

The latter implies that both $D - 3$ and $D - 4$ are early and so, by assumption, are differences. Therefore, by Theorem 2 there exists an integer b for which

$$\frac{b}{a} + \frac{1}{2} > D - 4 \quad \text{and} \quad \frac{b + 8}{a} + \frac{1}{2} < D + 1;$$

that is,

$$b - D + 8 - \frac{a}{2} < \frac{D}{a} < b - D + \frac{9}{2}a.$$

Comparing this with (4), we obtain

$$n - 8 < n - 4 - 2a < b - D < n - 7,$$

which is contrary to the fact that $b - D$ and n are integers.

ACKNOWLEDGMENT

The author would like to thank Gerald Bergum for his aid in increasing the readability of this paper.

REFERENCES

1. J. C. Butcher. "On a Conjecture Concerning a Set of Sequences Satisfying The Fibonacci Difference Equation." *The Fibonacci Quarterly* 16 (1978):81-83.
2. M. D. Hendy. "Stolarsky's Distribution of Positive Integers." *The Fibonacci Quarterly* 16 (1978):70-80.
3. V. E. Hoggatt, Jr. *Fibonacci and Lucas Numbers*. Boston: Houghton Mifflin, 1969. Pp. 34-35.
4. K. Stolarsky. "A Set of Generalized Fibonacci Sequences Such That Each Natural Number Belongs to Exactly One." *The Fibonacci Quarterly* 15 (1977): 224.

INITIAL DIGITS IN NUMBER THEORY

J. KNOPFMACHER

University of the Witwatersrand, Johannesburg, 2001, South Africa

INTRODUCTION

It has been observed empirically by various authors (cf. Raimi [5] and his references) that the numbers in "random" tables of physical or other data tend to begin with low digits more frequently than one might on first consideration expect. In fact, in place of the plausible-looking frequency of $1/9$, it is found that for the numbers with first significant digit equal to

$$a \in \{1, 2, \dots, 9\}$$

in any particular table the observed proportion is often approximately equal to

$$\log_{10} \left(1 + \frac{1}{a} \right).$$

A variety of explanations have been put forward for this surprising phenomenon.

Although more general cases have also been considered, most people might agree that it should suffice to consider only sets of positive integers, since empirical data are normally listed in terms of finite lists of numbers with finite decimal expansions (for which the signs or positions of decimal points are immaterial here). On accepting this simplification, the common tendency

would probably then be to seek an explanation in terms of the concept of *natural density* of a set T of positive integers, i.e.,

$$\lim_{x \rightarrow \infty} \frac{1}{x} \sum_{n \leq x, n \in T} 1.$$

Unfortunately, this density simply does not exist for the immediately relevant set $N(a)$ of all positive integers beginning with the digit a as above, and this fact seems to have led both to a search for alternative explanations and to a certain amount of controversy as to what should actually constitute a satisfactory "explanation." Ignoring the latter difficulty for the moment (regarding which some further comments are offered in Section 3 below), the situation may be summarized by noting that various explanations have been suggested in terms of extensions of the density concepts that do exist and take the experimentally observed value of $\log_{10}(1+1/a)$ for the set $N(a)$; the most general and convincing of such approaches is perhaps that of Cohen [1].

The main purpose of this note is to add to these explanations by showing that the same type of initial-digit phenomenon occurs in a variety of number-theoretical situations. A notable investigation of this phenomenon of specific number-theoretical interest is that of Whitney [7] regarding the set P of all *prime* numbers. Whitney employs perhaps the most commonly used extension of the density concept, *logarithmic* (or Dirichlet) density, and this will also be used below. His discussion uses a corollary of one of the deeper forms of the Prime Number Theorem.

Here, using only elementary methods, it will be shown that, for quite a wide class of sets T of positive integers possessing a natural density, the subset $T(a) = T \cap N(a)$ has the relative logarithmic density $\log_{10}(1+1/a)$ in T . More generally, for quite a wide class of arithmetical functions, f , the logarithmic average value of f over all positive integers compared with that over $N(a)$ is shown to be weighted in the ratio $1:\log_{10}(1+1/a)$. In the actual discussion below, 10 is replaced by an arbitrary base $q \geq 2$, and a is replaced by an arbitrary initial sequence a_1, a_2, \dots, a_r of digits $a_i \in \{0, 1, \dots, q-1\}$ with $a_1 \neq 0$.

