify this claim). Also, the Gauss-Kuzmin Theorem describes the behavior for  $n$  large; thus, it is worthwhile to throw away the first few digits so we only study regions where the error term is small.

There are many special functions in number theory. If we evaluate countably many special functions at countably many points, we again obtain a countable set of measure 0. Thus, all these numbers' digits could violate the Gauss-Kuzmin probabilities. Experiments have shown, however, that special values of  $\Gamma(s)$  at rational arguments ([Ta]) and the Riemann zeta function  $\zeta(s)$  at positive integers ([Kua]) seem to follow the Gauss-Kuzmin probabilities.

**Project 1.31.** *Consider the non-trivial zeros of  $\zeta(s)$ , or, more generally, the zeros of any  $L$ -function. Do the digits follow the Gauss-Kuzmin distribution? For the Fourier coefficients of an elliptic curve,  $a_p = 2 \cos(\theta_p)$ ; how are the digits of  $\theta_p$  distributed? How are the digits of  $\log n$  distributed? How are the digits of  $2^{\sqrt{n}}$  distributed for  $n$  square-free?*

We know quadratic irrationals are periodic, and hence cannot follow Kuzmin's Law. Only finitely many numbers occur in the continued fraction expansion. Thus, only finitely many numbers have a positive probability of occurring in the expansion, but the Gauss-Kuzmin probabilities are positive for all positive integers.

**Project 1.32.** *What if we consider a family of quadratic irrationals with growing period? As the size of the period grows, does the distribution of digits tend to the Gauss-Kuzmin probabilities? See the warnings in Project 1.21 for more details.*

1.4.4. *Famous continued fraction expansions.* Finally, we would be remiss if we did not mention some famous continued fraction expansions. Often a special number whose decimal expansion seems random has a continued fraction expansion with a very rich structure. For example, compare the first 25 digits for  $e$ :

$$\begin{aligned} e &= 2.718281828459045235360287 \dots \\ &= [2, 1, 2, 1, 1, 4, 1, 1, 6, 1, 1, 8, 1, 1, 10, 1, 1, 12, 1, 1, 14, 1, 1, 16, 1, \dots]. \end{aligned}$$

For  $\pi$ , the positive simple continued fraction does not look particularly illuminating:

$$\begin{aligned} \pi &= 3.141592653589793238462643 \dots \\ &= [3, 7, 15, 1, 292, 1, 1, 1, 2, 1, 3, 1, 14, 2, 1, 1, 2, 2, 2, 2, 1, 84, 2, 1, 1, \dots]. \end{aligned}$$

If, however, we drop the requirement that the expansions are simple, the story is quite different. One nice expression for  $\pi$  is

$$\frac{4}{\pi} = 1 + \frac{1^2}{2 + \frac{3^2}{2 + \frac{5^2}{2 + \dots}}}. \quad (1.34)$$

There are many different types of non-simple expansions, leading to some of the most beautiful formulas in mathematics. For example,

$$e = 2 + \frac{1}{1 + \frac{1}{2 + \frac{2}{3 + \frac{3}{4 + \dots}}}}. \quad (1.35)$$

For some nice articles and simple and non-simple continued fraction expansions, see the entry at <http://mathworld.wolfram.com/> (in particular, the entries on  $\pi$  and  $e$ ).

**Project 1.33.** *Try to generalize as many properties as possible from simple continued fractions to non-simple. Clearly, numbers do not have unique expansions unless we specify exactly what the “numerators” of the non-simple expansions must be. One often writes such expansions in the more economical notation*

$$\frac{a}{\alpha+} \frac{b}{\beta+} \frac{c}{\gamma+} \frac{d}{\delta+} \dots \quad (1.36)$$

*For what choices of  $a, b, c, \dots$  and  $\alpha, \beta, \gamma, \dots$  will the above converge? How rapidly will it converge? Are there generalizations of the recurrence relations? How rapidly do the numerators and denominators of the rationals formed by truncating these expansions grow?*

## 2. ADDITIVE AND ELEMENTARY NUMBER THEORY

One of the reasons I love number theory is how easy it is to state the problems. One does not need several graduate classes to understand the formulation (though these are useful in understanding partial results!). Some of the most famous that have defied solution to this day are Goldbach’s Problem (every sufficiently large even number is the sum of two primes, where sufficiently large is believed to mean at least 4) to the Twin Prime Conjecture (there are infinitely many primes  $p$  such that  $p + 2$  is also prime). The Circle Method provides a powerful way to conjecture answers for such questions; sieving (inclusion-exclusion) can often give bounds. For example, if  $\pi_2(x)$  is the number of primes at most  $x$ , Brun proved  $\pi_2(x) \leq Cx/\log^2 x$  for some  $c$ . This allows us to deduce that the sum of the reciprocals of the twin primes converges. We call this sum Brun’s constant, and it was how the pentium bug was discovered [Ni1, Ni2].

Below are a variety of problems related to additive and elementary number theory.

- (1) More sums than differences: some of the projects are very elementary, some require deep results from analysis for full generality. There are also numerical projects related to trying to find sets with certain projects.
- (2) Products being a perfect power: Some of these questions are quite elementary and require only factorization of polynomials, while others require knowledge of elliptic curves (especially the Mordell-Weil group of rational solutions and the Birch and Swinnerton-Dyer conjecture).
- (3)  $3x+1$  problem: An algorithm to help prove the  $3x+1$  conjecture was developed by two of my former students. Their paper has a lot of small errors and vague wording; it is a very doable project (I believe!) to clean this up. A rough draft is already written, with numerous comments from me on what needs to be fixed. Basic combinatorics should suffice, though being able to write computer programs would be a tremendous asset.

### 2.1. More sums than differences sets.



2.1.1. *Introduction.* Let  $S$  be a subset of the integers. We define the sumset  $S + S$  and difference set  $S - S$  by

$$\begin{aligned} S + S &= \{s_1 + s_2 : s_i \in S\} \\ S - S &= \{s_1 - s_2 : s_i \in S\}, \end{aligned} \tag{2.1}$$

and denote the cardinality of a set  $A$  by  $|A|$ . As addition is commutative and subtraction is not, a typical pair of integers generates two differences but only one sum. It is therefore reasonable to expect a generic finite set  $S$  will have a larger difference set than sumset. We say a set is *sum dominated* (such sets are also called *more sums than differences*, or MSTD, sets) if the cardinality of its sumset exceeds that of its difference set. If the two cardinalities are equal we say the set is *balanced*, otherwise *difference dominated*. Sum dominated sets exist: consider for example  $\{0, 2, 3, 4, 7, 11, 12, 14\}$ . Nathanson wrote “*Even though there exist sets  $A$  that have more sums than differences, such sets should be rare, and it must be true with the right way of counting that the vast majority of sets satisfies  $|A - A| > |A + A|$ .*”

Recently Martin and O’Bryant [MO] showed there are many sum dominated sets. Specifically, let  $I_N = \{0, \dots, N\}$ . They prove the existence of a universal constant  $\kappa_{SD} > 0$  such that, for any  $N \geq 14$ , at least  $\kappa_{SD} \cdot 2^{N+1}$  subsets of  $I_N$  are sum dominated (there are no sum dominated sets in  $I_{13}$ ). Their proof is based on choosing a subset of  $I_N$  by picking each  $n \in I_N$  independently with probability  $1/2$ . The argument can be generalized to independently picking each  $n \in I_N$  with any probability  $p \in (0, 1)$ , and yields the existence of a constant  $\kappa_{SD,p} > 0$  such that, as  $N \rightarrow \infty$ , a randomly chosen (with respect to this model) subset is sum dominated with probability at least  $\kappa_{SD,p}$ . Similarly one can prove there are positive constants  $\kappa_{DD,p}$  and  $\kappa_{B,p}$  for the probability of having a difference dominated or balanced set.

While the authors remark that, perhaps contrary to intuition, sum dominated sets are ubiquitous, their result is a consequence of how they choose a probability distribution on the space of subsets of  $I_N$ . Suppose  $p = 1/2$ , as in their paper. With high probability a randomly chosen subset will have  $N/2$  elements (with errors of size  $\sqrt{N}$ ). Thus the density of a generic subset to the underlying set  $I_N$  is quite high, typically about  $1/2$ . Because it is so high, when we look at the sumset (resp., difference set) of a typical  $A$  there are many ways of expressing elements as a sum (resp., difference) of two elements of  $A$ . For example (see [MO]), if  $k \in A + A$  then there are roughly  $N/4 - |N - k|/4$  ways of writing  $k$  as a sum of two elements in  $A$  (similarly, if  $k \in A - A$  there are roughly  $N/4 - |k|/4$  ways of writing  $k$  as a difference of two elements of  $A$ ). This enormous redundancy means almost all numbers which can be in the sumset or difference set are. In fact, using uniform density on the subsets of  $I_N$  (i.e., taking  $p = 1/2$ ), Martin and O’Bryant show that the average value of  $|A + A|$  is  $2N - 9$  and that of  $|A - A|$  is  $2N - 5$  (note each set has at most  $2N + 1$  elements). In particular, it is only for  $k$  near extremes that we have high probability of not having  $k$  in an  $A + A$  or an  $A - A$ . In [MO] they prove a positive percentage of subsets of  $I_N$  (with respect to the uniform density) are sum dominated sets by specifying the fringe elements of  $A$ . Similar conclusions apply for any value of  $p > 0$ .

Two fascinating questions to investigate are (1) what happens if  $p$  depends on  $N$ , and (2) can one come up with explicit constructions of MSTD sets?

2.1.2. *Sum dominated sets in non-uniform models.* At the end of their paper, Martin and O’Byrant conjecture that if, on the other hand, the parameter  $p$  is a function of  $N$  tending to zero arbitrarily slowly, then as  $N \rightarrow \infty$  the probability that a randomly chosen subset of  $I_N$  is sum dominated should also tend to zero. Recently Hegarty and Miller proved this conjecture. Specifically, they showed

**Theorem 2.1.** *Let  $p : \mathbb{N} \rightarrow (0, 1)$  be any function such that*

$$N^{-1} = o(p(N)) \quad \text{and} \quad p(N) = o(1). \quad (2.2)$$

*For each  $N \in \mathbb{N}$  let  $A$  be a random subset of  $I_N$  chosen according to a binomial distribution with parameter  $p(N)$ . Then, as  $N \rightarrow \infty$ , the probability that  $A$  is difference dominated tends to one.*

*More precisely, let  $\mathcal{S}, \mathcal{D}$  denote respectively the random variables  $|A + A|$  and  $|A - A|$ . Then the following three situations arise :*

(i)  $p(N) = o(N^{-1/2})$  : Then

$$\mathcal{S} \sim \frac{(N \cdot p(N))^2}{2} \quad \text{and} \quad \mathcal{D} \sim 2\mathcal{S} \sim (N \cdot p(N))^2. \quad (2.3)$$

(ii)  $p(N) = c \cdot N^{-1/2}$  for some  $c \in (0, \infty)$  : Define the function  $g : (0, \infty) \rightarrow (0, 2)$  by

$$g(x) := 2 \left( \frac{e^{-x} - (1-x)}{x} \right). \quad (2.4)$$

Then

$$\mathcal{S} \sim g\left(\frac{c^2}{2}\right)N \quad \text{and} \quad \mathcal{D} \sim g(c^2)N. \quad (2.5)$$

(iii)  $N^{-1/2} = o(p(N))$  : Let  $\mathcal{S}^c := (2N + 1) - \mathcal{S}$ ,  $\mathcal{D}^c := (2N + 1) - \mathcal{D}$ . Then

$$\mathcal{S}^c \sim 2 \cdot \mathcal{D}^c \sim \frac{4}{p(N)^2}. \quad (2.6)$$

Theorem 2.1 proves the conjecture in [MO] and re-establishes the validity of Nathanson’s claim in a broad setting. It also identifies the function  $N^{-1/2}$  as a *threshold function* for the ratio of the size of the difference- to the sumset for a random set  $A \subseteq I_N$ . Below the threshold, this ratio is almost surely  $2 + o(1)$ , above it almost surely  $1 + o(1)$ . Part (ii) tells us that the ratio decreases continuously (a.s.) as the threshold is crossed. Below the threshold, part (i) says that most sets are ‘nearly Sidon sets’, that is, most pairs of elements generate distinct sums and differences. Above the threshold, most numbers which can be in the sumset (resp., difference set) usually are, and in fact most of these in turn have many different representations as a sum (resp., a difference). However the sumset is usually missing about twice as many elements as the difference set. Thus if we replace ‘sums’ (resp., ‘differences’) by ‘missing sums’ (resp., ‘missing differences’), then there is still a symmetry between what happens on both sides of the threshold.

