

On the binary expansions of algebraic numbers

par DAVID H. BAILEY, JONATHAN M. BORWEIN,
RICHARD E. CRANDALL et CARL POMERANCE

RÉSUMÉ. En combinant des concepts de théorie additive des nombres avec des résultats sur les développements binaires et les séries partielles, nous établissons de nouvelles bornes pour la densité de 1 dans les développements binaires de nombres algébriques réels. Un résultat clef est que si un nombre réel y est algébrique de degré $D > 1$, alors le nombre $\#(|y|, N)$ de 1 dans le développement de $|y|$ parmi les N premiers chiffres satisfait

$$\#(|y|, N) > CN^{1/D}$$

avec un nombre positif C (qui dépend de y), la minoration étant vraie pour tout N suffisamment grand. On en déduit la transcendance d'une classe de nombres réels $\sum_{n \geq 0} 1/2^{f(n)}$ quand la fonction f , à valeurs entières, croît suffisamment vite, disons plus vite que toute puissance de n . Grâce à ces méthodes on redémontre la transcendance du nombre de Kempner–Mahler $\sum_{n \geq 0} 1/2^{2^n}$; nous considérons également des nombres ayant une densité sensiblement plus grande de 1. Bien que le nombre $z = \sum_{n \geq 0} 1/2^{n^2}$ ait une densité de 1 trop grande pour que nous puissions lui appliquer notre résultat central, nous parvenons à développer une analyse fine de théorie des nombres avec des calculs étendus pour révéler des propriétés de la structure binaire du nombre z^2 .

ABSTRACT. Employing concepts from additive number theory, together with results on binary evaluations and partial series, we establish bounds on the density of 1's in the binary expansions of real algebraic numbers. A central result is that if a real y has algebraic degree $D > 1$, then the number $\#(|y|, N)$ of 1-bits in the expansion of $|y|$ through bit position N satisfies

$$\#(|y|, N) > CN^{1/D}$$

Manuscrit reçu le 20 mars 2003.

Bailey's work is supported by the Director, Office of Computational and Technology Research, Division of Mathematical, Information, and Computational Sciences of the U.S. Department of Energy, under contract number DE-AC03-76SF00098.

Borwein's work is funded by NSERC and the Canada Research Chair Program.

for a positive number C (depending on y) and sufficiently large N . This in itself establishes the transcendency of a class of reals $\sum_{n \geq 0} 1/2^{f(n)}$ where the integer-valued function f grows sufficiently fast; say, faster than any fixed power of n . By these methods we re-establish the transcendency of the Kempner–Mahler number $\sum_{n \geq 0} 1/2^{2^n}$, yet we can also handle numbers with a substantially denser occurrence of 1's. Though the number $z = \sum_{n \geq 0} 1/2^{n^2}$ has too high a 1's density for application of our central result, we are able to invoke some rather intricate number-theoretical analysis and extended computations to reveal aspects of the binary structure of z^2 .

1. Introduction

Research into the statistical character of digit expansions is often focused on the concept of normality. We call a real number b -normal if its base- b digits are random in a certain technical sense (see [31], [21], [3], and references therein). Qualitatively speaking, b -normality requires every string of k consecutive base- b digits to occur, in the limit, $1/b^k$ of the time, as if the digits are generated by tossing a “fair” b -sided die. In spite of the known fact that almost all numbers are b -normal (in fact almost all are absolutely normal, meaning b -normal for *every* base $b = 2, 3, \dots$) not a single, shall we say “genuine” fundamental constant such as $\pi, e, \log 2$ is known to be b -normal for any b . Artificially constructed normals are known, such as the 2-normal binary Champernowne number [9]

$$C_2 = (0.11011100101110\dots)_2,$$

obtained by sheer concatenation of the binary of positive integers. Previous research that motivates the present work includes [3], where a certain “Hypothesis A” relevant to chaotic maps is shown to imply 2-normality of $\pi, \log 2, \zeta(3)$; and [4], where the historical work of Korobov, Stoneham and others is augmented to establish b -normality of, shall we say, “less artificial” constants such as the numbers $\sum_{n \geq 0} 1/(c^n b^{c^n})$ where $b, c > 1$ are coprime. Intriguing connections with yet other fields—such as ergodic theory—are presented in [22].

Of interest for the present work is that all real algebraic irrationals are widely believed—shall we say suspected—to be absolutely normal (and this belief is at least a half-century old; see for example [6, 7]). This suspicion is based on numerical and visual evidence that the digit expansions of algebraics do appear empirically “random.” Yet again, the mathematical situation is as bleak as can be: Not a single algebraic real is known to be b -normal, nor has a single algebraic real irrational been shown not to be b -normal; all of this regardless of the base b . Though we expect every irrational algebraic real is absolutely normal, for all we know it could even

be that some algebraics are absolutely abnormal, i.e. not b -normal for any b whatsoever (absolutely abnormal numbers do exist; see [28]).

Herein we focus on the binary scenario $b = 2$, and though we do not achieve normality results per se, we establish useful lower bounds on the occurrence of 1-bits in positive algebraics. Our central result is that if y is a real algebraic of degree $D > 1$, then there exists a positive number C (depending only on y) such that for sufficiently large N the number $\#(|y|, N)$ of 1's in the binary expansion of $|y|$ through the N -th bit position satisfies

$$\#(|y|, N) > CN^{1/D}.$$

To achieve this bound we borrow ideas from additive number theory; in particular we employ the notion of additive representations. This notion is combined with our own bounds on the count of 1-bits resulting from binary operations, and also with previous observations on “BBP tails” that arise from arbitrary left-shifts of infinite series. In Section 6 we define and elaborate on BBP tails.

Irrational numbers y for which $\#(|y|, N)$ cannot achieve the above bound for any degree D are necessarily transcendental. In this way we easily re-establish the transcendency of the Kempner–Mahler number

$$M = \sum_{n \geq 0} \frac{1}{2^{2^n}},$$

first shown to be transcendental by Kempner [19], but the transcendency cannot be established directly from the celebrated Thue–Siegel–Roth theorem on rational approximations to algebraics (there are interesting anecdotes concerning Mahler’s approach to such an impasse, including his results on p -adic Thue–Siegel theory and his functional methods; see [26, 27, 29]). Incidentally, the number M above is sometimes called the Fredholm number, but this attribution may be historically erroneous [35]. (See also [1] for more on the number M .)

We can also handle numbers having a higher density of 1's than does M . For example, by our methods the Fibonacci binary

$$X = \sum_{n \geq 0} \frac{1}{2^{F_n}}$$

having 1's at Fibonacci-number positions $0, 1, 1, 2, 3, 5 \dots$ is transcendental. Now X was proved transcendental some decades ago [27] and explicit irrationality measures and certain continued-fraction properties are known for X [36]. In the present treatment, we can handle numbers like X but where the growth of the exponents is more general than the classic growth of the F_n .

With our central result we establish the transcendency of numbers whose 1-bits are substantially more dense than in the above examples, an example of such a “denser” number being

$$\sum_{n \geq 3} \frac{1}{2^{\lfloor n \log \log n \rfloor}}.$$

Incidentally, in the late stages of the present research project we found that this notion of “digital thinking” to establish results in analysis had been foreshadowed by a specific, pedagogical proof by M. Knight [20] that for any base $b > 1$

$$\sum_{n \geq 0} \frac{1}{b^{2^n}}$$

is transcendental (note that $b = 2$ gives the number M above). The author used terms such as “islands” for flocks of digits guarded on both sides by enough zeros to avoid carry problems when integral powers of a real number are taken. As will be seen, such notions pervade also our own treatment; however our results pertain to general 1-bit densities and not to specific real numbers. Other historical foreshadowings of our approach exist [34] [25]. (See also our Section 11 on open problems.)

Aside from transcendency results, we can employ the central theorem to establish bounds on the algebraic degree. For example, we shall see that

$$\sum_{n \geq 0} \frac{1}{2^{n^k}}, \quad \sum_{p \text{ prime}} \frac{1}{2^{p^k}}$$

must have algebraic degrees at least $k, k + 1$ respectively. (In this context we think of a transcendental number as having infinite degree.) Thus for example, $\sum 1/2^{p^2}$ must be an at-least-cubic irrational.

There are interesting numbers that do not fall under the rubric of our central theorem, such as the “borderline” case:

$$z = \sum_{n \geq 0} \frac{1}{2^{n^2}} = \frac{1}{2} \left(1 + \theta_3 \left(\frac{1}{2} \right) \right),$$

where θ_3 is the standard Jacobi theta function. The problem is that $\#(z, N) \sim \sqrt{N}$, so our central theorem does not give any information on the algebraic degree of z . Yet we are able to use further number-theoretical analysis—notably the theory of representations of integers as sums of two squares—to establish quadratic irrationality for z . We further argue, on the basis of such analysis, that z^2 has almost all 0’s, and more precisely that the 1’s count through the N -th bit position has a certain asymptotic behavior. Incidentally the number z , being essentially the evaluation of a theta function at an algebraic argument, is known to be transcendental

by other methods [30, 5, 12]. We stress that our binary approach is an apparently new way to look at such issues.