1. LOGARITHMIC AVERAGES AND DENSITIES

In order to cover a variety of specific examples of arithmetical functions and sets of positive integers in a fairly wide setting, first consider any fixed integers $q \geq 2$ and

$$A = a_1 q^{r-1} + a_2 q^{r-2} + \dots + a_r,$$

with $a_i \in \{0, 1, \dots, q-1\}$ and $a_1 \neq 0$. Let $N(A)$ denote the set of all positive integers whose canonical q -adic expansions begin with the sequence of digits a_1, a_2, \dots, a_r . We first wish to present the following theorem.

Theorem 1.1: Let f denote a nonnegative, real-valued function of the positive integers such that

$$\sum_{n \leq x} f(n) = Bx^\delta + O(x^\eta) \text{ as } x \rightarrow \infty,$$

where B, δ , and η are constants with $0 < \delta, \eta < \delta$. Then

$$\lim_{x \rightarrow \infty} \frac{1}{\log x} \sum_{n \leq x, n \in N(A)} f(n) n^{-\delta} = \delta B \log_q \left(1 + \frac{1}{A} \right).$$

Before proving Theorem 1.1, we need the following lemma.

Lemma 1.2: Under the hypothesis of Theorem 1.1, there exists a constant $\gamma = \gamma_f$ such that

$$\sum_{n \leq x} f(n)n^{-\delta} = \delta B \log x + \gamma + O(x^{\delta-\eta}) \text{ as } x \rightarrow \infty.$$

Proof: This lemma is actually a special case of a result discussed in [3, p. 86]. However, for the reader's convenience, we outline a direct proof here.

Let

$$F(x) = \sum_{n \leq x} f(n).$$

Then by partial summation (cf. [2, Theorem 421]), one obtains

$$\begin{aligned} \sum_{n \leq x} f(n)n^{-\delta} &= F(x)x^{-\delta} + \delta \int_1^x F(t)t^{-\delta-1} dt \\ &= [Bx^\delta + O(x^\eta)]x^{-\delta} + \delta \int_1^x [Bt^\delta + O(t^\eta)]t^{-\delta-1} dt \\ &= B + \delta B \log x + I(x) + O(x^{\eta-\delta}), \end{aligned}$$

where

$$\begin{aligned} I(x) &= \delta \left(\int_1^\infty - \int_x^\infty \right) [F(t) - Bt^\delta] t^{-\delta-1} dt \\ &= I - O \left(\int_x^\infty t^{\eta-\delta-1} dt \right) = I - O(x^{\eta-\delta}), \end{aligned}$$

for some constant I . The lemma follows, with $\gamma = B + I$.

Proof of Theorem 1.1: In order to deduce Theorem 1.1, first consider

$$x_m = (A+1)q^m.$$

By Lemma 1.2 (using the convention that $Aq^0 - 1$ be replaced by 1 if $A = 1$), we have

$$\begin{aligned} \sum_{n < x_m, n \in N(A)} f(n)n^{-\delta} &= \sum_{t=0}^m \sum_{Aq^t \leq n < (A+1)q^t} f(n)n^{-\delta} \\ &= \sum_{t=0}^m \left\{ \delta B \log \frac{(A+1)q^t - 1}{Aq^t - 1} + O(q^{t(\eta-\delta)}) \right\}. \end{aligned}$$

Thus

$$\begin{aligned} \sum_{n < x_m, n \in N(A)} f(n)n^{-\delta} &= \delta B \sum_{t=0}^m \log \frac{(A+1)q^t - 1}{Aq^t - 1} + O(1) \\ &= \delta B \sum_{t=0}^m \left\{ \log \left(1 + \frac{1}{A} \right) + \log \left(1 - \frac{1}{(A+1)q^t} \right) \right. \\ &\quad \left. - \log \left(1 - \frac{1}{Aq^t} \right) \right\} + O(1) \\ &= (m+1)\delta B \log \left(1 + \frac{1}{A} \right) + O(1), \end{aligned}$$

since (for $c \geq 1$),

$$\sum_{t=1}^m \log \left(1 - \frac{1}{cq^t} \right) = \sum_{t=1}^m O(q^{-t}) = O(1).$$

Now let $x_{m-1} \leq x \leq x_m$. Then $\log x \sim m \log q$ as $m, x \rightarrow \infty$, and [for $g(n) \geq 0$],

$$\sum_{n < x_{m-1}, n \in N(A)} g(n) \leq \sum_{n \leq x, n \in N(A)} g(n) \leq \sum_{n < x_m, n \in N(A)} g(n).$$

The asymptotic formula implies that

$$\lim_{x \rightarrow \infty} \frac{1}{\log x} \sum_{n \leq x, n \in N(A)} f(n)n^{-\delta} = \delta B \frac{\log(1 + 1/A)}{\log q} = \delta B \log_q \left(1 + \frac{1}{A}\right),$$

and Theorem 1.1 is proved.