The proof in general uses recent strong concentration results, though if  $p(N) = o(N^{-1/2})$  Chebyshev’s theorem from probability suffices. The theorem can be generalized to arbitrary bilinear forms:

**Theorem 2.2.** *Let  $p : \mathbb{N} \rightarrow (0, 1)$  be a function satisfying (2.2). Let  $u, v$  be non-zero integers with  $u \geq |v|$ ,  $\text{GCD}(u, v) = 1$  and  $(u, v) \neq (1, 1)$ . Put  $f(x, y) := ux + vy$ . For a positive integer  $N$ , let  $A$  be a random subset of  $I_N$  obtained by choosing each  $n \in I_N$  independently with probability  $p(N)$ . Let  $\mathcal{D}_f$  denote the random variable  $|f(A)|$ . Then the following three situations arise :*

(i)  $p(N) = o(N^{-1/2})$  : Then

$$\mathcal{D}_f \sim (N \cdot p(N))^2. \quad (2.7)$$

(ii)  $p(N) = c \cdot N^{-1/2}$  for some  $c \in (0, \infty)$  : Define the function  $g_{u,v} : (0, \infty) \rightarrow (0, u + |v|)$  by

$$g_{u,v}(x) := (u + |v|) - 2|v| \left( \frac{1 - e^{-x}}{x} \right) - (u - |v|)e^{-x}. \quad (2.8)$$

Then

$$\mathcal{D}_f \sim g_{u,v} \left( \frac{c^2}{u} \right) N. \quad (2.9)$$

(iii)  $N^{-1/2} = o(p(N))$  : Let  $\mathcal{D}_f^c := (u + |v|)N - \mathcal{D}_f$ . Then

$$\mathcal{D}_f^c \sim \frac{2u|v|}{p(N)^2}. \quad (2.10)$$

Here is a sample of issues which could be the subject of further investigations.

**Project 2.3.** *One unresolved matter is the comparison of arbitrary difference forms in the range where  $N^{-3/4} = O(p)$  and  $p = O(N^{-3/5})$ . Here the problem is that the binomial model itself does not prove of any use. This provides, more generally, motivation for looking at other models. Obviously one could look at the so-called uniform model on subsets (see [JER]), but this seems a more awkward model to handle. Note that the property of one binary form dominating another is not monotone, or even convex.*

**Project 2.4.** *A very tantalizing problem is to investigate what happens while crossing a sharp threshold.*

**Project 2.5.** *One can ask if the various concentration estimates in Theorem 2.1 can be improved. When  $p = o(N^{-1/2})$  we have only used an ordinary second moment argument, and it is possible to provide explicit estimates. The range  $N^{-1/2} = o(p(N))$  seems more interesting, however. Here we proved that the random variable  $\mathcal{S}^c$  has expectation of order  $P(N)^2$ , where  $P(N) = 1/p(N)$ , and is concentrated within  $P(N)^{3/2} \log^2 P(N)$  of its mean. Now one can ask whether the constant  $3/2$  can be improved, or at the very least can one get rid of the logarithm?*

**Project 2.6.** *It is natural to ask for extensions of our results to  $\mathbb{Z}$ -linear forms in more than two variables. Let*

$$f(x_1, \dots, x_k) = u_1 x_1 + \dots + u_k x_k, \quad u_i \in \mathbb{Z}_{\neq 0}, \quad (2.11)$$

*be such a form. We conjecture the following generalization of Theorem 3.1 :*

**Conjecture 2.7.** *Let  $p : \mathbb{N} \rightarrow (0, 1)$  be a function satisfying (2.2). For a positive integer  $N$ , let  $A$  be a random subset of  $I_N$  obtained by choosing each  $n \in I_N$  independently with probability  $p(N)$ . Let  $f$  be as in (4.1) and assume that  $\text{GCD}(u_1, \dots, u_n) = 1$ . Set*

$$\theta_f := \#\{\sigma \in S_k : (u_{\sigma(1)}, \dots, u_{\sigma(k)}) = (u_1, \dots, u_k)\}. \quad (2.12)$$

Let  $\mathcal{D}_f$  denote the random variable  $|f(A)|$ . Then the following three situations arise :

(i)  $p(N) = o(N^{-1/k})$  : Then

$$\mathcal{D}_f \sim \frac{1}{\theta_f} (N \cdot p(N))^k. \quad (2.13)$$

(ii)  $p(N) = c \cdot N^{-1/k}$  for some  $c \in (0, \infty)$  : There is a rational function  $R(x_0, \dots, x_k)$  in  $k + 1$  variables, which is increasing in  $x_0$ , and an increasing function

$g_{u_1, \dots, u_k} : (0, \infty) \rightarrow (0, \sum_{i=1}^k |u_i|)$  such that

$$\mathcal{D}_f \sim g_{u_1, \dots, u_k}(R(c, u_1, \dots, u_k)) \cdot N. \quad (2.14)$$

(iii)  $N^{-1/k} = o(p(N))$  : Let  $\mathcal{D}_f^c := \left(\sum_{i=1}^k |u_i|\right) N - \mathcal{D}_f$ . Then

$$\mathcal{D}_f^c \sim \frac{2\theta_f \prod_{i=1}^k |u_i|}{p(N)^k}. \quad (2.15)$$

2.1.3. *Creating dense families of sum dominated sets.* Though MSTD sets are rare, they do exist (and, in the uniform model, are somewhat abundant by the work of Martin and O’Bryant). Examples go back to the 1960s. Conway is said to have discovered  $\{0, 2, 3, 4, 7, 11, 12, 14\}$ , while Marica gave  $\{0, 1, 2, 4, 7, 8, 12, 14, 15\}$  in 1969 and Freiman and Pigarev found  $\{0, 1, 2, 4, 5, 9, 12, 13, 14, 16, 17, 21, 24, 25, 26, 28, 29\}$  in 1973. Recent work includes infinite families constructed by Hegarty [He] and Nathanson [Na2], as well as existence proofs by Ruzsa [Ru1, Ru2, Ru3].

Most of the previous constructions<sup>1</sup> of infinite families of MSTD sets start with a symmetric set which is then ‘perturbed’ slightly through the careful addition of a few elements that increase the number of sums more than the number of differences; see [He, Na2] for a description of some previous constructions and methods. In many cases, these symmetric sets are arithmetic progressions; such sets are natural starting points because if  $A$  is an arithmetic progression, then  $|A + A| = |A - A|$ .<sup>2</sup>

We present a new method (by Miller-Orosz-Scheinerman) which takes an MSTD set satisfying certain conditions and constructs an infinite family of MSTD sets. While these families are not dense enough to prove a positive percentage of subsets of  $\{1, \dots, r\}$  are MSTD sets, we are able to elementarily show that the percentage is at least  $C/r^4$  for some constant  $C$ . Thus our families are far denser than those in [He, Na2]; trivial counting<sup>3</sup> shows all of their infinite families give at most  $f(r)2^{r/2}$  of the subsets

<sup>1</sup>An alternate method constructs an infinite family from a given MSTD set  $A$  by considering  $A_t = \{\sum_{i=1}^t a_i m^{i-1} : a_i \in A\}$ . For  $m$  sufficiently large, these will be MSTD sets; this is called the base expansion method. Note, however, that these will be very sparse. See [He] for more details.

<sup>2</sup>As  $|A + A|$  and  $|A - A|$  are not changed by mapping each  $x \in A$  to  $\alpha x + \beta$  for any fixed  $\alpha$  and  $\beta$ , we may assume our arithmetic progression is just  $\{0, \dots, n\}$ , and thus the cardinality of each set is  $2n + 1$ .

<sup>3</sup>For example, consider the following construction of MSTD sets from [Na2]: let  $m, d, k \in \mathbb{N}$  with  $m \geq 4$ ,  $1 \leq d \leq m - 1$ ,  $d \neq m/2$ ,  $k \geq 3$  if  $d < m/2$  else  $k \geq 4$ . Set  $B = [0, m - 1] \setminus \{d\}$ ,  $L =$

of  $\{1, \dots, r\}$  (for some polynomial  $f(r)$ ) are MSTD sets, implying a percentage of at most  $f(r)/2^{r/2}$ .

We first introduce some notation. The first is a common convention, while the second codifies a property which we've found facilitates the construction of MSTD sets.

- We let  $[a, b]$  denote all integers from  $a$  to  $b$ ; thus  $[a, b] = \{n \in \mathbb{Z} : a \leq n \leq b\}$ .
- We say a set of integers  $A$  has the property  $P_n$  (or is a  $P_n$ -set) if both its sumset and its difference set contain all but the first and last  $n$  possible elements (and of course it may or may not contain some of these fringe elements).<sup>4</sup> Explicitly, let  $a = \min A$  and  $b = \max A$ . Then  $A$  is a  $P_n$ -set if

$$[2a + n, 2b - n] \subset A + A \quad (2.16)$$

and

$$[-(b - a) + n, (b - a) - n] \subset A - A. \quad (2.17)$$

We can now state our construction and main result.

**Theorem 2.8** (Miller-Orosz-Scheinerman [MOS]). *Let  $A = L \cup R$  be a  $P_n$ , MSTD set where  $L \subset [1, n]$ ,  $R \subset [n + 1, 2n]$ , and  $1, 2n \in A$ ;<sup>5</sup> see Remark 2.9 for an example of such an  $A$ . Fix a  $k \geq n$  and let  $m$  be arbitrary. Let  $M$  be any subset of  $[n + k + 1, n + k + m]$  with the property that it does not have a run of more than  $k$  missing elements (i.e., for all  $\ell \in [n + k + 1, n + m + 1]$  there is a  $j \in [\ell, \ell + k - 1]$  such that  $j \in M$ ). Assume further that  $n + k + 1 \notin M$  and set  $A(M; k) = L \cup O_1 \cup M \cup O_2 \cup R'$ , where  $O_1 = [n + 1, n + k]$ ,  $O_2 = [n + k + m + 1, n + 2k + m]$  (thus the  $O_i$ 's are just sets of  $k$  consecutive integers), and  $R' = R + 2k + m$ . Then*

- (1)  $A(M; k)$  is an MSTD set, and thus we obtain an infinite family of distinct MSTD sets as  $M$  varies;
- (2) there is a constant  $C > 0$  such that as  $r \rightarrow \infty$  the percentage of subsets of  $\{1, \dots, r\}$  that are in this family (and thus are MSTD sets) is at least  $C/r^4$ .

**Remark 2.9.** *In order to show that our theorem is not trivial, we must of course exhibit at least one  $P_n$ , MSTD set  $A$  satisfying all our requirements (else our family is empty!).*

---

$\{m - d, 2m - d, \dots, km - d\}$ ,  $a^* = (k + 1)m - 2d$  and  $A = B \cup L \cup (a^* - B) \cup \{m\}$ . Then  $A$  is an MSTD set. The width of such a set is of the order  $km$ . Thus, if we look at all triples  $(m, d, k)$  with  $km \leq r$  satisfying the above conditions, these generate on the order of at most  $\sum_{k \leq r} \sum_{m \leq r/k} \sum_{d \leq m} 1 \ll r^2$ , and there are of the order  $2^r$  possible subsets of  $\{0, \dots, r\}$ ; thus this construction generates a negligible number of MSTD sets. Though we write  $f(r)/2^{r/2}$  to bound the percentage from other methods, a more careful analysis shows it is significantly less; we prefer this easier bound as it is already significantly less than our method. See for example Theorem 2 of [He] for a denser example.

<sup>4</sup>It is not hard to show that for fixed  $0 < \alpha \leq 1$  a random set drawn from  $[1, n]$  in the uniform model is a  $P_{\lfloor \alpha n \rfloor}$ -set with probability approaching 1 as  $n \rightarrow \infty$ .