2. Additive representations

For any real nonnegative number x we consider the binary expansion

$$x = (\dots x_{-3}x_{-2}x_{-1}x_0 . x_1x_2x_3 \dots)_2.$$

The assignment of (finitely many) nonpositive indices for bits x_i to the left of the decimal (or if you will, binary) point is a convenience, for we shall, of course, be concentrating a great deal on the bits to the right. We adopt the convention that no x can end with infinitely many successive 1's, and this forces uniqueness of the binary expansion. Next we denote the 1's-position set of x by

$$\mathcal{P}(x) = \{p : x_p = 1\},$$

and further define $r_1(x, p) = 1$ if $x_p = 1$, else 0. (The rationale for the notation “ r_1 ” will be momentarily evident.) Now the number of 1-bits through bit position N inclusive is

$$\#(x, N) = \sum_{m \leq N} r_1(x, m) = \sum_{p \in \mathcal{P}(x), p \leq N} 1.$$

Note that when x is a nonnegative integer, $\#(x, 0)$ is the number of 1's in x . So for example $\#(7, 0) = 3$. On the other hand, for $x = \sqrt{2} = (1.011010100\dots)_2$, say, we have $\#(x, 0) = 1$, $\#(x, 5) = 4$, and so on.

We next introduce the representation count

$$r_d(x, n) = \#\{(p_1, \dots, p_d) \in \mathcal{P}^d : p_1 + \dots + p_d = n\},$$

just as in additive number theory where one studies representations of integers n as sums of primes, or squares, and so on. It is evident that r_d can be expressed as an acyclic convolution:

$$r_d(x, n) = \sum_{i+j=n} r_{d-1}(x, i)r_1(x, j).$$

We shall also employ a step-function on integers r , namely $H(r) = 1$ if $r > 0$, else 0. Thus $H(r_d(n)) = 1$ signifies that n has at least one representation $p_1 + \dots + p_d$. For our analysis it is a simple but useful combinatorial observation that the count of representables, call it

$$\rho_d(N) = \sum_{n \leq N} H(r_d(x, n)),$$

satisfies

$$(1) \quad \rho_d(N) \leq \sum_{n \leq N} r_d(x, n) \leq \#(x, N)^d.$$

Also of use will be an attractive relation for positive integral powers of x :

$$x^d = \sum_n \frac{r_d(x, n)}{2^n}.$$

Unfortunately it is in general extremely difficult to convert partial knowledge of the representation sequence $(r_d(x, n))$ into precise results on the binary expansion of x^d . The problem is that of carry: A summand $r_d(x, n)/2^n$ possibly causes carry, about $\lg r_d$ positions to the left, and thus the summands interfere (herein $\lg x$ means the base-2 logarithm of x). It can be said that the goal of the present treatment is the circumvention of this carry problem.

An instructive digression is appropriate here. With a view to additive number theory, let us define the number

$$G = \sum_{p \text{ prime, odd}} \frac{1}{2^{(p-1)/2}} = (0.11101101 \dots)_2.$$

Note that $\#(G, N) = \pi(2N+2) - 1$, where π is the standard prime-counting function. Then $r_2(G, N)$ is precisely the number of representations of $2N+2$ as a sum of two odd primes. Even if we knew the truth of the Goldbach conjecture—in this scenario, that every $N > 2$ has $H(r_2(x, N)) = 1$ —we would still not immediately know the binary expansion of G^2 , because of the carry problem. For all we know, it could be that the question of irrationality for G^2 is more difficult than the Goldbach conjecture itself. Conversely, it is unclear whether complete knowledge of the binary expansion of G^2 would yield results on the celebrated conjecture. In fact, it is easy to see that $r_2(G, N)$ is unbounded, so arbitrarily long carries (deposition of bits arbitrarily far to the left of a given position) can be expected.

Similarly, for the number $z = \sum_{n \geq 0} 2^{-n^2}$ introduced earlier we know that z^4 has a representation sequence $(r_4(z^4, 0), r_4(z^4, 1), \dots)$ of *all* positive entries, on the basis of the Lagrange theorem that every nonnegative integer is a sum of four squares. Here again, little can be gleaned about the binary expansion of z^4 from this perspective, again because of carry. We study the number z further in Section 9.

Now back to positive powers of x and representation lists. A sum we later call a “tail component” defined

$$T_d(x, R) = \sum_{m \geq 1} \frac{r_d(x, R + m)}{2^m},$$

which we note is 2^R times a partial series for the power x^d , can be bounded via combinatorial observations, as in

Theorem 2.1. For $x \in (1, 2)$ (whether algebraic or not) and $d \geq 1$ we have

$$r_d(x, n) \leq \binom{n + d - 1}{d - 1}.$$

Moreover, for the sum T_d defined above, we have for $0 \leq R \leq N$ the upper bound

$$T_d(x, R) < \frac{(R + d)^d}{(d - 1)!(R + 1)} \leq \frac{(N + d)^d}{(d - 1)!(N + 1)}$$

and the average bound

$$\sum_{0 \leq R \leq N} T_d(x, R) < \frac{(N + d)^d}{(d - 1)!}.$$

Proof. From the convolution

$$r_d(x, n) = \sum_{i+j=n} r_{d-1}(x, i)r_1(x, j)$$

we have

$$r_d(x, n) \leq \sum_{i_1, \dots, i_d \in [0, n], \sum i_j = n} 1 = \binom{n + d - 1}{d - 1}.$$

Thus $T_d(x, R) \leq U_d(R)$ where

$$U_d(R) = \sum_{m \geq 1} \frac{1}{2^m} \binom{R + m + d - 1}{d - 1}.$$

This expression is seen to satisfy the recurrence relation

$$U_d(R) = 2U_{d-1}(R) + \binom{R + d - 1}{d - 1},$$

which can be used to establish a finite form for U_d :

$$U_d(R) = \sum_{j=0}^{d-1} \binom{R + d}{j}.$$

So we have

$$U_d(R) < \frac{(R + d)^{d-1}}{(d - 1)!} \sum_{n \geq 0} \left(\frac{d - 1}{R + d}\right)^n = \frac{(R + d)^d}{(d - 1)!(R + 1)}.$$

Thus, the first bound follows. The bound on the sum $\sum T_d$ is simply obtained by summing the first bound over the stated range of R . \square

Remark. The finite form for $U_d(R)$ noted above is a polynomial in R with nonnegative coefficients and with main term $R^{d-1}/(d - 1)!$, so that this expression is not only a lower bound for $U_d(R)$, but is also equal to

it asymptotically. Moreover, it is possible to express $U_d(R)$ as a hypergeometric integral:

$$U_d(R) = \frac{(R + d)!}{R!(d - 1)!} \int_0^1 (2 - x)^{d-1} x^R dx.$$

We admit that the bounds of Theorem 2.1 and the present remark are actually stronger than what we need here; however, such stronger bounds could be useful in future research.

3. Preliminary bound on 1's density

Let x be a real algebraic irrational. The Thue–Siegel–Roth theorem [33] says that for any $\epsilon > 0$ the inequality

$$\left| x - \frac{a}{b} \right| < \frac{1}{b^{2+\epsilon}}$$

has at most finitely many integer-pair solutions a, b . This means that the 1-bits of such an x cannot be too far apart, in the sense of

Theorem 3.1. *For a real positive algebraic irrational number x , and any $\delta > 0$, the 1's positions $p_i \in \mathcal{P}(x)$ satisfy, for sufficiently large m ,*

$$p_m < (2 + \delta)p_{m-1}.$$

Furthermore, for sufficiently large k , the interval

$$\left(\left\lfloor \frac{k}{2 + \delta} \right\rfloor, k \right)$$

always contains a 1's position. Finally, the 1's count through sufficiently large position n satisfies

$$\#(x, n) > (1 - \delta) \lg n.$$

Proof. When x is irrational, $\mathcal{P}(x)$ is an infinite set, so arbitrarily large p_i can be chosen, with

$$x - \sum_{p \in \mathcal{P}(x), p \leq p_i} \frac{1}{2^p} < \frac{2}{2^{p_{i+1}}}.$$

Now the sum is a rational a/b with $b = 2^{p_i}$, and so the first inequality of the theorem is clearly satisfied if p_i is large enough. The rest of the conclusions are immediate from said inequality. □

The bound $\#(x, n) > (1 - \delta) \lg n$ is admittedly weak, relative to what we aim to prove later. It does, however, establish the transcendency of any number

$$m_\alpha = \sum_{n \geq 0} \frac{1}{2^{\lfloor \alpha^n \rfloor}}$$

for any real $\alpha > 2$. Note that the Kempner–Mahler number $M = m_2 = \sum_{n \geq 0} 1/2^{2^n}$ lies just out of reach of the Thue–Siegel–Roth implications. We shall be able to use our binary approach to establish, in fact, the transcendency of m_α for any real $\alpha > 1$.

There is a curious aspect to Theorem 3.1, namely, however weak the bounds on 1’s counts may be, there is a crucial juncture in what follows (the central Theorem 7.1) where we need Theorem 3.1 to assail the ubiquitous problem of carry propagation.

4. Bounds for binary evaluations

For nonnegative integers n we have defined $\#(n, 0)$ as the number of 1’s in the binary expansion of n . We proceed to give convexity relations on binary evaluations, i.e. on sums and products of integers, starting with some simple observations:

Lemma 4.1. *For integers $n > 0, j \geq 0$, we have*

$$\#(n, 0) \leq 1 + \lg n \leq n,$$

$$\#(2^j n, 0) = \#(n, 0),$$

$$\#(n + 2^j, 0) = \#(n, 0) + 1 - k_j,$$

where in the last relation k_j is the number of consecutive 1’s in n counting from the $(-j)$ -th position inclusive, to the left.

