

COMPLEXITY OF NATURAL NUMBERS

JUAN ARIAS DE REYNA

English version of the paper: Complejidad de los números naturales, Gaceta de la Real Sociedad Matemática Española 3 (2000) 230–250.

1. INTRODUCTION.

1.1. Complexity of a natural number. Our purpose is to explore what seems to be a trivial question that may be understood by high-school students in their early teens, but with very deep relationships.

Lately I have been interested in one of the mathematical problems that I consider most important: the $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ problem. In this case the first difficulty is to explain the problem to a professional mathematician, say to an expert in Analysis. This is not a minor issue, I think that the problem $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ may be put as an inequality. Hence to explain the question adequately, so that it is understood by an expert in Analysis, maybe the first step in the solution of the problem.

The question I shall discuss here arose while trying to obtain this explanation.

We start with the main question: Given a natural number n , how many 1's are needed to write n ? For example

$$19 = 1 + (1 + 1)(1 + 1 + 1)(1 + 1 + 1)$$

so that nine 1's suffice to write 19. We shall say that the complexity of 19 is less than or equal to 9, and we shall write this as $\|19\| \leq 9$. Of course, the complexity of 19 will be the number of 1's in the most economical representation of 19. We only admit expressions with sums and products.

The first values of the complexity function may be easily computed

$$1, 2, 3, 4, 5, 5, 6, 6, 6, 7, 8, 7, 8, 8, 8, 8, 9, 8, 9, 9, \dots$$

We see that this is not a monotonic sequence: $8 = \|11\| > \|12\| = 7$.

When in our investigations we find any sequence of natural numbers, there is something we must do: look in *The On-Line Encyclopedia of Integer Sequences* of Sloane and Plouffe [8]. In it we find this sequence and a reference to a paper by Guy [4] where it is defined and analyzed.

Date: May, 2000.

2. COMPLEXITY OF A NATURAL NUMBER.

We have defined the complexity as a function $n \mapsto \|n\|$ of $\mathbb{N} \rightarrow \mathbb{N}$ such that for every pair of natural numbers m and n we have

$$\|1\| = 1, \quad \|m + n\| \leq \|m\| + \|n\|, \quad \|m \cdot n\| \leq \|m\| + \|n\|.$$

In fact it is the largest function satisfying these conditions. To prove this and other assertions it is useful to introduce the concept of *expression*.

2.1. Definition of expression. An expression is a sequence of symbols. The allowed symbols are x , $+$, $($, $)$. Not every sequence of these symbols is an expression. Examples of expressions are:

$$(x + x); \quad (x+(xx)); \quad (x+((x+x)((x+(x+x))(x+(x+x))))).$$

The formal definition is inductive:

- (a) x is an expression.
- (b) If A and B are expressions, then $(A+B)$ and (AB) are also expressions.
- (c) The only expressions are those obtained by repeated applications of rules (a) and (b).

We define the value of an expression A as the number $v(A)$ that results when replacing x by 1. Again we use induction to define the value of an expression: $v(x) = 1$, and if A and B are expressions then $v((A+B)) = v(A) + v(B)$ and $v((AB)) = v(A)v(B)$.

Given an expression we may define its complexity as the number of letters x it contains, for example $\|(x+(xx))\| = 3$. Let \mathcal{E} be the set of expressions. We may translate the definition of the complexity as

$$\|n\| = \inf\{\|A\| : A \in \mathcal{E} \text{ and } v(A) = n\}.$$

If we want to compute the value of $\|n\|$ we may use the following Proposition.

Proposition 1. *For each natural number $n > 1$*

$$\|n\| = \min_{\substack{2 \leq d \leq \sqrt{n}, d|n \\ 1 \leq j \leq n/2}} \{\|d\| + \|n/d\|, \quad \|j\| + \|n - j\|\}$$

Proof. Let E an optimal expression for n , i. e. one that gives its complexity $\|n\| = \|E\|$. As an expression that is not x we will have $E=(A+B)$ or $E=(AB)$. Let $a = v(A)$ and $b = v(B)$. Then either $n = a + b$ and $\|n\| = \|a\| + \|b\|$ or $n = ab$ and $\|n\| = \|a\| + \|b\|$. In the first case if j is the least of a and b we will have $1 \leq j \leq n/2$, and in the second case if d is the least of a and b , then d will be a divisor of n with $2 \leq d \leq \sqrt{n}$. Of course for the reasoning to be valid we must check that if E is an

optimal expression for n , then A and B must be optimal expressions for a and b respectively. We leave this check to the reader. \square

Using the above Proposition and the mathematical software Mathematica we have computed the values of $\|n\|$ for $1 \leq n \leq 200\,000$.