<sup>5</sup>Requiring  $1, 2n \in A$  is quite mild; we do this so that we know the first and last elements of  $A$ .

We may take the set<sup>6</sup>  $A = \{1, 2, 3, 5, 8, 9, 13, 15, 16\}$ ; it is an MSTD set as

$$\begin{aligned} A + A &= \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, \\ &\quad 22, 23, 24, 25, 26, 28, 29, 30, 31, 32\} \\ A - A &= \{-15, -14, -13, -12, -11, -10, -8, -7, -6, -5, -4, -3, -2, -1, \\ &\quad 0, 1, 2, 3, 4, 5, 6, 7, 8, 10, 11, 12, 13, 14, 15\} \end{aligned} \quad (2.18)$$

(so  $|A + A| = 30 > 29 = |A - A|$ ).  $A$  is also a  $P_n$ -set, as (2.16) is satisfied since  $[10, 24] \subset A + A$  and (2.17) is satisfied since  $[-7, 7] \subset A - A$ .

For the uniform model, a subset of  $[1, 2n]$  is a  $P_n$ -set with high probability as  $n \rightarrow \infty$ , and thus examples of this nature are plentiful. For example, of the 1748 MSTD sets with minimum 1 and maximum 24, 1008 are  $P_n$ -sets.

**Project 2.10.** Read [MOS]. Can their argument be improved to yield a positive percentage through explicit construction?

Instead of sums and differences of two sets, we can consider a more general problem. Instead of searching for  $A$  such that  $|A + A| > |A - A|$ , we now consider the more general problem of when

$$|\epsilon_1 A + \cdots + \epsilon_n A| > |\tilde{\epsilon}_1 A + \cdots + \tilde{\epsilon}_n A|, \quad \epsilon_i, \tilde{\epsilon}_i \in \{-1, 1\}. \quad (2.19)$$

Consider the generalized sumset

$$f_{j_1, j_2}(A) = A + A + \cdots + A - A - A - \cdots - A, \quad (2.20)$$

where there are  $j_1$  pluses<sup>7</sup> and  $j_2$  minuses, and set  $j = j_1 + j_2$ . Our notion of a  $P_n$ -set generalizes, and we find that if there exists one set  $A$  with  $|f_{j_1, j_2}(A)| > |f_{j'_1, j'_2}(A)|$ , then we can construct infinitely many such  $A$ . Note without loss of generality that we may assume  $j_1 \geq j_2$ .<sup>8</sup>

**Definition 2.11** ( $P_n^j$ -set.). Let  $A \subset [1, k]$  with  $1, k \in A$ . We say  $A$  is a  $P_n^j$ -set if any  $f_{j_1, j_2}(A)$  contains all but the first  $n$  and last  $n$  possible elements.

**Remark 2.12.** Note that a  $P_n^2$ -set is the same as what we called a  $P_n$ -set earlier.

We expect the following generalization of Theorem 2.8 to hold.

**Conjecture 2.13.** For any  $f_{j_1, j_2}$  and  $f_{j'_1, j'_2}$ , if there exists a finite set of integers  $A$  which is (1) a  $P_n^j$ -set; (2)  $A \subset [1, 2n]$  and  $1, 2n \in A$ ; and (3)  $|f_{j_1, j_2}(A)| > |f_{j'_1, j'_2}(A)|$ , then there exists an infinite family of such sets.

The difficulty in proving the above conjecture is that we need to find a set  $A$  satisfying  $|f_{j_1, j_2}(A)| > |f_{j'_1, j'_2}(A)|$ ; once we find such a set, we can mirror the construction from Theorem 2.8. Currently we can only find such  $A$  for  $j \in \{2, 3\}$ :

<sup>6</sup>This  $A$  is trivially modified from [?] by adding 1 to each element, as we start our sets with 1 while other authors start with 0. We chose this set as our example as it has several additional nice properties that were needed in earlier versions of our construction which required us to assume slightly more about  $A$ .

<sup>7</sup>By a slight abuse of notation, we say there are two sums in  $A + A - A$ , as is clear when we write it as  $\epsilon_1 A + \epsilon_2 A + \epsilon_3 A$ .

<sup>8</sup>This follows as we are only interested in  $|f_{j_1, j_2}(A)|$ , which equals  $|f_{j_2, j_1}(A)|$ . This is because  $B$  and  $-B$  have the same cardinality, and thus (for example) we see  $A + A - A$  and  $-(A - A - A)$  have the same cardinality.

**Theorem 2.14.** *Conjecture 2.13 is true for  $j \in \{2, 3\}$ .*

Similar to the original result, it is crucial that we have a set to start the process. The following set was obtained by taking elements in  $\{2, \dots, 49\}$  to be in  $A$  with probability<sup>9</sup>  $1/3$  (and, of course, requiring  $1, 50 \in A$ ); it took about 300000 sets to find the first one satisfying our conditions:

$$A = \{1, 2, 5, 6, 16, 19, 22, 26, 32, 34, 35, 39, 43, 48, 49, 50\}. \quad (2.21)$$

To be a  $P_{25}^3$ -set we need to have  $A+A+A \supset [n+3, 6n-n] = [28, 125]$  and  $A+A-A \supset [-n+2, 3n-1] = [-23, 74]$ . A simple calculation shows  $A+A+A = [3, 150]$ , all possible elements, while  $A+A-A = [-48, 99] \setminus \{-34\}$  (i.e., every possible element but -34). Thus  $A$  is a  $P_{25}^3$ -set satisfying  $|A+A+A| > |A+A-A|$ , and thus we have the example we need to prove Theorem 2.14. We could also have taken

$$A = \{1, 2, 3, 4, 8, 12, 18, 22, 23, 25, 26, 29, 30, 31, 32, 34, 45, 46, 49, 50\}, \quad (2.22)$$

which has the same  $A+A+A$  and  $A+A-A$ .

**Project 2.15.** *Find a set  $A$  that will work for  $|A+A+A+A| > |A+A-A-A|$  or  $|A+A+A+A| > |A+A+A-A|$ .*

**Project 2.16.** *Generalize the above to  $|a_1A + a_2A + a_3A| > |b_1A + b_2A + b_3A|$ .*

**2.2. Structure of MSTD sets.** Frequently in mathematics we are interested in subsets of a larger collection where the subsets possess an additional property. In this sense, they are no longer generic subsets; however, we can ask what other properties they have or omit.

We observed earlier (Footnote 4) that for a constant  $0 < \alpha \leq 1$ , a set randomly chosen from  $[1, 2n]$  is a  $P_{\lfloor \alpha n \rfloor}$ -set with probability approaching 1 as  $n \rightarrow \infty$ . MSTD sets are of course not random, but it seems logical to suppose that this pattern continues.

**Project 2.17.** *Prove or disprove:*

**Conjecture 2.18.** *Fix a constant  $0 < \alpha \leq 1/2$ . Then as  $n \rightarrow \infty$  the probability that a randomly chosen MSTD set in  $[1, 2n]$  containing 1 and  $2n$  is a  $P_{\lfloor \alpha n \rfloor}$ -set goes to 1.*

In our construction and that of [MO], a collection of MSTD sets is formed by fixing the fringe elements and letting the middle vary. The intuition behind both is that the fringe elements matter most and the middle elements least. Motivated by this it is interesting to look at all MSTD sets in  $[1, n]$  and ask with what frequency a given element is in these sets. That is, what is

$$\gamma(k; n) = \frac{\#\{A : k \in A \text{ and } A \text{ is an MSTD set}\}}{\#\{A : A \text{ is an MSTD set}\}} \quad (2.23)$$

as  $n \rightarrow \infty$ ? We can get a sense of what these probabilities might be from Figure 3.

Note that, as the graph suggests,  $\gamma$  is symmetric about  $\frac{n+1}{2}$ , i.e.  $\gamma(k, n) = \gamma(n+1-k, n)$ . This follows from the fact that the cardinalities of the sumset and difference set are unaffected by sending  $x \rightarrow \alpha x + \beta$  for any  $\alpha, \beta$ . Thus for each MSTD set  $A$  we get

<sup>9</sup>Note the probability is  $1/3$  and not  $1/2$ .

Estimated  $\gamma(k,n)$

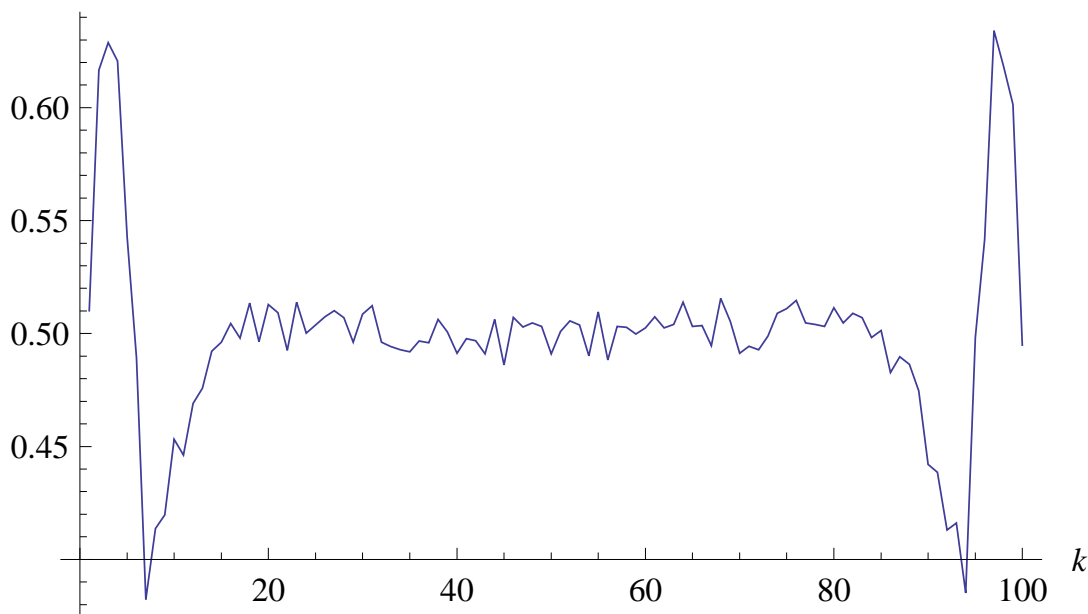


FIGURE 3. Estimation of  $\gamma(k, 100)$  as  $k$  varies from 1 to 100 from a random sample of 4458 MSTD sets.

a distinct MSTD set  $n + 1 - A$  showing that our function  $\gamma$  is symmetric. These sets are distinct since if  $A = n + 1 - A$  then  $A$  is sum-difference balanced.<sup>10</sup>

**Project 2.19.** *Make the following argument rigorous: From [MO] we know that a positive percentage of sets are MSTD sets. By the central limit theorem we then get that the average size of an MSTD set chosen from  $[1, n]$  is about  $n/2$ . This tells us that on average  $\gamma(k, n)$  is about  $1/2$ . The graph above suggests that the frequency goes to  $1/2$  in the center.*

The above leads us to the following conjecture:

**Project 2.20.**

**Conjecture 2.21.** *Fix a constant  $0 < \alpha < 1/2$ . Then  $\lim_{n \rightarrow \infty} \gamma(k, n) = 1/2$  for  $\lfloor \alpha n \rfloor \leq k \leq n - \lfloor \alpha n \rfloor$ . More generally, we could ask which non-decreasing functions  $f(n)$  have  $f(n) \rightarrow \infty$ ,  $n - f(n) \rightarrow \infty$  and  $\lim_{n \rightarrow \infty} \gamma(k, n) = 1/2$  for all  $k$  such that  $\lfloor f(n) \rfloor \leq k \leq n - \lfloor f(n) \rfloor$ .*

Note: Kevin O'Bryant may have some partial results along these lines; check with me before pursuing this.

**2.3. Catalan's conjecture and products of consecutive integers.** We can show that  $x(x + 1)(x + 2)(x + 3)$  is never a perfect square or cube for  $x$  a positive integer. One

<sup>10</sup>The following proof is standard (see, for instance, [Na2]). If  $A = n + 1 - A$  then

$$|A + A| = |A + (n + 1 - A)| = |n + 1 + (A - A)| = |A - A|. \quad (2.24)$$



proof involves using elliptic curves to handle some cases; without using elliptic curves, one can handle many cases by reducing to the Catalan equation, and in fact show it is never a perfect power.