Proof. The first inequality follows from the observation that the total number $B(n)$ of bits in n (counting 0’s and 1’s) satisfies $2^{B(n)-1} \leq n$, and $\#(n, 0) \leq B(n)$. The second statement is obvious (left-shifting by j bits introduces no new 1’s). The third statement follows by the simple rule of add-with-carry. □

This lemma leads to

Theorem 4.2 (Convexity relations). *For nonnegative integers m, n we have upper bounds on the 1’s counts of evaluations, as*

$$\#(m + n, 0) \leq \#(m, 0) + \#(n, 0),$$

$$\#(mn, 0) \leq \#(m, 0)\#(n, 0).$$

Proof. The first, additive relation follows by repeated application of the last equality of Lemma 4.1, one application for each 1-bit of m , say. The second, multiplicative relation follows in similar fashion, by writing $mn = (\sum 2^{-p})n$, where p runs through the 1’s positions of m , and using the second (shift) relation of Lemma 4.1. □

Back to the manageable case of upper bounds for binary evaluations, consider the polynomial

$$f(x) = A_D x^D + A_{D-1} x^{D-1} + \dots + A_1 x + A_0,$$

for integers A_i all nonnegative. Then from Lemma 4.1 and Theorem 4.2 we easily have, for nonnegative integers n , the following convexity relation for polynomial evaluations:

$$\#(f(n), 0) \leq \sum_{d=0}^D \max(0, 1 + \lg A_d) \#(n, 0)^d \leq f(\#(n, 0)).$$

This relation will next be applied to algebraic numbers whose minimum integer polynomial has all coefficients (except A_0) nonnegative.

5. Application of binary-evaluation bounds

Our strongest bounds on 1's density will be obtained for the class of real algebraic irrationals for which the coefficients of the minimum integer polynomial are nonnegative, except for the constant term. We begin with

Lemma 5.1. *For irrational $x \in (1, 2)$ and a given integral power d , the inequality*

$$0 < x^d - \frac{\lfloor 2^N x \rfloor^d}{2^{Nd}} < \frac{d2^d}{2^N}$$

holds for all sufficiently large N .

Proof. Setting $i = \lfloor 2^N x \rfloor$, we have $2^N \leq i < 2^{N+1}$, and $x = i/2^N + z$, where $z \in (0, 1/2^N)$. Now

$$x^d = \frac{i^d}{2^{Nd}} (1 + 2^N z/i)^d,$$

so that

$$0 < x^d - \frac{i^d}{2^{Nd}} < \frac{i^d}{2^{Nd}} ((1 + 1/i)^d - 1).$$

Choose M such that $d < i$ for $N > M$, whence

$$x^d - \frac{i^d}{2^{Nd}} < \frac{i^d}{2^{Nd}} \frac{2d}{i} < (e - 1) d \frac{2^{(N+1)(d-1)}}{2^{Nd}} < \frac{d2^d}{2^N}.$$

□

We are now in a position to state

Theorem 5.2. *Let y be a real algebraic of degree $D > 1$ and assume for $x = |y|/2^{\lfloor \lg |y| \rfloor}$ a minimum integer polynomial equation*

$$A_D x^D + A_{D-1} x^{D-1} + \dots + A_1 x + A_0 = 0,$$

where $A_D > 0$ and A_{D-1}, \dots, A_1 are nonnegative integers. Then for any $\epsilon > 0$ we have

$$\#(|y|, N) > (1 - \epsilon)(1 + \lg A_D)^{-1/D} N^{1/D}$$

for sufficiently large N (with threshold depending on y, ϵ).

Proof. Note that $x \in (1, 2)$ and because x is a shift of y , the counts $\#(x, N), \#(y, N)$ differ only by an integer constant, so we may concentrate on x . Observe that A_0 is a negative integer. From Lemma 5.1 we can select N and assign $i = \lfloor 2^N x \rfloor$ so that

$$x^d = \frac{i^d}{2^{Nd}} + z_d,$$

with $z_d \in (0, d2^d/2^N)$, for $1 \leq d \leq D$. Now define the integer

$$Y_N = 2^{ND} \sum_{d=1}^D A_d i^d 2^{-Nd} = 2^{ND} \sum_{d=1}^D A_d (x^d - z_d),$$

so that

$$-A_0 = \frac{Y_N}{2^{ND}} + z_N,$$

where

$$0 < z_N = \sum_{d=1}^D A_d z_d < \frac{1}{2^N} \sum_{d=1}^D d2^d A_d.$$

(It is this last inequality where the signs of the $A_i, i > 0$ are essential.) Thus for sufficiently large N we have a fractional part

$$\left\{ \frac{Y_N}{2^{ND}} \right\} = \{-A_0 - z_N\} = \{1 - z_N\} \geq 1 - \frac{C}{2^N} > 0$$

for a strictly positive constant $C = \sum_{d=1}^D d2^d A_d$ independent of N . But this means for some constant C' (also independent of N) that the integer Y_N has more than $N - C'$ 1-bits. Since $\#(i, 0) = \#(x, N)$, on using Theorem 4.2, we have (again using in an essential way that $A_i \geq 0$ for $i > 0$)

$$\begin{aligned} (2) \quad N - C' &< \#(Y_N, 0) \leq \sum_{d=1}^D \max(0, 1 + \lg A_d) \#(x, N)^d \\ (3) \quad &\leq \#(x, N)^D \left((1 + \lg A_D) + \frac{A_{D-1}}{\#(x, N)} + \dots \right), \end{aligned}$$

and since $\#(x, N)$ is unbounded (x is irrational) the result follows. □

A side result is

Corollary 5.3. *If $y > 0$ is irrational, and there exists an integer $d > 1$ such that for every $\eta > 0$ we have*

$$\#(y, N) < \eta N^{1/d}$$

for infinitely many N , then y^d is also irrational.

Proof. Assuming y^d is rational then for $x = y/2^{\lfloor \lg y \rfloor}$ there is a polynomial $Ax^d - B$, with positive integers A, B . This polynomial is of the required form for application of Theorem 5.2, whose conclusion contradicts $\liminf \#(y, N)N^{-1/d} = 0$. \square

So for example the number

$$\left(\sum_{n \geq 0} \frac{1}{2^{n^5}} \right)^4$$

is irrational; the number being 4-th-powered does not, in the sense of Corollary 5.3, have enough 1-bits.

Theorem 5.2 reveals that the assignments $y = \sqrt{2}$ or $y = (-1 + \sqrt{5})/2$ (the golden mean) each have $\#(y, N) > (1 - \epsilon)\sqrt{N}$ for large enough N ; in the latter case one may use the polynomial equation $x^2 + 2x - 4 = 0$, whose root $-1 + \sqrt{5}$ is in $(1, 2)$. On a historical note: J. Samborski, in a published problem [34], asked for a proof that $\#(y, N) < 5 \cdot 2^{N-2}$ —an interesting, hard bound but asymptotically very much weaker than our square-root density.

6. Bounds on BBP tails

Now we desire to lift all restrictions on the coefficient signs, except the high coefficient $A_D > 0$ and contemplate the following representation relation (in this section we assume $x \in (1, 2)$ is algebraic of degree $D > 1$, see the remarks opening the proof of Theorem 5.2):

$$A_D x^D + \dots + A_1 x + A_0 = 0 = A_0 + \sum_{n \geq 0} \frac{1}{2^n} \sum_{d=1}^D A_d r_d(x, n).$$

Consider a shift by R bits of all entities, so that

$$-2^R A_0 = I(x, R) + T(x, R),$$

where $I(x, R)$ is an integer and the BBP tail is defined

$$T(x, R) = \sum_{m \geq 1} \frac{1}{2^m} \sum_{d=1}^D A_d r_d(x, R + m) = \sum_{d=1}^D A_d T_d(x, R),$$

where as in Section 2 we identify a tail component

$$T_d(x, R) = \sum_{m \geq 1} \frac{r_d(x, R + m)}{2^m}.$$

The concept of BBP tail comes from the Bailey–Borwein–Plouffe formalism [2], whereby one may rapidly compute isolated bits of certain binary expansions—such as for $\pi, \log 2$ —by rapid computation of the integer $I(x, R)$ and rigorous control of the “tail” $T(x, R)$.