We say that a function f of positive integers has *mean-value* (respectively, *logarithmic mean-value*) α over a set T of positive integers if and only if

$$\alpha = \lim_{x \rightarrow \infty} \frac{1}{x} \sum_{n \leq x, n \in T} f(n)$$

(respectively, $\alpha = \lim_{x \rightarrow \infty} \frac{1}{\log x} \sum_{n \leq x, n \in T} \frac{f(n)}{n}$).

It is shown by Wintner [8, p. 52] that the existence of the logarithmic mean-value over a set T follows from that of the mean-value over T and the values are equal. The converse is false. Applying Theorem 1.1, we have

Corollary 1.3: Let f denote a nonnegative, real-valued function that possesses the mean-value B over all positive integers in the strong sense that there exists a constant $\eta < 1$ such that

$$\sum_{n \leq x} f(n) = Bx + O(x^\eta) \text{ as } x \rightarrow \infty.$$

Then f possesses the logarithmic mean-value $B \log_q(1 + 1/A)$ over $N(A)$.

A subset S of a set T of positive integers is said to have the relative *logarithmic density* Δ in T if and only if

$$\Delta = \lim_{x \rightarrow \infty} \left(\sum_{n \leq x, n \in S} \frac{1}{n} \right) / \left(\sum_{n \leq x, n \in T} \frac{1}{n} \right).$$

If the function f of Corollary 1.3 is replaced by the characteristic function of the set T in the set N of all natural numbers, we obtain

Corollary 1.4: Let T denote a set of positive integers having natural density B in the strong sense that there exists a constant $\eta < 1$ such that

$$\sum_{n \leq x, n \in T} 1 = Bx + O(x^\eta) \text{ as } x \rightarrow \infty.$$

Then the set $T(A) = T \cap N(A)$ has the relative logarithmic density $\log_q(1 + 1/A)$ in T .

2. APPLICATIONS TO SPECIFIC SETS AND FUNCTIONS

In addition to the set N of all natural numbers, the following natural examples of sets T satisfying the hypothesis and hence the conclusion of Corollary 1.4 may be noted:

(2.1) Let $T_{m,r}$ denote the arithmetical progression

$$r, r + m, r + 2m, \dots \quad (0 \leq r < m).$$

Then clearly

$$\sum_{n \leq x, n \in T_{m,r}} 1 = \sum_{r + km \leq x} 1 = \left[\frac{x - r}{m} \right] = \frac{x}{m} + O(1) \text{ as } x \rightarrow \infty.$$

(2.2) Given any integer $k \geq 2$, let $N_{[k]}$ denote the set of all k -free positive

integers, i.e., integers not divisible by any k th power $r^k \neq 1$. (Thus $N_{[2]}$ is the familiar set of all *square-free* numbers.) Then it is known (see, e.g., [3, p. 108]) that

$$\sum_{n \leq x, n \in N_{[k]}} 1 = \frac{x}{\zeta(k)} + O(x^{1/k}) \text{ as } x \rightarrow \infty,$$

where

$$\zeta(k) = \sum_{n=1}^{\infty} n^{-k}.$$

(2.3) Let $T_{m,r,k} = T_{m,r} \cap N_{[k]}$, where $T_{m,r}$ and $N_{[k]}$ are the sets defined above. If m, r are coprime, it is known (see, e.g., [4, p. 112]) that as $x \rightarrow \infty$,

$$\sum_{n \leq x, n \in T_{m,r,k}} 1 = \frac{x}{m\zeta(k)} \prod_{\text{prime } p|m} (1 - p^{-k})^{-1} + O(x^{1/k}),$$

where $\zeta(k)$ is as before.

Many naturally occurring arithmetical functions f satisfy the hypothesis and hence the conclusion of Corollary 1.3. Out of examples of such functions treated in books, we mention only two:

(2.4) Let $r(n)$ denote the number of lattice points (a, b) such that $a^2 + b^2 = n$. Then (see, e.g., [2, Theorem 339]),

$$\sum_{n \leq x} r(n) = \pi x + O(x^{1/2}) \text{ as } x \rightarrow \infty.$$

(2.5) Let $\alpha(n)$ denote the total number of nonisomorphic abelian groups of finite order n . A theorem of Erdős and Szekeres (see, e.g., [3, p. 117]) states that

$$\sum_{n \leq x} \alpha(n) = x \prod_{k=2}^{\infty} \zeta(k) + O(x^{1/2}) \text{ as } x \rightarrow \infty.$$

Next we mention a few examples of concrete arithmetical functions f satisfying the slightly more general hypothesis of Theorem 1.1:

(2.6) The Euler function

$$\phi(n) = \sum_{r \leq n, n(r,n)=1} 1$$

has the property that

$$\sum_{n \leq x} \phi(n) = \frac{3x^2}{\pi^2} + O(x \log x) \text{ as } x \rightarrow \infty$$

(see, e.g., [2, Theorem 330]).