Catalan's conjecture is that the only adjacent non-trivial perfect powers are 8 and 9 (we say  $n$  is a perfect power if  $n = m^a$  for some  $a \geq 2$ ). Catalan's theorem was proved in 2002. Explicitly

**Theorem 2.22** (Mihalescu 2002). *Let  $a, b \in \mathbb{Z}$  and  $n, m \geq 2$  positive integers. Consider the equation*

$$a^n - b^m = \pm 1. \quad (2.25)$$

*The only solutions are  $3^2 - 2^3 = 1$ ,  $2^3 - 3^2 = -1$ ,  $1^n - 0^m = 1$ , and  $0^n - 1^m = -1$ .*

Consider

$$x(x+1)(x+2)(x+3) = y^3. \quad (2.26)$$

We can re-group the factors and obtain

$$x(x+3) \cdot (x+1)(x+2) = (x^2 + 3x) \cdot (x^2 + 3x + 2) = y^3. \quad (2.27)$$

Letting  $z = x^2 + 3x + 1$ , we find that

$$(z-1)(z+1) = y^3. \quad (2.28)$$

We may re-write this as

$$z^2 - y^3 = 1. \quad (2.29)$$

The only solution is  $z = 3, y = 2$ , and this does not correspond to  $x$  a positive integer.

We now consider the obvious generalization to showing that  $x(x+1)(x+2)(x+3)$  is never a perfect power. The only change in the previous argument is that we now have  $y^m$  instead of  $y^3$  for some positive integer  $m \geq 2$ . We again obtain

$$z^2 - y^m = 1, \quad (2.30)$$

and again  $z = x^2 + 3x + 1 = 3$ , which has no solution. Note this also handles the case  $m = 2$  (ie,  $x(x+1)(x+2)(x+3)$  is never a square). This immediately gives

$$z^2 - 1 = y^2 \quad (2.31)$$

or equivalently

$$z^2 = y^2 + 1, \quad (2.32)$$

and there are no adjacent perfect squares other than 0 and 1; note  $z = 0$  yields a non-integral  $x$ .

**Project 2.23.** *Can this be generalized to products of more factors? What if we replace a perfect power by twice a perfect power?*

*Note: I have a lot of notes (joint with Cosmin Roman and Warren Sinnott) about elementary approaches that do not use Mihalescu's theorem, or only uses it in some cases. For example, let's consider the question of whether*

$$x(x+1)(x+2)(x+3) = y^2 \quad (2.33)$$

has any solutions in positive integers. (We find that it does not.) Let

$$u = 2x + 3, \quad z = u^2 \quad (2.34)$$

so that

$$\begin{aligned}
 (4y)^2 &= 2x(2x+2)(2x+4)(2x+6) \\
 &= (u-3)(u-1)(u+1)(u+3) \\
 &= (u^2-1)(u^2-9) \\
 &= (z-1)(z-9).
 \end{aligned} \tag{2.35}$$

The difference between  $z-1$  and  $z-9$  is 8, so the factors  $z-1$  and  $z-9$  have at most a power of 2 in common; since the left-hand side of the equations above is a square we may write

$$z-1 = 2^a v^2, \quad z-9 = 2^b w^2, \tag{2.36}$$

where  $a, b$  are either 0 or 1 and  $a+b$  is even, i.e., either  $a=b=0$  or  $a=b=1$ .

**Case One:**  $a=b=0$ . Here we have

$$z = 1 + v^2 = 9 + w^2, \tag{2.37}$$

so

$$8 = v^2 - w^2 = (v-w)(v+w). \tag{2.38}$$

So  $v-w$  and  $v+w$  are divisors of 8, the second larger than the first; also  $v-w$  and  $v+w$  must have the same parity. The only possibility is then

$$v-w = 2, \quad v+w = 4, \tag{2.39}$$

which implies that  $v=3$ , and  $z=10$ . But  $z=u^2$  is a square, so there are no solutions in this case.

**Case Two:**  $a=b=1$ . Here we have

$$z = 1 + 2v^2 = 9 + 2w^2, \tag{2.40}$$

so

$$4 = v^2 - w^2 = (v-w)(v+w). \tag{2.41}$$

So  $v-w$  and  $v+w$  are divisors of 4, the second larger than the first, and both of the same parity; so there are no solutions in this case either.

To see how elliptic curves can arise in questions such as this, consider now

$$x(x+1)(x+2)(x+3) = y^3. \tag{2.42}$$

Letting  $u = x-1$  we may re-write the above as

$$(u-1)u(u+1)(u+2) = y^3. \tag{2.43}$$

The only divisors any of the four factors can have in common are 2 and 3.

Assume that 3 divides at most one of the factors. Thus, 3 divides either  $u$  or  $u+1$ . Split the multiplication into two parts,  $(u-1)(u+1)$  and  $u(u+2)$ . All the factors of 2 occur in either the first multiplication or the second, but not both. As we are assuming 3 divides  $u$  or  $u+1$ , this implies that each of the two multiplications must be a perfect cube. In particular, we have

$$(u-1)(u+1) = w^3. \tag{2.44}$$

This simplifies to

$$u^2 - w^3 = 1. \quad (2.45)$$

This is the Catalan Equation, which is now known to have just one solution, namely  $u = 3$  and  $w = 2$ . Substituting in for  $u$  gives

$$(3 - 1)(3)(3 + 1)(3 + 2) = 120 = 2^3 \cdot 3 \cdot 5, \quad (2.46)$$

which is not a perfect square.

We are left with the case when  $3|u$  and  $3|(u + 2)$ . Clearly  $2|u(u + 1)$ . If, however, 4 does not divide  $u(u + 1)$ , then we must have

$$u(u + 1) = 2w^3, \quad (u - 1)(u + 2) = 2^2v^3. \quad (2.47)$$

Multiplying the first equation by 4 gives

$$(2u)(2u + 2) = (2w)^3. \quad (2.48)$$

Let  $z = 2u + 1$ . Then the above equation becomes

$$(z - 1)(z + 1) = (2w)^3, \quad (2.49)$$

which may be re-written as

$$z^2 - (2w)^3 = 1. \quad (2.50)$$

We again obtain the Catalan equation, which now has the unique solution  $z = 3, w = 1$ . If  $z = 3$  then  $u = 1$ , and  $(u - 1)u(u + 1)(u + 2) = 0$ , implying there are no solutions.

Thus, we are left with the case when  $3|u$ ,  $3|(u + 2)$ , and  $4|u(u + 1)$ . We could use elliptic curve arguments again. If  $(u - 1)(u + 1) \equiv 9 \pmod{27}$ , we would have

$$(u - 1)(u + 1) = 9w^3. \quad (2.51)$$

This leads to the elliptic curve

$$u^2 = 9w^3 + 1. \quad (2.52)$$

Letting  $u_2 = \frac{u}{2}$  and  $w_2 = \frac{w}{2}$  we obtain the elliptic curve

$$E : u_2^2 = w_2^3 + 81. \quad (2.53)$$

As  $L(E, 1) \approx 2.02$ , this curve has rank 0, and the only rational solutions are the torsion points. Direct calculation gives the torsion group is  $\mathbb{Z}/6\mathbb{Z}$ , generated by  $[0, 9]$ . Further computation should yield none of these give valid solutions to the original equation. Unfortunately, if  $(u - 1)(u + 1) \equiv 3 \pmod{27}$ , we obtain a rank 2 elliptic curve, which is a little harder to analyze. Fortunately, if this is the case than instead of looking at  $(u - 1)(u + 1)$ , we can look at  $u(u + 2)$ , which is equivalent to  $9 \pmod{27}$ . Letting  $z = u - 1$ , this gives us

$$(z - 1)(z + 1) = 9v^3, \quad (2.54)$$

and this is the same equation as before. It will also have zero rank, and torsion group  $\mathbb{Z}/6\mathbb{Z}$  generated by  $[0, 9]$ . Direct calculation will finish the proof.

**Project 2.24.** *See how far arguments like this can be pushed for this and related problems. Note: if you decide to work on this, email me and I'll send you my work in progress with Roman and Sinnott.*

**2.4. The  $3x + 1$  Problem.** Let  $x$  be a positive odd integer. Then  $3x + 1$  is even, and we can find a unique  $k > 0$  such that  $(3x + 1)/2^k$  is an odd number not divisible by 3. We denote this map by  $T$ , which is defined on  $\Pi = \{\ell > 0 : \ell \equiv_6 1 \text{ or } 5\}$  (the set of positive integers not divisible by 2 or 3). The famous  $3x + 1$  Conjecture states that for any  $x \in \Pi$  there is an  $n$  such that  $T^n(x) = 1$  (where  $T^2(x) = T(T(x))$  and so forth). As of February 1<sup>st</sup>, 2007, the conjecture has been numerically verified up to  $13 \cdot 2^{58} \approx 3.7 \cdot 10^{18}$ ; see [?, ?] for details.

People working on the Syracuse-Kakutani-Hasse-Ulam-Hailstorm-Collatz- $(3x + 1)$ -Problem (there have been a few) often refer to two striking anecdotes. One is Erdős' comment that "Mathematics is not yet ready for such problems." The other is Kakutani's communication to Lagarias: "For about a month everybody at Yale worked on it, with no result. A similar phenomenon happened when I mentioned it at the University of Chicago. A joke was made that this problem was part of a conspiracy to slow down mathematical research in the U.S." Coxeter has offered \$50 for its solution, Erdős \$500, and Thwaites, £1000. The problem has been connected to holomorphic solutions to functional equations, a Fatou set having no wandering domain, Diophantine approximation of  $\log_2 3$ , the distribution mod 1 of  $\left\{\left(\frac{3}{2}\right)^k\right\}_{k=1}^{\infty}$ , ergodic theory on  $\mathbb{Z}_2$ , undecidable algorithms, and geometric Brownian motion, to name a few (see [Lag1, Lag2]).

The following definition is a useful starting point for investigations of elements of  $\Pi = \{\ell > 0 : \ell \equiv_6 1 \text{ or } 5\}$  (the set of positive integers not divisible by 2 or 3) under the  $3x + 1$  map.

**Definition 2.25** (*m*-path). *The m-path of an  $x \in \Pi$  is the m-tuple of positive integers  $(k_1, \dots, k_m)$  such that*

$$T^i(x) = \frac{3T^{i-1}(x) + 1}{2^{k_i}}, \quad i \in \{1, \dots, m\}. \quad (2.55)$$

We often write  $\gamma_m(x)$  for the *m*-path of  $x$ .

For example, the first few iterates of 41 are 31, 47, 71, and 107. Thus 41 has a 1-path of (2), a 2-path of (2, 1), a 3-path of (2, 1, 1) and a 4-path of (2, 1, 1, 1). Similarly the 4-path of 11 (which iterates to 17, 13, 5 and then 1) is (1, 2, 3, 4). The *m*-paths are useful in studying the  $3x + 1$  problem. For example, if the sum of the elements in the *m*-path of  $x$  is "close" to  $m$  then the  $m^{\text{th}}$  iterate of  $x$  is "large" relative to  $x$  (as we see in our example with  $x = 41$ ); if the sum of the elements in the *m*-path of  $x$  is "large" relative to  $m$  then the  $m^{\text{th}}$  iterate of  $x$  is "small" relative to  $x$  (as we see in our example with  $x = 11$ ; in fact, all further iterates are 1, so 11 has an *m*-path of (1, 2, 3, 4, 2, 2,  $\dots$ , 2), where there are  $m - 4$  twos at the end). Crucial in our investigations is the Structure Theorem of Sinai and Kontorovich-Sinai [Si, KonSi].