Remarkably, it is a fact that for the algebraic x in question, $T(x, R)$ is always an integer, for the simple reason that $T(x, R) = -2^R A_0 - I(x, R)$. To facilitate further analysis, we shall require a bound on the average absolute value of the tails $T(x, R)$ in terms of one value of $\#(x, N)$:

Lemma 6.1. *Let x be an algebraic number in $(1, 2)$ of degree $D > 1$ with minimum integer polynomial $A_D x^D + A_{D-1} x^{D-1} + \cdots + A_0$, so $A_D > 0$. Let $N \geq 2D$ and set $K = \lceil 2D \lg N \rceil$. Then for $1 \leq d \leq D$ we have*

$$\sum_{1 \leq R \leq N-K} T_d(x, R) < \#(x, N)^d + 1,$$

and so

$$\sum_{1 \leq R \leq N-K} |T(x, R)| < \sum_{d=1}^D |A_d| \left(\#(x, N)^d + 1 \right).$$

Proof. We have

$$\begin{aligned} \sum_{R \leq N-K} T_d(x, R) &= \sum_{m \geq 1} 2^{-m} \sum_{R \leq N-K} r_d(x, R + m) \\ &\leq \sum_{m=1}^K 2^{-m} \sum_{R \leq N} r_d(x, R) \\ &\quad + 2^{-K} \sum_{m > K} 2^{K-m} \sum_{R \leq N-K} r_d(x, R + m) \\ &< \sum_{R \leq N} r_d(x, R) + 2^{-K} \sum_{K \leq R \leq N} T_d(x, R). \end{aligned}$$

Using (1) and Theorem 2.1 we have

$$\sum_{1 \leq R \leq N-K} T_d(x, R) \leq \#(x, N)^d + N^{-2D} (N + 1)^d,$$

and the lemma is proved. \square

We shall use Lemma 6.1 to show that if $\#(x, N)$ is small, then not too many values of $T(x, R)$ are positive. Counter to this, the following lemma gives conditions on when there are many positive tails $T(x, R)$.

Lemma 6.2. *Let x be an algebraic number in $(1, 2)$ of degree $D > 1$. Suppose that $R_0 < R_1$ are positive integers with $r_{D-1}(x, R) = 0$ for all integers $R \in (R_0, R_1]$ and $T(x, R_1) > 0$. Then $T(x, R) > 0$ for every integer $R \in [R_0, R_1]$.*

Proof. Say the minimum integer polynomial for x is $A_D x^D + A_{D-1} x^{D-1} + \dots + A_0$. As the 0-bit of x is 1 it follows that $r_d(x, N) \geq r_{d-1}(x, N)$ for $d \geq 2$. Thus the hypothesis implies that for each $d = 1, 2, \dots, D-1$ we have $r_d(x, R) = 0$ for each integer $R \in (R_0, R_1]$. From the general recurrence relation on tails,

$$T(x, R - 1) = \frac{1}{2}T(x, R) + \frac{1}{2} \sum_{d=1}^D A_d r_d(x, R) = \frac{1}{2}T(x, R) + \frac{1}{2}A_D r_D(x, R).$$

Assuming inductively that $T(x, R) > 0$, and using $A_D > 0$, we get $T(x, R - 1) > 0$. □

7. The central theorem regarding general real algebraic numbers

We have established that for a certain restricted class of algebraics y of degree $D \geq 2$,

$$\#(|y|, N) > (1 - \epsilon)(1 + \lg A_D)^{-1/D} N^{1/D}$$

for sufficiently large N , where A_D is the leading coefficient of the minimum integer polynomial for the normalized algebraic $x = |y|/2^{\lfloor \lg |y| \rfloor}$. Now we move to general algebraics, so that there will be no coefficient constraints except for the natural $A_D > 0$. Fortunately, we shall achieve a bound which is weaker only by an overall constant factor.

Theorem 7.1. *For real algebraic y of degree $D > 1$ and for any $\epsilon > 0$ we have for sufficiently large N (with threshold depending on y, ϵ)*

$$\#(|y|, N) > (1 - \epsilon)(2A_D)^{-1/D} N^{1/D},$$

where $A_D > 0$ is the leading coefficient of the minimum integer polynomial of $x = |y|/2^{\lfloor \lg |y| \rfloor}$.

Proof. As in the proof of Theorem 5.2, we use the normalized algebraic $x \in (1, 2)$, observing that $\#(x, N), \#(|y|, N)$ differ only by an integer constant, so again we may concentrate on the bit-counting for x . Suppose $\#(x, N) \leq cN^{1/D}$. Then from (1) applied for $d = D - 1$, and the fact that each $r_{D-1}(x, R)$ is a nonnegative integer, we have that the number of integers $R \leq N$ with $r_{D-1}(x, R) > 0$ is at most $c^{D-1}N^{1-1/D}$. Say these R 's are $0 = R_1 < R_2 < \dots < R_M$, where $M \leq c^{D-1}N^{1-1/D}$. Let $R_{M+1} = N$. Trivially we have

$$\sum_{i=1}^M (R_{i+1} - R_i) = N.$$

For $\delta > 0$, let I denote the set of numbers $i \leq M$ such that $R_{i+1} - R_i \geq \frac{\delta}{3}c^{1-D}N^{1/D}$. (Ultimately we transform δ into the ϵ of the theorem.) We have

$$\sum_{i \in I} (R_{i+1} - R_i) \geq N - \frac{\delta}{3}c^{1-D}N^{1/D}M \geq \left(1 - \frac{\delta}{3}\right)N.$$

Now we wish to show, if $i \in I$ and if integer $R \in (R_i, R_{i+1} - D \log N]$ has $r_D(x, R) > 0$, then $T(x, R - 1) > 0$:

$$\begin{aligned} T(x, R - 1) &\geq \frac{1}{2}A_D - \sum_{d=1}^{D-1} |A_d| \sum_{m \geq 1} 2^{-m} r_d(x, R - 1 + m) \\ &= \frac{1}{2}A_D - \sum_{d=1}^{D-1} |A_d| \sum_{m > R_{i+1} - R} 2^{-m} r_d(x, R - 1 + m) \\ &= \frac{1}{2}A_D - \sum_{d=1}^{D-1} |A_d| 2^{R - R_{i+1}} T_d(x, R_{i+1} - 1) \\ &\geq \frac{1}{2} - N^{-D} \sum_{d=1}^{D-1} |A_d| (N + d)^d / (d - 1)! N, \end{aligned}$$

where this last inequality follows from Theorem 2.1. Thus, for sufficiently large N , the positivity of the tail $T(x, R - 1)$ for such an R is established. Now if $r_1(x, j) > 0$ and $i \leq M$ then $r_D(x, R_i + j) > 0$. A key observation now is: By the Thue–Siegel–Roth implication Theorem 3.1, for N sufficiently large and for any $i \in I$, there is some integer j_i with

$$j_i \in \left(\frac{1}{2 + \delta/2} (R_{i+1} - R_i - D \log N), R_{i+1} - R_i - D \log N \right)$$

and $r_1(x, j_i) > 0$. We conclude that $r_D(x, R_i + j_i) > 0$, so it follows from our previous reasoning that $T(x, R_i + j_i - 1) > 0$. Then from Lemma 6.2 we have $T(x, R) > 0$ for every integer $R \in [R_i, R_i + j_i]$. Hence $T(x, R) > 0$ for at least

$$\frac{1}{2 + \delta/2} \sum_{i \in I} (R_{i+1} - R_i - D \log N)$$

values of $R \leq N$. But this last expression is at least

$$\frac{1}{2 + \delta/2} \left(1 - \frac{\delta}{3}\right) N - Dc^{D-1}N^{1-1/D} \log N$$

which is at least $(\frac{1}{2} - \frac{\delta}{3})N$ for all sufficiently large values of N .

We now show that if c is too small, this last conclusion is impossible. By Lemma 6.1 we have (with K as in the lemma)

$$\begin{aligned} \sum_{R \leq N-K} |T(x, R)| &\leq \sum_{d=1}^D |A_d| (\#(x, N)^d + 1) \\ &\leq \sum_{d=1}^D |A_d| (c^d N^{d/D} + 1) \\ &= A_D c^D N + O(N^{1-1/D}). \end{aligned}$$

Suppose now that $c \leq ((2 + \delta)A_D)^{-1/D}$. It follows from this last calculation and the fact that each $T(x, R)$ is an integer that $T(x, R) > 0$ for at most $\frac{1}{2+\delta}N + O(N^{1-1/D})$ values of $R \leq N$. So for N sufficiently large, this assertion is incompatible with the assertion that $T(x, R) > 0$ for at least $(\frac{1}{2} - \frac{\delta}{3})N$ values of $R \leq N$. Finally, for the arbitrary positive δ we set $\epsilon = 1 - (1 + \delta/2)^{-1/D}$ to obtain the statement of the theorem. \square

8. Implications of the central theorem

Theorem 7.1 can be used to establish transcendency of a class of binary expansions, as in

Theorem 8.1. *Let a function $f : \mathcal{R} \rightarrow \mathcal{R}$ be strictly increasing, with f attaining integer values for integer arguments. If for any $\epsilon > 0$ the inverse of f satisfies*

$$f^{-1}(y) = O(y^\epsilon)$$

then the number

$$x = \sum_{n \geq 0} \frac{1}{2^{f(n)}}$$

is transcendental.

Proof. Note that the bit positions $f(n)$ are distinct, so the observation

$$f^{-1}(N) = \#\{n > 0 : f(n) \leq N\} = \#(x, N)$$

means $\#(x, N) = O(N^\epsilon)$, which for algebraic x is incompatible with Theorem 7.1. \square

Corollary 8.2. *For any real $\alpha > 1$ the number*

$$m_\alpha = \sum_{n \geq 0} \frac{1}{2^{\lfloor \alpha^n \rfloor}}$$

is transcendental. So the Kempner–Mahler number $M = m_2$ and the Fibonacci binary involving the Fibonacci numbers $(F_n) = (0, 1, 1, 2, 3, 5, \dots)$

$$X = \sum_{n \geq 0} \frac{1}{2^{F_n}}$$

are transcendental. Finally, there are transcendental numbers of still greater 1-bit densities, such as

$$Y = \sum_{n \geq 3} \frac{1}{2^{\lceil n^{\log \log n} \rceil}}.$$

Remark. Recall that the Thue–Siegel–Roth implication Theorem 3.1 handles $\alpha > 2$.