(2.7) The divisor-sum function

$$\sigma(n) = \sum_{d|n} d$$

has the property that

$$\sum_{n \leq x} \sigma(n) = \frac{1}{12} \pi^2 x^2 + O(x \log x) \text{ as } x \rightarrow \infty$$

(cf. [2, Theorem 324]).

(2.8) Given any positive integer k , let $T_{m,r}^k$ denote the set of all k th powers of numbers in the arithmetical progression $T_{m,r}$ of (2.1). Then,

$$\sum_{n \leq x, n \in T_{m,r}^k} 1 = \sum_{n \leq x^{1/k}, n \in T_{m,r}} 1 = \frac{1}{m} x^{1/k} + O(1),$$

by (2.1). Thus Theorem 1.1 applies to the characteristic function of the set $T_{m,r}^k$ in N .

Finally, it may be remarked that the applicability of Theorem 1.1 carries over to the restrictions to $T_{m,r}$ of arithmetical functions of the above kinds, when m, r are coprime. (For preliminary theorems that make such applications possible, see, e.g., [3, Ch. 9], [4, Ch. II], and Smith [6].)

3. "SCIENTIFIC" VERSUS MATHEMATICAL EXPLANATIONS

In [5] Raimi expresses some reservations about purely mathematical explanations of the initial-digit phenomenon in numerical tables of empirical data and calls for a more "scientific" discussion (e.g., in terms of statistical distribution functions). However, in this direction, general agreement does not seem to have been reached or even to be imminent. By way of contrast, even if it is theoretically correct to have done so, one might query whether such a problem would ever have been seriously raised in practice if it had not been for the nonexistence of certain desired natural densities.

For, suppose that a detailed examination of "random" tables of numerical data was found to show that, in most cases, approximately 1/10 of the numbers considered *end* in a particular digit $b \in \{0, 1, \dots, 9\}$, or approximately 10 of them end in a particular sequence of digits $b_1, b_2, \dots, b \in \{0, 1, \dots, 9\}$. In view of the elementary example (2.1) above, surely very few people would be surprised by this or be led to call seriously for a "scientific" explanation, even though it is theoretically as legitimate to do so here as in the original problem.

Although the nonexistence of natural densities does on first consideration seem to lend an element of confusion to the initial-digit problem, the preceding remarks suggest that (unless overwhelming experimental evidence* warrants otherwise) it is perhaps nevertheless adequate for most purposes to accept an explanation in terms of one or more reasonable mathematical substitutes for natural density. In showing the quite widespread nature of this phenomenon in number theory, the earlier theorem and various mathematical examples perhaps lend further weight to this suggestion. After all, what can be scientifically interesting about the purely *numerological* properties of a list of street addresses, or areas of rivers, and so on?

REFERENCES

1. D. I. A. Cohen. "An Explanation of the First Digit Phenomenon." *J. Combin. Theory* A20 (1976):367-379.
2. G. H. Hardy & E. M. Wright. *An Introduction to the Theory of Numbers*. Oxford University Press, 1960.
3. J. Knopfmacher. *Abstract Analytic Number Theory*. North-Holland Publishing Company, 1975.
4. J. Knopfmacher. "Arithmetical Properties of Finite Rings and Algebras, and Analytic Number Theory, V." *J. reine angew. Math.* 271 (1974):95-121.
5. R. A. Raimi. "The First Digit Problem." *Amer. Math. Monthly* 83 (1976):521-538.
6. R. A. Smith. "The Circle Problem in an Arithmetical Progression." *Canad. Math. Bull.* 11 (1968):175-184.
7. R. E. Whitney. "Initial Digits for the Sequence of Primes." *Amer. Math. Monthly* 79 (1972):150-152.
8. A. Wintner. *The Theory of Measure in Arithmetical Semigroups*. The Waverly Press, 1944.

*The status of Raimi's anomalous population data $PP(n)$ [5, p. 522] is difficult to evaluate without further investigation, but his anomalous data $V(n)$ do not seem surprising if one remembers that telephone numbers normally have favoured initial digits.