**Theorem 2.26** (The Structure Theorem). *Let  $k_1, \dots, k_m$  be given positive integers, and  $\varepsilon \in \{1, 5\}$ . Then there exists a  $q_m \in [0, 6 \cdot 2^{k_1 + \dots + k_m})$  with  $q_m \equiv_6 \varepsilon$  such that*

$$\{x \in \Pi : \gamma_m(x) = (k_1, \dots, k_m)\} = \{6 \cdot 2^{k_1 + \dots + k_m} p + q_m : p \in \mathbb{N}\}. \quad (2.56)$$

Hence, for given  $k_1, \dots, k_m$ , we have two *full* arithmetic progressions, one for  $\varepsilon = 1$  and one for  $\varepsilon = 5$ . Further, we need only find the minimal representatives in order to completely determine the solutions. Two of my former students, Bruce Adcock and

Sucheta Soundarajan, constructing a nice algorithm to determine the minimal representative of an arbitrary  $m$ -path. We investigated the properties of paths associated to elements of  $\Pi$  which do not iterate to 1 (i.e., elements which eventually iterate into a cycle or diverge to infinity). Our main result is that

if an element  $x > 1$  of  $\Pi$  is the minimal element of a cycle of length  $m$ , then  $m > 6,586,818,669$ .

**(Check and see if our number has been improved, and give statements on what we could improve it to if we improve regions where we know  $3x + 1$  holds.)** The two main ingredients of the proof are (1) an analysis of paths of iterates which always remain above the starting seed, and (2) knowing the  $3x + 1$  conjecture is true for all  $x \in \Pi$  with  $x \leq B$ . Our results extend those of other researches (see [Br, Sim, SimWe, Sin] and the references therein), though one must be careful in comparing the strength of bounds from paper to paper, as often different variants of the  $3x + 1$  problem are used<sup>11</sup>.

**Project 2.27.** *Take the preprint that I have and clean it up. This was a really nice project which the authors never finished writing up, and which I've been saving for a student. There are a few small mistakes throughout the paper, lots of places where the explanations are unclear. I have made numerous comments throughout the paper to help whomever looks at this complete the project. While it will take some time to clean everything up, there are some nice ideas here, and it is definitely a significant contribution to join the team and get this paper to publication.*

### 3. DIFFERENTIAL EQUATIONS

In [KP] the following equation is shown to be related to the propagation of infections:

$$f_n \left( \begin{pmatrix} x \\ y \end{pmatrix} \right) = \begin{pmatrix} 1 - (1 - ax)(1 - by)^n \\ 1 - (1 - ay)(1 - bx) \end{pmatrix} \quad (3.1)$$

(where we have replaced  $d$  with  $1 - a$ ). We study  $f_n : [0, 1]^2 \rightarrow [0, 1]^2$ .

The model is as follows. We have a central hub and  $n$  satellite vertices forming a graph. There are only edges from each satellite to the central hub; thus the satellite vertices communicate with each other only through the hub. The goal is to understand how viruses propagate in such a network; this is not an unreasonable model for certain situations (such as airlines). We have a complete solution when  $n = 1$  (which is fairly trivial) and a conjecture as to what is happening for general  $n$ , namely that the critical threshold is comparing  $b$  to  $(1 - a)/\sqrt{n}$ . If  $b < (1 - a)/\sqrt{n}$  then the behavior is trivial, and all initial configurations collapse to the trivial fixed point; we conjecture that if  $b > (1 - a)/\sqrt{n}$  all iterates end up at a unique non-trivial fixed point (so long as we don't start off at the trivial fixed point of the system).

I have a large draft of a paper on this with several colleagues (Leo Kontorovich and Amitabha Roy); the paper is available on the webpage (it is poorly written, more as a free association of results as we attempt to understand the system, so perhaps it's worth

---

<sup>11</sup>We pull out all powers of 2 in the same step as multiplying by 3 and adding 1. Some authors use instead the map  $T_1(x) = 3x + 1$  for  $x$  odd and  $x/2$  for  $x$  even, while others use  $T_2(x) = (3x + 1)/2$  for  $x$  odd and  $x/2$  for  $x$  even.

reading as an insight into how we chip away at a problem). We have lots of numerics supporting our conjecture. It might be possible to prove our results with sufficiently delicate coding and error analysis. To prove our claims in many regions involves nice applications of linear algebra and multivariable calculus, and an introduction to fixed point theorems.

#### 4. PROBABILITY

- (1) Products of Poisson Random Variables: elementary number theory and probability theory (some Fourier analysis is useful in understanding the applications).
- (2) Sabermetrics: elementary probability theory, though statistics would help if you are interested in numerical investigations / comparisons.
- (3) Die battles: elementary probability and combinatorics.
- (4) Beyond the pidgeonhole principle: elementary probability and combinatorics.
- (5) Differentiating identities: elementary probability and combinatorics.

**4.1. Products of Poisson Random Variables.** We live in an age where we are constantly bombarded with massive amounts of data. Satellites orbiting the earth daily transmit more information than is in the entire Library of Congress; researchers must quickly sort through these data sets to find the relevant pieces. It is thus not surprising that people are interested in patterns in data. One of the more interesting, and initially surprising, is Benford's law.

At some point in secondary school, we are introduced to scientific notation: any positive number  $x$  may be written as  $M(x) \cdot 10^k$ , where  $M(x) \in [1, 10)$  is the mantissa and  $k$  is an integer. Thus 1701.24601 would be written as  $1.70124601 \cdot 10^3$  and .00729735257 would be  $7.29735257 \cdot 10^{-3}$ ; the first has a leading digit (or first digit) of 1 while the second has a leading digit of 7.

**Definition 4.1** (Benford's Law). *Benford's law states that for many natural sets of data, the probability of observing a first digit of  $d$  is  $\log_{10} \left( \frac{d+1}{d} \right)$ .*

Other useful definitions:

**Definition 4.2** (Modular (or clock) arithmetic). *We say  $a \equiv b \pmod{n}$  if  $a-b$  is a multiple of  $n$ . This is frequently called clock arithmetic, as this is the most common example; on a clock, 13 o'clock and 1 o'clock are both represented by 1.*

**Definition 4.3** (Equidistributed modulo 1). *A sequence  $\{z_n\}_{n=-\infty}^{\infty}$  is equidistributed modulo 1 if*

$$\lim_{N \rightarrow \infty} \frac{\#\{n : |n| \leq N, z_n \bmod 1 \in [a, b]\}}{2N + 1} = b - a \quad (4.1)$$

*for all  $[a, b] \subset [0, 1]$ . A similar definition holds for  $\{z_n\}_{n=0}^{\infty}$ .*

We may generalize Benford's law in many ways. The two most common are:

- (1) Instead of studying the distribution of the first digit, we may study the distribution of the first two, three, or more generally the mantissa of our number. Benford's law becomes the probability of observing a mantissa of at most  $s$  is  $\log_{10} s$ .

- (2) Instead of working base 10, we may work base  $B$ , in which case the Benford probabilities become  $\log_B \left( \frac{d+1}{d} \right)$  for the distribution of the first digit, and  $\log_B s$  for a mantissa of at most  $s$ .

It has been shown (see for example [JKKKM, ?]) that products of random variables converge to Benford's law. The mathematics behind this in full generality typically uses Fourier or Mellin transforms and lead to terrific estimates on the rate of convergence, but the rough idea is not hard to explain. Benford's law is really equivalent to the statement that  $\{x_n\}$  is Benford if and only if  $\{\log_{10} x_n \bmod 1\}$  is equidistributed in  $[0, 1]$ . The proof follows from the following two lemmas:

**Lemma 4.4.** *The first digits of  $10^u$  and  $10^v$  are the same in base  $b$  if and only if  $u \equiv v \pmod 1$ .*

Consider the unit interval  $[0, 1)$ . For  $d \in \{1, \dots, 9\}$ , define  $p_d$  by

$$10^{p_d} = d \quad \text{or equivalently} \quad p_d = \log_{10} d. \quad (4.2)$$

For  $d \in \{1, \dots, 9\}$ , let

$$I_d = [p_d, p_{d+1}) \subset [0, 1). \quad (4.3)$$

**Lemma 4.5.** *The first digit of  $10^y$  is  $d$  if and only if  $y \bmod 1 \in I_d$ .*

Why does this imply products converge to Benford's law? Let  $X_1, \dots, X_n$  be independent, identically distributed random variables with mean  $\mu$ , variance  $\sigma^2$  and finite higher moments (the result holds under far weaker conditions). Then  $X_1 + \dots + X_n$  converges to being normally distributed with mean  $n\mu$  and variance  $n\sigma^2$ . We, however, want to study the product  $Y_n = X_1 \cdots X_n$ . Whenever we see a product our first thought should be to take logarithms. Thus  $\log_{10} Y_n = \log_{10} X_1 + \dots + \log_{10} X_n$ . If we let  $\tilde{\mu}$  and  $\tilde{\sigma}^2$  be the mean and variance of  $\log_{10} X_i$ , we see that  $Y_n$  converges to a Gaussian with mean  $n\tilde{\mu}$  and variance  $n\tilde{\sigma}^2$ ; however, to obtain Benford behavior we just need to understand the distribution of  $\log_{10} Y_n$  modulo 1. It is not hard to show that as the variance of a normal distribution tends to infinity, modulo 1 the probability converges to the uniform distribution, and this is where we obtain the Benford behavior. The proof is a nice application of Fourier analysis (in particular, Poisson summation), though it could probably be proved by a careful use of Taylor's theorem with remainder.

For reasons I don't want to get into in a public post, it is of interest to study products of Poisson random variables. Recall  $X$  is said to be a Poisson random variable with parameter  $\lambda$  if the probability  $X$  equals  $n$  is  $\lambda^n e^{-\lambda} / n!$ .

**Project 4.6.** *Derive a closed form expression for the probability density of  $X_1 \cdots X_n$ , where each  $X_i$  is a Poisson random variable with parameter  $\lambda$ . Obtain as tractable expressions as possible. This will require some number theory. For example, say  $n = 2$ . The probability that  $X_1 X_2 = 42$  is much larger than the probability the product is either 41 or 43, as the latter two are primes and 42 is composite. For our purposes, it would suffice to have a good formula for the probability the product is in  $d \cdot 10^k$  to  $(d+1) \cdot 10^k$  for any  $d \in \{1, \dots, 9\}$  and  $k$  a non-negative integer.*

4.2. **Sabermetrics.** There are numerous fun problems in sabermetrics (applying math/stats to baseball). Here are two of my favorites.

4.2.1. *The Pythagorean Won-Loss Theorem.* It has been noted that in many professional sports leagues a good predictor of a team's end of season won-loss percentage is Bill James' Pythagorean Formula  $\frac{RS_{\text{obs}}^\gamma}{RS_{\text{obs}}^\gamma + RA_{\text{obs}}^\gamma}$ , where  $RS_{\text{obs}}$  (resp.  $RA_{\text{obs}}$ ) is the observed average number of runs scored (allowed) per game and  $\gamma$  is a constant for the league; for baseball the best agreement is when  $\gamma$  is about 1.82. This formula is often used in the middle of a season to determine if a team is performing above or below expectations, and estimate their future standings.

I provided a theoretical justification for this formula and value of  $\gamma$  by modeling the number of runs scored and allowed in baseball games as independent random variables drawn from Weibull distributions with the same  $\beta$  and  $\gamma$  but different  $\alpha$ ; the probability density is

$$f(x; \alpha, \beta, \gamma) = \begin{cases} \frac{\gamma}{\alpha} ((x - \beta)/\alpha)^{\gamma-1} e^{-((x-\beta)/\alpha)^\gamma} & \text{if } x \geq \beta \\ 0 & \text{otherwise.} \end{cases}$$

This model leads to a predicted won-loss percentage of  $\frac{(RS-\beta)^\gamma}{(RS-\beta)^\gamma + (RA-\beta)^\gamma}$ ; here  $RS$  (resp.  $RA$ ) is the mean of the Weibull random variable corresponding to runs scored (allowed), and  $RS - \beta$  (resp.  $RA - \beta$ ) is an estimator of  $RS_{\text{obs}}$  (resp.  $RA_{\text{obs}}$ ). An analysis of the 14 American League teams from the 2004 baseball season shows that (1) given that the runs scored and allowed in a game cannot be equal, the runs scored and allowed are statistically independent; (2) the best fit Weibull parameters attained from a least squares analysis and the method of maximum likelihood give good fits. Specifically, least squares yields a mean value of  $\gamma$  of 1.79 (with a standard deviation of .09) and maximum likelihood yields a mean value of  $\gamma$  of 1.74 (with a standard deviation of .06), which agree beautifully with the observed best value of 1.82 attained by fitting  $\frac{RS_{\text{obs}}^\gamma}{RS_{\text{obs}}^\gamma + RA_{\text{obs}}^\gamma}$  to the observed winning percentages.