Proof. As for m_α , take $n_0 = \lceil -\log(\alpha - 1)/\log \alpha \rceil$ so that there is a strictly monotone function whose integer evaluations are $f(n) = \lceil \alpha^{n+n_0} \rceil$, with $f^{-1}(N) = O(\log N)$, so that Theorem 8.1 applies and the partial binary sum for m_α starting from index n_0 , hence m_α itself, is transcendental. As for the Fibonacci binary, the n -th Fibonacci number can be written $f(n) = ((1 + \tau)^n - (-\tau)^n)/\sqrt{5}$, where $\tau = (\sqrt{5} - 1)/2$, so the growth of 1's positions is essentially that of $m_{1+\tau}$ and again incompatible with Theorem 7.1 if X is assumed algebraic. For the number Y it is evident that $\#(Y, N) \sim N^{1/\log \log N}$ which is of slower growth than any positive power of N . \square

We can also use Theorem 7.1 to generate results on algebraic degrees for certain constants, as in the following (as before let us stipulate that the algebraic degree of a transcendental is ∞):

Theorem 8.3. For positive integer k the number

$$X_k = \sum_{n \geq 0} \frac{1}{2^{n^k}}$$

has algebraic degree at least k , while the number

$$P_k = \sum_{p \text{ prime}} \frac{1}{2^{p^k}}$$

has algebraic degree at least $k + 1$.

Proof. In the first case, $\#(X_k, N) = \#\{n^k \leq N\} < CN^{1/k}$, so by Theorem 7.1 we must have degree $D \geq k$. In the second case we have $\#(P_k, N) = \#\{p^k \leq N\} = \pi(N^{1/k}) < AN^{1/k}/(\log N)$ for a constant A , so that again by Theorem 7.1, we must have $D > k$. \square

Thus for example neither P_2 nor X_3 is a quadratic irrational. The case of X_2 is just the previously mentioned number $z = \sum 1/2^{n^2}$, on which number we focus attention in the next section.

9. Study of a “borderline” number

The number

$$z = \sum_{n \geq 0} \frac{1}{2^{n^2}} = \frac{1}{2} \left(\theta_3 \left(\frac{1}{2} \right) + 1 \right)$$

is, with respect to the present treatment, a “borderline” case because, as we have seen, a square-root density of 1’s is beyond reach of our methods. Recall also as in Section 4 that there are numbers with the same essential density of 1’s as z but for which products of such numbers can be rational. Note that

$$z' = 2z - 1 = \theta_3 \left(\frac{1}{2} \right)$$

so that

$$z'^2 = \sum_{n \geq 0} \frac{r_2(n)}{2^n},$$

where now we are using the standard notation of $r_2(n)$ for the number of representations $n = a^2 + b^2$ for $a, b \in \mathbf{Z}$, counting sign and order. It will be convenient therefore to study z' , from which algebraic properties of z follow. Incidentally z'^2 has some interesting numerological features; for one thing it is very close to $\pi/\log 2$; in fact the approximation

$$z'^2 \approx \frac{\pi}{\log 2} \left(1 + 2e^{-\pi^2/\log 2} \right)^2 = 4.53237201425897410082795 \dots$$

can be obtained via Jacobi θ -transformation, and remarkably is correct to the implied 23 decimal places in the above display. It is fascinating that such relations between z'^2 and fundamental constants exist even though, as we shall prove, almost all of the binary bits of z'^2 are 0’s.

It is one of the earliest results in additive number theory, due to Jacobi, that

$$r_2(n) = 4 \sum_{d|n, d \text{ odd}} (-1)^{(d-1)/2}.$$

It turns out that the representation count $r_2(n)$ is positive if and only if every prime $p \equiv 3 \pmod{4}$ dividing n does so to an even power. Thus, the representation sequence $(r_2(0), \dots)$ has zeros in any position $n = 3k$ with $(3, k)$ coprime, and so on. Deeper results on r_2 include the fact that the number of representable integers not exceeding N behaves according to the Landau theorem:

$$\sum_{n \leq N} H(r_2(n)) \sim L \frac{N}{\sqrt{\log N}},$$

where the Landau constant is

$$L = \left(\frac{1}{2} \prod_{p \equiv 3 \pmod{4}} \left(1 - \frac{1}{p^2} \right)^{-1} \right)^{1/2} = 0.764223653 \dots$$

(See [15] for descriptions of this and other facets of sums of squares.) The Landau density of representable numbers does not on the face of it imply a similar density of 1-bits in the expansion of z'^2 .

Evidently we have

$$z'^2 = 4 \sum_{d \text{ odd}} \frac{(-1)^{(d-1)/2}}{2^d - 1}.$$

This form is reminiscent of the Erdős–Borwein number

$$E = \sum_{n>0} \frac{1}{2^n - 1} = \sum_{m>0} \frac{d(m)}{2^m},$$

where $d(m)$ denotes the number of divisors of m . The constant E was proven irrational by Erdős [13] who used number-theoretical arguments (outlined in [4]) which did, in fact, motivate our present analysis of z'^2 . Later the irrationality of such forms was established via Padé approximants, by P. Borwein [8].

What we shall show is that z' is not a quadratic irrational, and so neither is z . In one sense this is stronger than the quoted irrationality results for the number E . On the other hand, it is already known that theta functions at algebraic arguments, hence z, z' , are transcendental [5, 12]. To effect our nonquadratic-irrationality proof, we shall follow the same basic prescription as for Theorem 7.1; namely, we establish upper bounds on the size of representations, and employ some knowledge of zero-runs. As for upper bounds, it is known [17] that for any fixed $\epsilon > 0$ we have

$$r_2(n) < 2^{(\frac{1}{2} + \epsilon) \frac{\log n}{\log \log n}}$$

for sufficiently large n . Note that this bound is much tighter than the general one of Theorem 2.1. This tighter bound works well with what we can show about zero-runs:

Theorem 9.1. *Let $\epsilon > 0$ be arbitrary, but fixed, and define*

$$u_\epsilon(x) = \frac{1 - \epsilon}{2L} \frac{\log x}{\sqrt{\log \log x}}$$

where L is the Landau constant. Then for sufficiently large x there is a square integer M with $M < x$ and an integer $a < M$ such that $r_2(n) = 0$ whenever

$$n \equiv a + i \pmod{M}, \quad 1 \leq i \leq u_\epsilon(x).$$

Proof. Let x be large and let

$$u = \left\lfloor \frac{1 - \epsilon}{2L} \frac{\log x}{\sqrt{\log \log x}} \right\rfloor.$$

Let $v_p(n)$ denote the exponent on the prime p in the prime factorization of n . Cast out from $[1, u]$ any integer o with $v_p(o)$ odd for some prime $p < u/\log u$, $p \equiv 3 \pmod{4}$. Let \mathcal{E} denote the set of remaining numbers n , and let E denote the cardinality of \mathcal{E} . Also, let \mathcal{E}_1 denote the number of integers in $[1, u]$ which are the sum of two squares, and let E_1 denote the cardinality of \mathcal{E}_1 . By the Landau theorem,

$$E_1 \sim \frac{Lu}{\sqrt{\log u}}.$$

Clearly, $\mathcal{E}_1 \subset \mathcal{E}$. In particular, $E - E_1$ is at most the number of integers $n \in [1, u]$ divisible by some prime p with $u/\log u \leq p \leq u$ and $p \equiv 3 \pmod{4}$. Then

$$E - E_1 \leq u \sum_{u/\log u \leq p \leq u} \frac{1}{p} = O(u \log \log u / \log u).$$

We conclude that

$$E \sim \frac{Lu}{\sqrt{\log u}}.$$

Label the members of \mathcal{E} as n_1, n_2, \dots, n_E .

Next, let $M_1 = \prod p^{a_p}$, where p runs over the primes with $p \equiv 3 \pmod{4}$, $p < u/\log u$, and $a_p = 2\lceil (\log u)/(2 \log p) \rceil$. (Thus, a_p is the least even integer with $p^{a_p} \geq u$.) We have $\log M_1 = O(u/\log u)$.

Let

$$v = \left\lfloor \frac{\log x}{1 + \epsilon} \right\rfloor,$$

and let $M_2 = \prod p^2$ where p runs over the primes $p \equiv 3 \pmod{4}$ with $u/\log u \leq p \leq v$. Then $\log M_2 \sim v$, so that for x sufficiently large we have $M := M_1 M_2 < x$. Label the prime factors of M_2 as p_1, p_2, \dots, p_F , where $F \sim v/(2 \log v)$. We have

$$E \sim \frac{Lu}{\sqrt{\log u}} \sim \frac{1 - \epsilon}{2} \frac{\log x}{\log \log x},$$

$$F \sim \frac{v}{2 \log v} \sim \frac{1}{2(1 + \epsilon)} \frac{\log x}{\log \log x},$$

so that for x sufficiently large we have $F \geq E$.

For $1 \leq i \leq E$ let r_i be a solution to

$$n_i + r_i M_1 \equiv p_i \pmod{p_i^2}.$$

Further, let the integer r satisfy

$$r \equiv r_i \pmod{p_i^2}, \quad \text{for } 1 \leq i \leq E.$$

Let h be an arbitrary integer. If n is an integer in $[1, u]$ that is not in \mathcal{E} , then $v_p(n + rM_1 + hM)$ is odd for some prime $p|M_1$ and so $r_2(n) = 0$. Suppose $n = n_i \in \mathcal{E}$. Then $v_{p_i}(n + rM_1 + hM) = 1$, and so $r_2(n) = 0$. Thus, with $a = rM_1$ we have that $r_2(a + i + hM) = 0$ for $1 \leq i \leq u$. This completes the proof of the theorem. \square

Corollary 9.2. *For integer n sufficiently large, the interval $(n^2, n^2 + n)$ contains a zero-run of the r_2 representation of length at least $u_\epsilon(n)$.*

Proof. Take $x = n/3$ in Theorem 9.1. Then for relevant M and a , the position

$$n^2 + (a + 1 + M - (n^2 \bmod M)) \leq n^2 + 2n/3$$

is the start of a zero-run of length $u_\epsilon(n/3) \sim u_\epsilon(n)$, which run for sufficiently large n is contained in $(n^2, n^2 + n)$. \square

We are now in a position to use representation bounds and the zero-run bound of Theorem 9.1, to establish

Theorem 9.3. *The number $z = \sum_{n \geq 0} 1/2^{n^2}$ is not a quadratic irrational.*

Proof. We shall focus on the number $z' = \sum_{n \in \mathbb{Z}} 1/2^{n^2}$ whence the result will follow for z . Assume that

$$A_2 z'^2 + A_1 z' + A_0 = 0,$$

Consider the interval $[n^4, (n^2 + 1)^2]$ and within that, positions

$$n^4, n^4 + f, n^4 + f + Z, n^4 + n^2, (n^2 + 1)^2.$$

By Corollary 9.2, for sufficiently large n , these positions are in order, with a zero-run length $Z = \lfloor u_\epsilon(n^2) \rfloor$, so that $(r_2(n^4 + f + 1), \dots, r_2(n^4 + f + Z))$ is a length- Z zero-vector. Note also that $r_1(n^4) = r_1((n^2 + 1)^2) = 1$, yet every other r_1 in the entire interval is zero. Thus

$$T(z', n^4 + n^2 - 1) \geq \frac{1}{2} - \frac{|A_1|}{2n^2} > 0.$$

Thus any tail $T(z', n^4), \dots, T(z', n^4 + n^2 - 1)$ is positive. However, using the upper bound on $r_2(n)$ to bound the tail component $T_2(z', n^4 + f)$, we get

$$T(z', n^4 + f) \leq \frac{2A_2}{2Z} 2^{(2+4\epsilon) \log n / \log \log n} + \frac{|A_1|}{2^{2n^2 - f}}.$$

Since Z has the $\sqrt{\log \log n}$ denominator, we have for sufficiently large n

$$0 < T(z', n^4 + f) < 1,$$

a contradiction. \square

We now state the following result, which was first suggested to us by numerical computation.

Theorem 9.4. *Almost all bits of z^2 are 0's; in fact the 1's-count has asymptotic behavior*

$$\#(z^2, N) \sim C_0 \frac{N}{\sqrt{\log N}},$$

for an absolute constant $C_0 \approx 0.7996\dots$ (we give a formula for C_0 in the proof).

Clearly, Theorem 9.4 implies that z^2 is irrational (because of arbitrarily long zero-runs), and Theorem 9.3 may well follow also from the asymptotic 1's density (although see Section 11). Incidentally the asymptotic density also holds for z'^2 , as follows from a slight modification (actually simplification) of the proof. In spite of the paucity of 1's for these squared numbers, higher powers such as z^3, z^4 are likely 2-normal. Indeed, all such higher powers will involve interfering carries. For example, it is known that $r_3(n) > 0$ for a limiting fraction 5/6 of all n (see [15]), so the carry problem for z^3 is already formidable.

The proof of Theorem 9.4 is based on the following two lemmas.

Lemma 9.5. *There is an absolute constant c such that for any integers $N, B \geq 2$, the number of integers $n \leq N$ with $r_2(z, n) > 0$ and $r_2(z, m) > 0$ for some integer m with $0 < |n - m| < B$ is at most $cBN/\log N$.*

Lemma 9.6. *For any positive integers B, N , the number of integers $n \leq N$ with $r_2(z, n) \geq B$ is at most $(\sqrt{N} + 1)^2/B$.*

Note that Lemma 9.6 is very easy. The assertion follows instantly from the inequality $\sum_{n \leq N} r_2(z, n) \leq (\sqrt{N} + 1)^2$. We postpone the proof of Lemma 9.5 until later. First we see how Theorem 9.4 follows from the lemmas.

Proof. (Theorem 9.4.) Let $b(m) = \#(m, 0)$ denote the number of 1's in the binary representation of the nonnegative integer m , and let $b(0) = 0$. It follows from Theorem 4.2 and the fact that $r_2(z, n) \leq n^{o(1)}$ that for N large,

$$\#(z^2, N) \leq \sum_{n \leq N + \log N} b(r_2(z, n)).$$

The goal is to get a similar-looking lower bound. Let S_N denote the set of natural numbers $n \leq N$ such that n is not a square and

$$\begin{aligned} r_2(z, n) &> 0, \\ r_2(z, m) &= 0 \text{ for } 0 < |n - m| < 3 \lg \lg N, \\ r_2(z, m) &< (\lg N)^2 \text{ for } |n - m| < 2 \lg N. \end{aligned}$$

Note that if $n \in S_N$ and N is sufficiently large then

$$(4) \quad \sum_{m>n} \frac{r_2(z, m)}{2^m} < \frac{1}{2^n}.$$

Indeed, we first note that

$$\sum_{m \geq n+2 \lg N} \frac{r_2(z, m)}{2^m} < \sum_{m \geq n+2 \lg N} \frac{m}{2^m} = O\left(\frac{1}{N2^n}\right).$$

Next note that

$$\begin{aligned} \sum_{n+2 \lg N > m > n} \frac{r_2(z, m)}{2^m} &= \sum_{n+2 \lg N > m \geq n+3 \lg \lg N} \frac{r_2(z, m)}{2^m} \\ &\leq \sum_{n+2 \lg N > m \geq n+3 \lg \lg N} \frac{(\lg N)^2}{2^m} \\ &= O\left(\frac{1}{2^n \lg N}\right). \end{aligned}$$

Thus, we have (4). Further, for $n \in S_N$ and N large we have

$$\sum_{m \geq n} \frac{r_2(z, m)}{2^m} < \frac{r_2(z, n) + 1}{2^n} < \frac{(\lg N)^2 + 1}{2^n} < \frac{1}{2^{n'}},$$

where $n' < n$ is the largest number with $r_2(z, n') > 0$. We conclude from these estimates that appearing in the bit stream for z^2 we see intact all of the bits of the numbers $r_2(z, n)$ for $n \in S_N$, when N is large. Thus, we have for large N that

$$\#(z^2, N) \geq \sum_{n \in S_N} b(r_2(z, n)).$$

It follows from the lemmas that the number of integers $n \leq N$ with $r_2(z, n) > 0$ that are *not* in S_N is $O(N \log \log N / \log N)$. The number of 1-bits contributed to $\#(z^2, N)$ from $n \leq N$ with $n \notin S_N$ and $r_2(z, n) < 2^{(\log N)^{1/4}}$ is at most

$$O\left((\log N)^{1/4} \frac{N \log \log N}{\log N}\right) = o\left(\frac{N}{\sqrt{\log N}}\right).$$

And, by Lemma 9.6 there are at most $O(N/2^{(\log N)^{1/4}})$ values of $n \leq N$ with $r_2(z, n) > 2^{(\log N)^{1/4}}$. Since $b(r_2(z, n)) = o(\log n)$, the contribution of these values of n to $\#(z^2, N)$ is also $o(N/\sqrt{\log N})$. It follows that

$$\#(z^2, N) = \sum_{n \leq N} b(r_2(z, n)) + o(N/\sqrt{\log N}).$$

Using the identity $z'^2 = 4z^2 - 4z + 1$ and that $r_2(n) = 4r_2(z, n)$ when n is not a square, and $r_2(n) = 4r_2(z, n) - 4 \geq 0$ when n is a positive square, we further see that

$$\#(z^2, N) = \sum_{n \leq N} b(r_2(n)) + o(N/\sqrt{\log N}).$$

Hence it is sufficient to estimate this last sum.

Suppose $n = n_1 n_2 n_3$ where n_i is the largest divisor of n composed of primes that are congruent to $i \pmod{4}$. We have $r_2(n) > 0$ if and only if n_3 is a square. And if n_3 is a square, then $r_2(n)/4 = d(n_1)$, where d is the standard divisor function. It follows that if n_3 is a square and if $g(n)$ denotes the largest squarefull divisor of n_1 then $r_2(n)/d(g(n))$ is a power of 2, so that

$$b(r_2(n)) = b(d(g(n))).$$

Incidentally by squarefull is meant an integer none of whose prime factors appears to the power 1.