The main calculation is as follows. We determine the mean of a Weibull distribution with parameters  $(\alpha, \beta, \gamma)$ , and then use this to prove our main result, the Pythagorean Formula. Let  $f(x; \alpha, \beta, \gamma)$  be the probability density of a Weibull with parameters  $(\alpha, \beta, \gamma)$ :

$$f(x; \alpha, \beta, \gamma) = \begin{cases} \frac{\gamma}{\alpha} \left(\frac{x-\beta}{\alpha}\right)^{\gamma-1} e^{-((x-\beta)/\alpha)^\gamma} & \text{if } x \geq \beta \\ 0 & \text{otherwise.} \end{cases} \quad (4.4)$$

For  $s \in \mathbb{C}$  with the real part of  $s$  greater than 0, recall the  $\Gamma$ -function (see [?]) is defined by

$$\Gamma(s) = \int_0^\infty e^{-u} u^{s-1} du = \int_0^\infty e^{-u} u^s \frac{du}{u}. \quad (4.5)$$

Letting  $\mu_{\alpha, \beta, \gamma}$  denote the mean of  $f(x; \alpha, \beta, \gamma)$ , we have

$$\begin{aligned} \mu_{\alpha, \beta, \gamma} &= \int_\beta^\infty x \cdot \frac{\gamma}{\alpha} \left(\frac{x-\beta}{\alpha}\right)^{\gamma-1} e^{-((x-\beta)/\alpha)^\gamma} dx \\ &= \int_\beta^\infty \alpha \frac{x-\beta}{\alpha} \cdot \frac{\gamma}{\alpha} \left(\frac{x-\beta}{\alpha}\right)^{\gamma-1} e^{-((x-\beta)/\alpha)^\gamma} dx + \beta. \end{aligned} \quad (4.6)$$



We change variables by setting  $u = \left(\frac{x-\beta}{\alpha}\right)^\gamma$ . Then  $du = \frac{\gamma}{\alpha} \left(\frac{x-\beta}{\alpha}\right)^{\gamma-1} dx$  and we have

$$\begin{aligned}\mu_{\alpha,\beta,\gamma} &= \int_0^\infty \alpha u^{\gamma-1} \cdot e^{-u} du + \beta \\ &= \alpha \int_0^\infty e^{-u} u^{1+\gamma-1} \frac{du}{u} + \beta \\ &= \alpha \Gamma(1 + \gamma^{-1}) + \beta.\end{aligned}\quad (4.7)$$

A similar calculation determines the variance. We record these results:

**Lemma 4.7.** *The mean  $\mu_{\alpha,\beta,\gamma}$  and variance  $\sigma_{\alpha,\beta,\gamma}^2$  of a Weibull with parameters  $(\alpha, \beta, \gamma)$  are*

$$\begin{aligned}\mu_{\alpha,\beta,\gamma} &= \alpha \Gamma(1 + \gamma^{-1}) + \beta \\ \sigma_{\alpha,\beta,\gamma}^2 &= \alpha^2 \Gamma(1 + 2\gamma^{-1}) - \alpha^2 \Gamma(1 + \gamma^{-1})^2.\end{aligned}\quad (4.8)$$

We can now prove our main result:

Let  $X$  and  $Y$  be independent random variables with Weibull distributions  $(\alpha_{RS}, \beta, \gamma)$  and  $(\alpha_{RA}, \beta, \gamma)$  respectively, where  $X$  is the number of runs scored and  $Y$  the number of runs allowed per game. As the means are RS and RA, by Lemma 4.7 we have

$$\begin{aligned}\text{RS} &= \alpha_{RS} \Gamma(1 + \gamma^{-1}) + \beta \\ \text{RA} &= \alpha_{RA} \Gamma(1 + \gamma^{-1}) + \beta.\end{aligned}\quad (4.9)$$

Equivalently, we have

$$\begin{aligned}\alpha_{RS} &= \frac{\text{RS} - \beta}{\Gamma(1 + \gamma^{-1})} \\ \alpha_{RA} &= \frac{\text{RA} - \beta}{\Gamma(1 + \gamma^{-1})}.\end{aligned}\quad (4.10)$$

We need only calculate the probability that  $X$  exceeds  $Y$ . Below we constantly use the integral of a probability density is 1. We have

$$\begin{aligned}\text{Prob}(X > Y) &= \int_{x=\beta}^\infty \int_{y=\beta}^x f(x; \alpha_{RS}, \beta, \gamma) f(y; \alpha_{RA}, \beta, \gamma) dy dx \\ &= \int_{x=\beta}^\infty \int_{y=\beta}^x \frac{\gamma}{\alpha_{RS}} \left(\frac{x-\beta}{\alpha_{RS}}\right)^{\gamma-1} e^{-((x-\beta)/\alpha_{RS})^\gamma} \frac{\gamma}{\alpha_{RA}} \left(\frac{y-\beta}{\alpha_{RA}}\right)^{\gamma-1} e^{-((y-\beta)/\alpha_{RA})^\gamma} dy dx \\ &= \int_{x=0}^\infty \frac{\gamma}{\alpha_{RS}} \left(\frac{x}{\alpha_{RS}}\right)^{\gamma-1} e^{-(x/\alpha_{RS})^\gamma} \left[ \int_{y=0}^x \frac{\gamma}{\alpha_{RA}} \left(\frac{y}{\alpha_{RA}}\right)^{\gamma-1} e^{-(y/\alpha_{RA})^\gamma} dy \right] dx \\ &= \int_{x=0}^\infty \frac{\gamma}{\alpha_{RS}} \left(\frac{x}{\alpha_{RS}}\right)^{\gamma-1} e^{-(x/\alpha_{RS})^\gamma} [1 - e^{-(x/\alpha_{RA})^\gamma}] dx \\ &= 1 - \int_{x=0}^\infty \frac{\gamma}{\alpha_{RS}} \left(\frac{x}{\alpha_{RS}}\right)^{\gamma-1} e^{-(x/\alpha)^\gamma} dx,\end{aligned}\quad (4.11)$$

where we have set

$$\frac{1}{\alpha^\gamma} = \frac{1}{\alpha_{RS}^\gamma} + \frac{1}{\alpha_{RA}^\gamma} = \frac{\alpha_{RS}^\gamma + \alpha_{RA}^\gamma}{\alpha_{RS}^\gamma \alpha_{RA}^\gamma}.\quad (4.12)$$

Therefore

$$\begin{aligned}
 \text{Prob}(X > Y) &= 1 - \frac{\alpha^\gamma}{\alpha_{\text{RS}}^\gamma} \int_0^\infty \frac{\gamma}{\alpha} \left(\frac{x}{\alpha}\right)^{\gamma-1} e^{-(x/\alpha)^\gamma} dx \\
 &= 1 - \frac{\alpha^\gamma}{\alpha_{\text{RS}}^\gamma} \\
 &= 1 - \frac{1}{\alpha_{\text{RS}}^\gamma} \frac{\alpha_{\text{RS}}^\gamma \alpha_{\text{RA}}^\gamma}{\alpha_{\text{RS}}^\gamma + \alpha_{\text{RA}}^\gamma} \\
 &= \frac{\alpha_{\text{RS}}^\gamma}{\alpha_{\text{RS}}^\gamma + \alpha_{\text{RA}}^\gamma}. \tag{4.13}
 \end{aligned}$$

Substituting the relations for  $\alpha_{\text{RS}}$  and  $\alpha_{\text{RA}}$  of (4.10) into (4.13) yields

$$\text{Prob}(X > Y) = \frac{(\text{RS} - \beta)^\gamma}{(\text{RS} - \beta)^\gamma + (\text{RA} - \beta)^\gamma}, \tag{4.14}$$

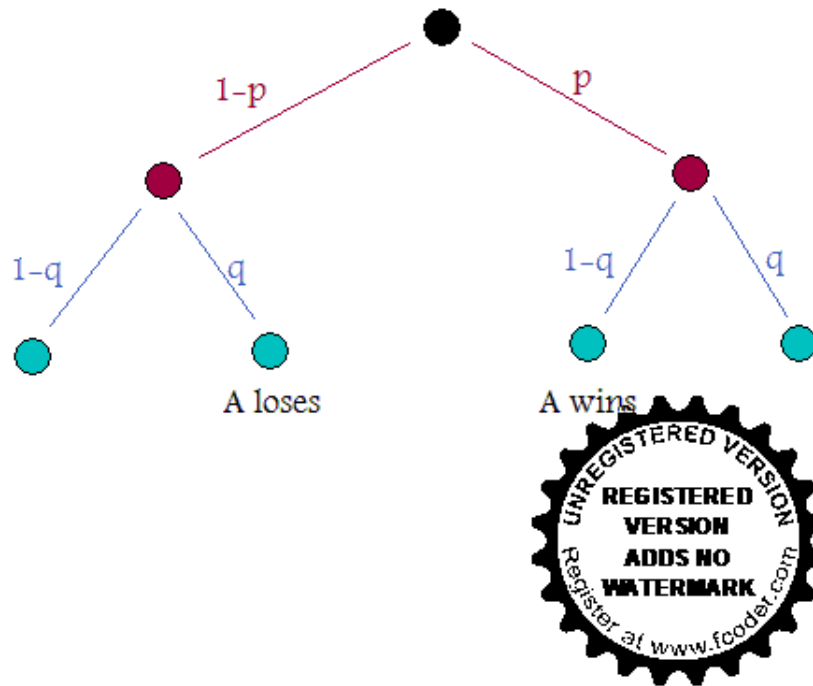
**Project 4.8.** *Obviously, I don't feel that baseball players sit down and talk about how to score and allow runs Weibullishly; I chose the three parameter Weibull distribution as it is quite flexible and fits a variety of 'one-hump' distributions and all the needed integrals can be done in closed form. This means we get a nice formula for the winning percentage in terms of the parameters of the teams, and thus we can quickly predict how much a team would improve by working on various parts. This is why explicit formulas are so useful; it is trivial to do lots of numerical simulations, but difficult in general to obtain a closed form. Can you find other distributions that lead to closed form expressions? (The generalized Gamma should work for some values of its parameters.)*

4.2.2. *The log 5 Rule.* Let  $p$  and  $q$  denote the winning percentages of teams  $A$  and  $B$ . The following formula has numerically been observed to provide a terrific estimate of the probability that  $A$  beats  $B$ :  $(p - pq)/(p + q - 2pq)$ .

When we say  $A$  has a winning percentage of  $p$ , we mean that if  $A$  were to play an average team many times, then  $A$  would win about  $p\%$  of the games (for us, an average team is one whose winning percentage is .500). Let us image a third team, say  $C$ , with a .500 winning percentage. We image  $A$  and  $C$  playing as follows. We randomly choose either 0 or 1 for each team; if one team has a higher number then they win, and if both numbers are the same then we choose again (and continue indefinitely until one team has a higher number than the other). For  $A$  we choose 1 with probability  $p$  and 0 with probability  $1 - p$ , while for  $C$  we choose 1 and 0 with probability  $1/2$ . It is easy to see that this method yields  $A$  beating  $C$  exactly  $p\%$  of the time.

The probability that  $A$  wins the first time we choose numbers is  $p \cdot 1/2$  (the only way  $A$  wins is if we choose 1 for  $A$  and 0 for  $C$ , and the probability this happens is just  $p \cdot 1/2$ ). If  $A$  were to win on the second iteration then we must have either chosen two 1's initially (which happens with probability  $p \cdot 1/2$ ) or two 0's initially (which happens with probability  $(1 - p) \cdot 1/2$ ), and then we must choose 1 for  $A$  and 0 for  $B$  (which happens with probability  $p \cdot 1/2$ ). Continuing this process, we see that the probability  $A$  wins on the  $n^{\text{th}}$  iteration is

$$\left(p \cdot \frac{1}{2} + (1 - p) \cdot \frac{1}{2}\right)^{n-1} \cdot \left(p \cdot \frac{1}{2}\right) = \frac{p}{2^n}. \tag{4.15}$$

FIGURE 4. Probability tree for  $A$  beats  $B$  in one iteration.