We now count the number $T_g(N)$ of integers $n \leq N$ with $r_2(n) > 0$ and such that $g(n) = g$, where g is a given squarefull integer all of whose primes are congruent to 1 (mod 4). It is not too difficult to see that

$$T_g(N) \sim L \frac{N}{\sqrt{\log N}} \frac{\alpha}{g} \prod_{p|g} \left(1 - \frac{1}{p}\right) \left(1 - \frac{1}{p^2}\right)^{-1},$$

where

$$\alpha = \prod_{p \equiv 1 \pmod{4}} \left(1 - \frac{1}{p^2}\right) = \frac{16L^2}{\pi^2},$$

and where p in these formulae runs over primes. Letting $\psi(g) = g \times \prod_{p|g} (1 + 1/p)$, we thus have that

$$T_g(N) \sim \frac{16L^3}{\pi^2} \frac{N}{\psi(g)\sqrt{\log N}}.$$

Hence, we have Theorem 9.4 with

$$C_0 = \frac{16L^3}{\pi^2} \sum_g \frac{b(d(g))}{\psi(g)},$$

where g runs over the squarefull integers divisible solely by primes that are congruent to 1 (mod 4). Note that this sum is convergent, which convergence partially justifies the adding of the asymptotic relations for $T_g(N)$. □

We do not give a proof of the asymptotic relations for $T_g(N)$, but these can be achieved as corollaries of the Landau asymptotic formula. In Section

10 we give numerical verification of Theorem 9.4. We close the present section with a proof of Lemma 9.5.

Proof. (Lemma 9.5.) Let $r'(n)$ denote the number of coprime representations of n as the sum of two squares. First we count the number of integers $n \leq N$ for which $r'(n) > 0$ and for which $r'(m) > 0$ for some integer m with $0 < |n - m| < B$. Note that for $r'(m)$ to be positive it is necessary that m is not divisible by any prime congruent to 3 (mod 4), that is, that $m_3 = 1$. (This condition is almost sufficient: to make it sufficient it should also be the case that m_2 , the 2-power in m , is not a power of 4.) For a given integer $k > 0$, the number of integers $n \leq N$ with both $r'(n) > 0, r'(n + k) > 0$ is, by Theorem 2.3 in [16], at most

$$c'\psi(k)\frac{N}{\log N},$$

where c' is an absolute constant and where ψ is defined in the proof of Theorem 9.4. (Actually one can have the smaller factor $\psi(d_3)$, but this is unimportant.) Since

$$\sum_{k \leq B} \psi(k) = O(B),$$

as is easily seen by elementary methods (see [17], Ch. 18), it follows that the number of $n \leq N$ with $r'(n) > 0, r'(n + k) > 0$ for some integer k with $0 < |k| < B$ is $O(BN/\log N)$. This proves the lemma for the function r' . To get it for $r_2(z, n)$ we generalize the above proof for the case $u^2|n, v^2|n+k$, where uv is divisible only by primes that are congruent to 3 (mod 4) and where $r'(n/u^2) > 0, r'((n+k)/v^2) > 0$. For any fixed choice for u, v we get an estimate of $O(\psi(k)N/(u^2v^2 \log N))$ for the number of such $n \leq N$. Now we sum over k, u, v getting the lemma. \square

10. Numerical experiments for C_0

The intricacies of the borderline number z and its powers show that global bit-density arguments alone are insufficient to handle low 1's-density cases: We required number theory to focus on certain details of the bit pattern. Later in the research, we found that computational aspects—such as bit-counting—for z^2 are nontrivial. In attempts to verify Theorem 9.4 empirically—in particular, to justify the value of C_0 —the present authors were met with considerable computational consternation. There are two basic difficulties that need be overcome. Note that calculation of C_0 from the sum formula is not too hard, and gives us the cited 0.7996... value that we obtained by summing over squarefull $g \leq 10^5, 10^6, 10^7$ in succession, then extrapolating on the assumption of a reasonable form for the series-truncation error. The remaining difficulties all pertain to the actual counting of representations up through some large n .

The first difficulty is that the Landau asymptotic formula is, for all practical purposes, generally below the mark, in the sense that a more accurate formula, also due to Landau, is [23]

$$\sum_{n \leq N} H(r_2(n)) = L \frac{N}{\sqrt{\log N}} \left(1 + \frac{C_1}{\log N} \right) + o\left(\frac{N}{(\log N)^{3/2}} \right),$$

where $C_1 = 0.5819\dots$ is yet another constant. One might also use the Ramanujan form [18]

$$\sum_{n \leq N} H(r_2(n)) \sim L \int_0^N \frac{dx}{\sqrt{\log x}}$$

which is reminiscent of the logarithmic integral $\text{Li}(x)$ which, as is well known, stands as a better approximation to $\pi(x)$ than the classic $x/\log x$. Remarkably, for the original Landau expression $LN/\sqrt{\log N}$ to be accurate to say 1 per cent of an empirical count of representables, one has to go up to about $N = e^{50} \approx 10^{22}$, or about a mole of bits. That is comparable to the total digital storage presently available on the entire planet.

The second computational difficulty is that the proof of Theorem 9.4 basically tells us that most r_2 values eventually “separate” so that carries do not interfere. When does separation become significant? A very rough heuristic runs as follows. For very large N the mean separation between positive representation counts is about $\sqrt{\log N}/L$, and this should be greater than the base-2 logarithm of the largest r_2 values of the region. So, and again this is quite heuristic, for significant separation we should have

$$\frac{\sqrt{\log N}}{\log \log N} \approx \frac{2}{L},$$

which leads to the estimate of $N \approx 10^{80}$ bits, which is an oft-quoted estimate on the number of protons in the visible universe.

So these difficulties required the authors to calculate entities that converge to reasonable values for N well below the aforementioned astronomical thresholds. It turns out that the quantity

$$C(N) = L \frac{\sum_{n \leq N} b(r_2(n))}{\sum_{n \leq N} H(r_2(n))} \sim C_0$$

is relatively well-behaved, and gives an excellent empirical value for C_0 (although, still, N has to be taken painfully far and extrapolation techniques were required when our machinery reached its limit, as described below). Notice that this quantity $C(N)$ essentially measures—up to the L factor—the number of bits per *positive* representation. The numerical results reported below suggest (and Theorem 9.4 implies) an amusing principle: The average number of bits in a positive representation count $r_2(n)$

is about $C_0/L = 1.05$. That is, it is a very good bet that a random positive r_2 value is a power of 2.

In order to calculate entities for n up to and beyond say 10^8 , we employed a sieving expedient described in [11] for rapidly obtaining long strings of representation counts. The algorithm is quite simple:

(1) Create $N + 1$ bins, say $B[0], \dots, B[N]$ intended to hold representation counts,

then set $B[0] = 1$ and zero all other bins.

(2) for(odd $d \leq N$) {
 if($d \equiv 1 \pmod{4}$) add 4 to every $B[kd]$ with $kd \leq N$;
 else subtract 4 from every $B[kd]$ with $kd \leq N$;
 }

The result of this algorithm is that $r_2(n)$ is sitting in bin $B[n]$ for every $n \in [0, N]$. In the following table we denote by $\sum r_2$ the sum of counts $r_2(n)$ through $n = N$, by $\sum H$ the count of representables, by c_2 the number of $r_2(n)$ being a power-of-two, and by $\sum b$ the sum of all bit counts $b(r_2(n))$:

N	$LN/\sqrt{\log N}$	$\sum H$	$\sum r_2$	c_2	$\sum b$	$\#(z^2, N)$	$C(N)$
10^6	205605.6	216342	3141552	204082	228646	213480	0.807683
10^7	1903573.9	1985460	31416028	1877532	2093748	1968680	0.805901
10^8	17805966.8	18457848	314159056	17482500	19436147	18353248	0.804725

Note the following interesting features of this table. The 1st-order Landau estimate (2nd column) indeed lags behind the representation count $\sum H$. Next, the similarity of $\sum r_2$ and the decimal digits of π is, of course, not a coincidence: The celebrated Gauss circle problem starts with the proven estimate $\sum r_2 = \pi N + o(N)$, whence one focuses attention on the little- o term. Thus the present sieve technique for counting representations might be useful in numerical studies of the circle error. We see that indeed all but a few per cent of representables are a power-of-two (i.e., c_2 is close to $\sum H$). We stress that the $\#(z^2, N)$ column comes from processing of the r_2 with carry. Actually, rather than work out the carry chain for the $B[]$ bins, we instead obtained the $\#$ column by simply squaring high-precision z values. For both sheer arithmetic of that sort, or for our divisor-sieve, our machinery was not able to go up to $N = 10^9$. Incidentally this is not because of CPU power—the sieve is quite fast, as are convolution methods of squaring reals—the problem is memory. Thus we are forced to extrapolate from the last column of data. Under Romberg extrapolation and the assumption of exponential approach, we estimate the final column's limit to be about 0.802. This extrapolation is within 0.3 per cent of the theoretical $C_0 = 0.7996\dots$ and so we deem this numerical exercise successful. We should also mention that the $C(N)$ results, when other N values are

included in a larger table, are remarkably smooth, which in itself suggests the validity of extrapolation.