Summing these probabilities gives a geometric series:

$$\sum_{n=1}^{\infty} \frac{p}{2^n} = p, \quad (4.16)$$

proving the claim.

Imagine now that  $A$  and  $B$  are playing. We choose 1 for  $A$  with probability  $p$  and 0 with probability  $1-p$ , while for  $B$  we choose 1 with probability  $q$  and 0 with probability  $1-q$ . If in any iteration one of the teams has a higher number than the other, we declare that team the winner; if not, we randomly choose numbers for the teams until one has a higher number.

The probability  $A$  wins on the first iteration is  $p \cdot (1-q)$  (the probability that  $A$  is 1 and  $B$  is 0). The probability that  $A$  neither wins or loses on the first iteration is  $(1-p)(1-q) + pq = 1-p-q+2pq$  (the first factor is the probability we chose 0 twice, while the second is the probability we chose 1 twice). Thus the probability  $A$  wins on the second iteration is just  $(1-p-q+2pq) \cdot p(1-q)$ ; see Figure 4.

Continuing this argument, the probability  $A$  wins on the  $n^{\text{th}}$  iteration is just

$$(1-p-q+2pq)^{n-1} \cdot p(1-q). \quad (4.17)$$

Summing<sup>12</sup> we find the probability  $A$  wins is just

$$\begin{aligned} \sum_{n=1}^{\infty} (1-p-q+2pq)^{n-1} \cdot p(1-q) &= p(1-q) \sum_{n=0}^{\infty} (1-p-q+2pq)^n \\ &= \frac{p(1-q)}{1-(1-p-q+2pq)} \\ &= \frac{p(1-q)}{p+q-2pq}. \end{aligned} \quad (4.18)$$

It is illuminating to write the denominator as  $p(1-q)+q(1-p)$ , and thus the formula becomes

$$\frac{p(1-q)}{p(1-q)+q(1-p)}. \quad (4.19)$$

This variant makes the extreme cases more apparent. Further, there are only two ways the process can terminate after one iteration:  $A$  wins (which happens with probability  $p(1-q)$ ) or  $B$  wins (which happens with probability  $(1-p)q$ ). Thus this formula is the probability that  $A$  won given that the game was decided in just one iteration.

**Project 4.9.** *Can you find other simple, elegant formulas to predict the probability one team beats another? The more information one uses, the more accurate the formula should be but the harder it will be to apply.*

**4.3. Die battles.** Two players roll die with  $k$  sides, with each side equally likely of being rolled. Player one rolls  $m$  dice and player two rolls  $n$  dice. If player one's highest roll exceeds the highest roll of player two then player one wins, otherwise player two wins. We can calculate the probability that player one wins, giving a concise summation and integral version, as well as estimating the probability that player one wins for many triples  $(m, n, k)$ . The answer involves numerous useful techniques (adding zero, multiplying by one, telescoping series), as well as some beautiful formulas (formulas for sums of powers, the binomial theorem, order statistics, partial summation).

**Project 4.10.** *Read the paper on the webpage. This is the first of many problems one can ask about player one and player two. Other natural questions are:*

- (1) *For fixed  $k$ , what is the probability that player one wins as  $m$  and  $n$  tend to infinity? Does it matter how they tend to infinity? For example, is the answer different if  $m = n$  or  $m = n^2$ ?*
- (2) *What is the probability that player one's top two rolls exceed the top two rolls of player two? Or, more generally, compare the largest  $c$  rolls of player one and two. Such a calculation is useful in the board game RISK, where often the attacker uses three die and the defender two die.*

---

<sup>12</sup>To use the geometric series formula, we need to know that the ratio is less than 1 in absolute value. Note  $1-p-q+2pq = 1-p(1-q)-q(1-p)$ . This is clearly less than 1 in absolute value (as long as  $p$  and  $q$  are not 0 or 1). We thus just need to make sure it is greater than -1. But  $1-p(1-q)-q(1-p) > 1-(1-q)-(1-p) = p+q-1 > -1$ . Thus we may safely use the geometric series formula.

**4.4. Beyond the Pidgeonhole Principle.** Everyone has experience with the Pidgeonhole Principle; what if we ask about having at least  $k$  pidgeons in a box when there are  $N$  pidgeons?

Specifically, consider  $N$  boxes and  $m$  balls, with each ball equally likely to be in each box. For fixed  $k$ , we can bound the probability of at least  $k$  balls being in the same box, as  $N$  and  $m$  tend to infinity. In particular, we can show that if  $m = N^{\frac{k-1}{k}}$  then this probability is at least  $\frac{1}{k!} - \frac{1}{2 \cdot k!^2} + O(N^{-1/k})$  and at most  $\frac{1}{k!} + O(N^{-1/k})$ . We investigated what happens when  $k$  grows with  $N$  and  $m$ , and showed there is negligible probability of having at least  $N$  balls in the same box when  $m = N^{2-\epsilon}$ .

**Project 4.11.** *The arguments in my notes were written a few years ago and in response to a question asked by a colleague. I haven't carefully gone through all my approximations, and almost surely some are wrong, but the general flavor should be right. One should make these arguments rigorous and see how far they can be pushed.*

**4.5. Differentiating identities.** Identities are the bread and butter of mathematics. Thus, if there is a way to generate infinitely more identities from one, then this is a technique one should study! For example, what is  $\sum_{n=1}^{\infty} n2^{-n}$ ?

The starting point in the method of differentiating identities is some known identity, for example, in this case the geometric series formula

$$\sum_{n=0}^{\infty} x^n = (1-x)^{-1}. \quad (4.20)$$

A nice exercise is to show that we can interchange a derivative with respect to  $x$  with the infinite sum. We apply the operator  $xd/dx$  to both sides (we use  $xd/dx$  and not  $d/dx$  so that we end up with  $x^n$  and not  $x^{n-1}$ ), and find

$$\sum_{n=0}^{\infty} nx^n = \frac{x}{(1-x)^2}. \quad (4.21)$$

Taking  $x = 1/2$  we see the answer to our original question is just 2.

This is but one of many formulas which can be proved using this technique (other classic examples are means, variances and moments of probability distributions). Here is another fun example.

Using Induction, it is possible to prove results such as

**Theorem 4.12.** *For  $p$  a positive integer*

$$\sum_{k=1}^n k^p = f_p(n), \quad (4.22)$$

where  $f_p(x)$  is a polynomial of degree  $p+1$  in  $x$  with rational coefficients, and the leading term is  $\frac{x^{p+1}}{p+1}$ .

Everyone is very familiar with the  $p = 1$  case, and perhaps the  $p = 2$  case. One way to figure out the polynomial is to compute the answer for  $p$  or  $p+1$  values of  $n$  and then solve a system of equations. Other ways to prove these results are through Bernoulli polynomials.

It is also possible to prove these results *without* resorting to induction! Namely, we can prove these results by differentiating identities. We need the following result about finite geometric series:

**Lemma 4.13.** For any  $x \in \mathbb{R}$ ,

$$1 + x + x^2 + \cdots + x^n = \frac{x^{n+1} - 1}{x - 1}. \quad (4.23)$$

*Proof.* If  $x = 1$  we evaluate the right hand side by L'Hospital's Rule, which gives  $\frac{n+1}{1} = n + 1$ . For other  $x$ , let  $S = 1 + x + \cdots + x^n$ . Then

$$\begin{aligned} S &= 1 + x + x^2 + \cdots + x^n \\ xS &= x + x^2 + \cdots + x^n + x^{n+1}. \end{aligned} \quad (4.24)$$

Therefore

$$xS - S = x^{n+1} - 1 \quad (4.25)$$

or

$$S = \frac{x^{n+1} - 1}{x - 1}. \quad (4.26)$$

□

We now show how to sum the  $p^{\text{th}}$  powers of the first  $n$  integers. We first investigate the case when  $p = 1$ . Consider the identity

$$\sum_{k=0}^n x^k = \frac{x^{n+1} - 1}{x - 1}. \quad (4.27)$$

We apply the operator  $x \frac{d}{dx}$  to each side and obtain

$$\begin{aligned} x \frac{d}{dx} \sum_{k=0}^n x^k &= x \frac{d}{dx} \frac{x^{n+1} - 1}{x - 1} \\ \sum_{k=0}^n kx^k &= x \frac{(n+1)x^n \cdot (x-1) - 1 \cdot (x^{n+1} - 1)}{(x-1)^2} \\ \sum_{k=0}^n kx^k &= x \frac{nx^{n+1} - (n+1)x^n + 1}{(x-1)^2}. \end{aligned} \quad (4.28)$$

If we set  $x = 1$ , the left hand side becomes the sum of the first  $n$  integers. To evaluate the right hand side we use L'Hospital's rule, as when  $x = 1$  we get  $1 \cdot \frac{0}{0}$ . As long as one of the factors has a limit, the limit of a product is the product of the limits. As  $x \rightarrow 1$ , the factor of  $x$  becomes just 1 and we must study  $\lim_{x \rightarrow 1} \frac{nx^{n+1} - (n+1)x^n + 1}{(x-1)^2}$ . We find

$$\lim_{x \rightarrow 1} \frac{nx^{n+1} - (n+1)x^n + 1}{(x-1)^2} = \lim_{x \rightarrow 1} \frac{n(n+1)x^n - n(n+1)x^{n-1}}{2(x-1)}. \quad (4.29)$$

As the right hand side is  $\frac{0}{0}$  when  $x = 1$  we apply L'Hospital again and find

$$\begin{aligned} \lim_{x \rightarrow 1} \frac{nx^{n+1} - (n+1)x^n + 1}{(x-1)^2} &= \lim_{x \rightarrow 1} \frac{n^2(n+1)x^{n-1} - n(n+1)(n-1)x^{n-1}}{2} \\ &= \frac{n(n+1)}{2}. \end{aligned} \quad (4.30)$$

Therefore, by differentiating the finite geometric series and using L'Hospital's rule we were able to prove the formula for the sum of integers *without* resorting to induction. The reason we used the operator  $x \frac{d}{dx}$  and not  $\frac{d}{dx}$  is this leaves the power of  $x$  unchanged. While this flexibility is not needed to compute sums of first powers of integers, if we want to calculate sums of  $k^p$  for  $p > 1$ , this will simplify the formulas.

**Theorem 4.14.** *For  $n$  a positive integer,*

$$\sum_{k=0}^n k^2 x^k = \frac{n(n+1)(2n+1)}{6}. \quad (4.31)$$

*Proof.* To find the sum of  $k^2$  we apply  $x \frac{d}{dx}$  twice to (4.27) and get

$$\begin{aligned} x \frac{d}{dx} \left[ x \frac{d}{dx} \sum_{k=0}^n x^k \right] &= x \frac{d}{dx} \left[ x \frac{d}{dx} \frac{x^{n+1} - 1}{x - 1} \right] \\ x \frac{d}{dx} \sum_{k=0}^n kx^k &= x \frac{d}{dx} \left[ x \frac{nx^{n+1} - (n+1)x^n + 1}{(x-1)^2} \right] \\ \sum_{k=0}^n k^2 x^k &= x \frac{d}{dx} \left[ \frac{nx^{n+2} - (n+1)x^{n+1} + x}{(x-1)^2} \right] \\ \sum_{k=0}^n k^2 x^k &= x \frac{[n(n+2)x^{n+1} - (n+1)^2 x^n + 1] \cdot (x-1)^2}{(x-1)^4} \\ &\quad - x \frac{[nx^{n+2} - (n+1)x^{n+1} + x] \cdot 2(x-1)}{(x-1)^4}. \end{aligned} \quad (4.32)$$

Simple algebra (multiply everything out on the right hand side and collect terms) yields

$$\sum_{k=0}^n k^2 x^k = x \frac{n^2 x^{n+2} - (2n^2 + 2n - 1)x^{n+1} + (n^2 + 2n + 1)x^n - x - 1}{(x-1)^3} \quad (4.33)$$

The left hand side is the sum we want to evaluate; however, the right hand side is  $\frac{0}{0}$  for  $x = 1$ . As the denominator is  $(x-1)^3$  it is reasonable to expect that we will need to apply L'Hospital's rule three times; we provide a proof of this in Remark 4.15.

Applying L'Hospital's rule three times to the right hand side we find the right hand side is

$$\frac{n^2(n+2)(n+1)nx^{n-1} - (2n^2 + 2n - 1)(n+1)n(n-1)x^{n-2} + (n^2 + 2n + 1)n(n-1)(n-2)x^{n-3}}{3 \cdot 2 \cdot 1}. \quad (4.34)$$

Taking the limit as  $x \rightarrow 1$  we obtain

$$\begin{aligned} \sum_{k=0}^n k^2 x^k &= \frac{n^2(n+2)(n+1)n - (2n^2+2n-1)(n+1)n(n-1) + (n^2+2n+1)n(n-1)(n-2)}{6} \\ &= \frac{n(n+1)(2n+1)}{6}, \end{aligned} \tag{4.35}$$

where the last line follows from simple algebra.  $\square$

**Remark 4.15.** *While we are able to obtain the correct formula for the sum of squares without resorting to induction, the algebra is starting to become tedious, and will get more so for sums of higher powers. After applying  $x \frac{d}{dx}$  twice we had  $\frac{g(x)}{(x-1)^3}$ , where  $g(x)$  is a polynomial of degree  $n+2$  and  $g(1) = 0$ . It is natural to suppose that we need to apply L'Hospital's rule three times as we have a factor of  $(x-1)^3$  in the denominator. However, if  $g'(1)$  or  $g''(1)$  is not zero, then we do not apply L'Hospital's rule three times but rather only once or twice. Thus we really need to check and make sure that  $g'(1) = g''(1) = 0$ . While a straightforward calculation will show this, a moment's reflection shows us that both of these derivatives must vanish. If one of them was non-zero, say equal to  $a$ , then we would have  $\frac{a}{0}$  which is undefined; however, clearly the sum of the first  $n$  squares is finite. Therefore these derivatives will be zero and we do have to apply L'Hospital's rule three times.*

**Remark 4.16.** *For those concerned about the legitimacy of applying L'Hospital's rule and these formulas when  $x = 1$ , we can consider a sequence of  $x$ 's, say  $x_N = 1 - \frac{1}{N}$  with  $N \rightarrow \infty$ . Everything is then well-defined, and it is of course natural to use L'Hospital's rule to evaluate  $\lim_{N \rightarrow \infty} \frac{g(x_N)}{(x_N-1)^3}$ .*

**Project 4.17.** *Can you find a way to make the algebra work in general, or at least prove that one does get a polynomial of the claimed degree?*

## REFERENCES

- [AZ] M. Aigner and G. M. Ziegler, *Proofs from THE BOOK*, Springer-Verlag, Berlin, 1998.
- [AB] U. Andrews IV and J. Blatz, *Distribution of digits in the continued fraction representations of seventh degree algebraic irrationals*, Junior Thesis, Princeton University, Fall 2002.
- [Ap] R. Apéry, *Irrationalité de  $\zeta(2)$  et  $\zeta(3)$* , Astérisque **61** (1979) 11–13.
- [Apo] T. Apostol, *Introduction to Analytic Number Theory*, Springer-Verlag, New York, 1998.
- [Bec] M. Beceanu, *Period of the continued fraction of  $\sqrt{n}$* , Junior Thesis, Princeton University, 2003.
- [Br] T. Brox, *Collatz cycles with few descents*, Acta Arithm. **92** (2000), 181–188.
- [Da1] H. Davenport, *The Higher Arithmetic: An Introduction to the Theory of Numbers*, 7th edition, Cambridge University Press, Cambridge, 1999.
- [Da2] H. Davenport, *Multiplicative Number Theory*, 2nd edition, revised by H. Montgomery, Graduate Texts in Mathematics, Vol. 74, Springer-Verlag, New York, 1980.
- [Da3] H. Davenport, *On the distribution of quadratic residues (mod  $p$ )*, London Math. Soc. **6** (1931), 49–54.
- [Da4] H. Davenport, *On character sums in finite fields*, Acta Math. **71** (1939), 99–121.
- [Di] T. Dimofte, *Rational shifts of linearly periodic continued fractions*, Junior Thesis, Princeton University, 2003.
- [Gl] A. Gliga, *On continued fractions of the square root of prime numbers*, Junior Thesis, Princeton University, 2003.



- [He] P. V. Hegarty, *Some explicit constructions of sets with more sums than differences* (2007), *Acta Arithmetica* **130** (2007), no. 1, 61–77.
- [HM] P. V. Hegarty and S. J. Miller, *When almost all sets are difference dominated*, to appear in *Random Structures and Algorithms*. <http://arxiv.org/abs/0707.3417>
- [JKKKM] D. Jang, J. U. Kang, A. Kruckman, J. Kudo and S. J. Miller, *Chains of distributions, hierarchical Bayesian models and Benford's Law*, to appear in the *Journal of Algebra, Number Theory: Advances and Applications*.
- [JLR] S. Janson, T. Łuczak and A. Ruciński, *Random Graphs*, Wiley, 2000.
- [Ka] S. Kapnick, *Continued fraction of cubed roots of primes*, Junior Thesis, Princeton University, Fall 2002.
- [Kh] A. Y. Khinchin, *Continued Fractions*, 3rd edition, University of Chicago Press, Chicago, 1964.
- [KonSi] A. Kontorovich and Ya. G. Sinai, *Structure theorem for  $(d, g, h)$ -maps*, *Bull. Braz. Math. Soc. (N.S.)* **33** (2002), no. 2, 213–224.
- [KP] L. Kontorovich and P. Ravikumar, *Virus Propagation: Progress Report*, working notes.
- [Kua] F. Kuan, *Digit distribution in the continued fraction of  $\zeta(n)$* , Junior Thesis, Princeton University, Fall 2002.
- [Lag1] J. Lagarias, *The  $3x+1$  problem and its generalizations*. Pages 305–334 in *Organic mathematics (Burnaby, BC, 1995)*, CMS Conf. Proc., vol. 20, AMS, Providence, RI, 1997.
- [Lag2] J. Lagarias, *The  $3x+1$  problem: An annotated bibliography*, preprint.
- [Law1] J. Law, *Kuzmin's theorem on algebraic numbers*, Junior Thesis, Princeton University, Fall 2002.
- [Le] P. Lévy, *Sur les lois de probabilité dont dependent les quotients complets et incomplets d'une fraction continue*, *Bull. Soc. Math.* **57** (1929), 178–194.
- [MO] G. Martin and K. O'Bryant, *Many sets have more sums than differences*, *Additive combinatorics*, 287–305, CRM Proc. Lecture Notes **43**, Amer. Math. Soc., Providence, RI, 2007.
- [Mic1] M. Michelini, *Independence of the digits of continued fractions*, Junior Thesis, Princeton University, Fall 2002.
- [Mic2] M. Michelini, *Kuzmin's extraordinary zero measure set*, Senior Thesis, Princeton University, Spring 2004.
- [MOS] S. J. Miller, B. Orosz and D. Scheinerman, *Constructing infinite families of sum dominated sets*, preprint. <http://arxiv.org/pdf/0809.4621>
- [MT-B] S. J. Miller and R. Takloo-Bighash, *An Invitation to Modern Number Theory*, Princeton University Press, Princeton, NJ, 2006.
- [MN1] S. J. Miller and M. Nigrini, *The Modulo 1 Central Limit Theorem and Benford's Law for Products*, *International Journal of Algebra* **2** (2008), no. 3, 119–130.
- [Na1] M. B. Nathanson, *Problems in additive number theory, 1*. To appear in the Proceedings of CRM-Clay Conference on Additive Combinatorics, Montréal 2006.
- [Na2] M. B. Nathanson, *Sets with more sums than differences*, *Integers : Electronic Journal of Combinatorial Number Theory* **7** (2007), Paper A5 (24pp).
- [Ni1] T. Nicely, *The pentium bug*, <http://www.trnicely.net/pentbug/pentbug.html>
- [Ni2] T. Nicely, *Enumeration to  $10^{14}$  of the Twin Primes and Brun's Constant*, *Virginia J. Sci.* **46** (1996), 195–204.
- [vdP1] A. van der Poorten, *An introduction to continued fractions*. Pages 99–138 in *Diophantine Analysis (Kensington, 1985)*, London Mathematical Society Lecture Note Series, Vol. 109, Cambridge University Press, Cambridge, 1986.
- [vdP2] A. van der Poorten, *Notes on continued fractions and recurrence sequences*. Pages 86–97 in *Number theory and cryptography (Sydney, 1989)*, London Mathematical Society Lecture Note Series, Vol. 154, Cambridge University Press, Cambridge, 1990.
- [vdP3] A. van der Poorten, *Continued fractions of formal power series*. Pages 453–466 in *Advances in Number Theory (Kingston, ON, 1991)*, Oxford Science Publications, Oxford University Press, New York, 1993.

- [vdP4] A. van der Poorten, *Fractions of the period of the continued fraction expansion of quadratic integers*, Bull. Austral. Math. Soc. **44** (1991), no. 1, 155–169.
- [vdP5] A. van der Poorten, *Continued fraction expansions of values of the exponential function and related fun with continued fractions*, Nieuw Arch. Wisk. (4) **14** (1996), no. 2, 221–230.
- [vdP6] A. van der Poorten, *Notes on Fermat’s Last Theorem*, Canadian Mathematical Society Series of Monographs and Advanced Texts, Wiley-Interscience, New York, 1996.
- [PS1] A. van der Poorten and J. Shallit, *Folded continued fractions*, J. Number Theory **40** (1992), no. 2, 237–250.
- [PS2] A. van der Poorten and J. Shallit, *A specialised continued fraction*, Canad. J. Math. **45** (1993), no. 5, 1067–1079.
- [RV1] G. Rhin and C. Viola, *On the irrationality measure of  $\zeta(2)$* , Ann. Inst. Fourier (Grenoble) **43** (1993), no. 1, 85–109.
- [RV2] G. Rhin and C. Viola, *On a permutation group related to  $\zeta(2)$* , Acta Arithm. **77** (1996), 23–56.
- [Rie] H. J. J. te Riele, *On the sign of the difference  $\pi(x) - \text{Li}(x)$* , Mathematics of Computation **48** (1987), no. 177, 323–328.
- [Ru1] I. Z. Ruzsa, *On the cardinality of  $A + A$  and  $A - A$* , Combinatorics year (Keszthely, 1976), vol. 18, Coll. Math. Soc. J. Bolyai, North-Holland-Bolyai Társulat, 1978, 933–938.
- [Ru2] I. Z. Ruzsa, *Sets of sums and differences*, Séminaire de Théorie des Nombres de Paris 1982–1983 (Boston), Birkhäuser, 1984, 267–273.
- [Ru3] I. Z. Ruzsa, *On the number of sums and differences*, Acta Math. Sci. Hungar. **59** (1992), 439–447.
- [Sk] S. Skewes, *On the difference  $\pi(x) - \text{Li}(x)$* , J. London Math. Soc. **8** (1933), 277–283.
- [Sim] J. L. Simmons, *Post-transcendence conditions for the existence of  $m$ -cycles for the  $3x + 1$  problem*, preprint.
- [SimWe] J. L. Simmons and B. M. M. Weger, *Theoretical and computational bounds for  $m$ -cycles of the  $3n + 1$  problem*, Acta Arithm. **117** (2005), 51–70.
- [Si] Ya. G. Sinai, *Statistical  $(3x + 1)$  problem*, Comm. Pure Appl. Math. **56** (2003), no. 7, 1016–1028.
- [Sin] M. K. Sinisalo, *On the minimal cycle lengths of the Collatz sequences*, preprint, Univ. of Oulu, Finland.
- [So] K. Soundararajan, *Small gaps between prime numbers: The work of Goldston-Pintz-Yildirim*, Bull. of the AMS **44** (2007), no. 1, 1–18.
- [Ta] C. Taylor, *The Gamma function and Kuzmin’s theorem*, Junior Thesis, Princeton University, Fall 2002.
- [Wir] E. Wirsing, *On the theorem of Gauss-Kuzmin-Lévy and a Frobenius-type theorem for function spaces*, Acta Arith. **24** (1974) 507–528.

*E-mail address:* Steven.J.Miller@williams.edu

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