11. Open problems

Finally, we state some open problems:

- Is there a quantifiable sense in which the binary representation of $\sqrt{2}$ is *not* truly random? That is, we observe using the present techniques, that the representation list $(r_2(\sqrt{2}, 0), \dots, r_2(\sqrt{2}, N))$ for large N *cannot* have any zero-run of more than $2 \lg N$ consecutive zeros. Evidently this is a hard constraint that one would not want to put on the representation sequence for a “truly random” bit generator.
- As for the Fibonacci binary X of Corollary 8.2, what is the 1’s density of X^2 ?
- What can be said about Fourier representations and bit densities? For example, for $x = \sqrt{2}$ the simple fact of $x^2 = 2$ can be recast as

$$2 = \int_0^1 \left(\sum_{p \in \mathcal{P}(x)} \frac{e^{2\pi i p t}}{2^{\lambda p}} \right)^2 \frac{dt}{1 - 2^{\lambda-1} e^{-2\pi i t}},$$

where $\lambda \in (0, 1)$ is a free parameter, and in principle such an integral representation should convey *some* information about the 1’s positions p in the expansion of $\sqrt{2}$.

- For bases $b > 2$ there is the difficulty of having more than two possible digits. What kinds of bounds might be placed on counts of 1’s and 2’s for ternary expansions of algebraic numbers?
- We have mentioned that the nonquadratic-irrationality Theorem 9.3 may well follow from the density Theorem 9.4. But there is an impasse which would have to be overcome. Namely, it turns out that, whereas $\#(z, N) \sim \sqrt{N}$, there exist reals y with much greater than the square-root 1’s density but such that *still* we have $\#(z + y, N) \sim \sqrt{N}$. That is, adding y to z does not improve the 1’s count. To see this, define

$$y = \sum_{n>1} \frac{2^{k_n} - 1}{2^{n^2}}$$

with arbitrary positive integers k_n except for the constraint $k_n < 2n - 1$. When z is added to such a y , the sets of k_n 1’s are each obliterated by carry. This is an impasse because one cannot just infer the 1’s density of $Az^2 + Bz$ merely by observing that the (rather high) density in Theorem 9.4 dominates the \sqrt{N} density.

- The C_0 constant of Theorem 9.4 we have estimated as $0.7996\dots$. A fast algorithm for calculating the Landau constant L itself is given

in [14]. (In doing so, they effectively solved research Problem 1.88 of [10], which asks for a fast method.) Might there be a similar, fast construction for C_0 ?

- In recent times has emerged the field of “experimental mathematics,” wherein one uses high-precision numerical relations such as linear reduction to suggest exact algebraic identities, in this way igniting a profusion of new results and theorems. One might say that the results of the present paper amount to a kind of “digitally motivated analysis” (we might abbreviate DMA), in which computers were *not* used (except to check various claims), yet results in the analysis field are obtained by thinking digitally, in our case thinking in binary. (And, we acknowledge the historical foreshadowing of “DMA”, as in [24, 20, 34, 25].) A question, then, would be: What other aspects of analysis apart from transcendency might succumb to the “DMA” approach?

12. Acknowledgments

We are grateful to our colleagues J. Buhler, G. Fee, S. Wagon, and S. Wolfram for aid in this research. We owe special debts to J. Shallit and A. J. van der Poorten who provided us with historical perspective. We would also like to thank a referee who provided valuable comments and alerted us to several important references.

References

- [1] J.-P. ALLOUCHE, J. SHALLIT, *Automatic Sequences; Theory, Applications, Generalizations*. Cambridge University Press, 2003.
- [2] DAVID H. BAILEY, PETER B. BORWEIN, SIMON PLOUFFE, *On The Rapid Computation of Various Polylogarithmic Constants*. *Mathematics of Computation* **66** no. **218** (1997), 903–913.
- [3] DAVID H. BAILEY, RICHARD E. CRANDALL, *On the Random Character of Fundamental Constant Expansions*. *Experimental Mathematics* **10** (2001), 175–190.
- [4] DAVID H. BAILEY, RICHARD E. CRANDALL, *Random generators and normal numbers*. *Experimental Mathematics*, to appear.
- [5] D. BERTRAND, *Theta functions and transcendence*. *Ramanujan Journal* **1** (1997), 339–350.
- [6] É. BOREL, *Sur les chiffres décimaux de $\sqrt{2}$ et divers problèmes de probabilités en chaîne*. *C. R. Acad. Sci. Paris* **230** (1950), 591–593.
- [7] É. BOREL, *Oeuvres d'É. Borel Vol. 2*. Éditions du CNRS, Paris, 1972, 1203–1204.
- [8] PETER BORWEIN, *On the Irrationality of Certain Series*. *Mathematical Proceedings of the Cambridge Philosophical Society*, **112** (1992), 141–146.
- [9] D. G. CHAMPERNOWNE, *The Construction of Decimals Normal in the Scale of Ten*. *Journal of the London Mathematical Society* **8** (1933), 254–260.
- [10] R. CRANDALL, C. POMERANCE, *Prime Numbers: A Computational Perspective*. Springer-Verlag, New York, 2002.
- [11] R. CRANDALL, S. WAGON, *Sums of squares: computational aspects*. Manuscript, 2002.

- [12] D. DUVERNEY, KEIJI. NISHIOKA, KUMIKO NISHIOKA, I. SHIOKAWA, *Transcendence of Jacobi's theta series and related results*. Number Theory. Diophantine, Computational and Algebraic Aspects, Kálmán Yöry (ed.) et al. , Walter de Gruyter, Berlin (1998), 157–168.
- [13] P. ERDŐS, *On Arithmetical Properties of Lambert Series*. Journal of the Indian Mathematical Society (N.S.) **12** (1948), 63–66.
- [14] P. FLAJOLET, I. VARDI, *Zeta Function Expansions of Classical Constants*. Manuscript (1996), available at <http://pauillac.inria.fr/algo/flajolet/Publications/Landau.ps>
- [15] E. GROSSWALD, *Representations of Integers as Sums of Squares*. Springer-Verlag, New York, 1985.
- [16] H. HALBERSTAM, H.-E. RICHERT, *Sieve Methods*. Academic Press, London, 1974.
- [17] G. H. HARDY, E. M. WRIGHT, *An Introduction to the Theory of Numbers*. Oxford University Press, 1979.
- [18] G. H. HARDY, *Ramanujan: Twelve Lectures on Subjects Suggested by His Life and Work*, 3rd ed. New York: Chelsea, 1999, 9–10, 55, and 60–64.
- [19] AUBREY J. KEMPNER, *On Transcendental Numbers* Transactions of the American Mathematical Society **17** (1916), 476–482.
- [20] M. J. KNIGHT, *An 'Ocean of Zeroes' Proof That a Certain Non-Liouville Number is Transcendental*. American Mathematical Monthly **98** (1991), 947–949.
- [21] L. KUIPERS, H. NIEDERREITER, *Uniform Distribution of Sequences*. Wiley-Interscience, New York, 1974.
- [22] J. LAGARIAS, *On the Normality of Fundamental Constants*. Experimental Mathematics **10**, no. **3** (2001), 353–366.
- [23] E. LANDAU, *Handbuch der Lehre von der Verteilung der Primzahlen, Bd. II*, 2nd ed. New York: Chelsea, 1953, 641–669.
- [24] S. LEHR, *Sums and rational multiples of q -automatic sequences are q -automatic*. Theor. Comp. Sci. **108** (1993) 385–391.
- [25] S. LEHR, J. SHALLIT, J. TROMP, *On the vector space of the automatic reals*. Theoretical Computer Science **163** (1996), 193–210.
- [26] J. H. LOXTON, *A method of Mahler in transcendence theory and some of its applications*. Bulletin of the Australian Mathematical Society **29** (1984), 127–136.
- [27] J. H. LOXTON, A. J. VAN DER POORTEN, *Arithmetic properties of certain functions in several variables III*. Bull. Austral. Math. Soc. **16** (1977), 15–47.
- [28] G. MARTIN, *Absolutely Abnormal Numbers*. American Mathematical Monthly **108** no. **8** (2001), 746–754.
- [29] W. MILLER, *Transcendence measures by a method of Mahler*. Journal of the Australian Mathematical Society (Series A) **32** (1982), 68–78.
- [30] YU. V. NESTERENKO, *Modular functions and transcendence questions*. Mat. Sb. **187** (1996), 65–96; translation in Sb. Math. **187** (1996), 1319–1348 [MR 97m:11102]
- [31] I. NIVEN, *Irrational Numbers*. Carus Mathematical Monographs, no. **11**, Wiley, New York, 1956.
- [32] P. RIBENBOIM, *The New Book of Prime Number Records*. Springer-Verlag, New York, 1996.
- [33] K. ROTH, *Rational Approximations to Algebraic Numbers*. Mathematika **2** (1955), 1–20. Corrigendum, pp. 168.
- [34] J. SAMBORSKI, *Problem E2667*. American Mathematical Monthly **84** (1977), pp. 567.
- [35] J. SHALLIT, private communication.
- [36] J. SHALLIT, A. VAN DER POORTEN, *A specialised continued fraction*. Canadian Journal of Mathematics **45** (1993), 1067–1079.

David H. BAILEY
Lawrence Berkeley National Laboratory
1 Cyclotron Road
Berkeley, CA 94720, USA
E-mail : dhbailey@lbl.gov

Jonathan M. BORWEIN
Dalhousie University
Department of Computer Science
Halifax, NS B3H 4R2, Canada
E-mail : jborrow@cs.dal.ca

Richard E. CRANDALL
Center for Advanced Computation
Reed College
Portland, OR 97202, USA
E-mail : crandall@reed.edu

Carl POMERANCE
Dartmouth College
Department of Mathematics
6188 Bradley Hall
Hanover, NH 03755-3551, USA
E-mail : carl.pomerance@dartmouth.edu