3. BOUNDS.

Proposition 2. *Let $P: \mathbb{N} \rightarrow \mathbb{R}$ be a function satisfying*

$$P(1) = 1, \quad P(n + m) \leq P(n) + P(m), \quad P(n \cdot m) \leq P(n) + P(m).$$

Then for each $n \in \mathbb{N}$ we have $P(n) \leq \|n\|$.

Proof. It is easy to see by induction that for each expression A , we have $P(v(A)) \leq \|A\|$. It is true for $A = x$, and, if it is true for A and B then it is true for $(A+B)$ and (AB) . For example, for the product:

$$P(v((AB))) = P(v(A)v(B)) \leq P(v(A)) + P(v(B)) \leq \|A\| + \|B\| = \|(AB)\|,$$

and a similar argument is valid for the sum. (Observe that by the definition of v we have $v((A+B)) = v(A) + v(B)$ and $v((AB)) = v(A)v(B)$).

Now in $P(v(A)) \leq \|A\|$ we take the minimum over all expressions A such that $n = v(A)$. In this way we get $P(n) \leq \|n\|$. \square

Corollary 3. *For each natural number n we have $\log_2(1 + n) \leq \|n\|$.*

Proof. It is sufficient to check the properties of $P(n) = \log_2(1 + n)$. \square

Later, in Corollary 9, we will obtain a better inequality.

3.1. Upper bounds. Now we get an upper bound. To this end we define a new function $L: \mathbb{N} \rightarrow \mathbb{N}$.

Definition 4. We define the function L inductively:

- (a) $L(1) = 1$.
- (b) If p is a prime number, then $L(p) = 1 + L(p - 1)$.
- (c) If $n = p_1 p_2 \cdots p_k$ is a product of prime numbers (may be repeated), then $L(p_1 p_2 \cdots p_k) = L(p_1) + L(p_2) + \cdots + L(p_k)$.

It is clear from this definition that if $n = ab$ with a and $b \geq 2$ then we will have $L(ab) = L(a) + L(b)$.

Proposition 5. *For each $n \in \mathbb{N}$ we have*

$$\|n\| \leq L(n).$$

Proof. We may prove this by induction. For $n = 1$ we have $\|1\| = L(1) = 1$. Assume that $\|k\| \leq L(k)$ for each $k < n$. There are two possibilities: if $n = p$ is a prime number

$$\|p\| \leq \|p-1\| + \|1\| = \|p-1\| + 1 \leq L(p-1) + 1 = L(p).$$

If n is composite $n = ab$ with a and $b > 2$,

$$\|n\| \leq \|a\| + \|b\| \leq L(a) + L(b) = L(ab) = L(n).$$

□

Proposition 6. *For each $n \geq 2$ we have*

$$L(n) \leq \frac{3}{\log 2}(\log n).$$

Proof. Since $L(2) = 2$ and $L(3) = 3$ the result is true for $n = 2$ and $n = 3$.

Assume now that $n > 3$ and that the Proposition is true for all natural numbers strictly less than n .

If $n = p$ is a prime number we have

$$(1) \quad L(p) = 1 + L(p-1) = 1 + 2 + L\left(\frac{p-1}{2}\right) \leq 3 + \frac{3}{\log 2} \log\left(\frac{p-1}{2}\right).$$

We want this to be

$$\leq \frac{3}{\log 2}(\log p).$$

Hence we must check that

$$(2) \quad 3 \leq \frac{3}{\log 2} \log\left(\frac{2p}{p-1}\right),$$

which is easily proved for $p \geq 3$.

If $n = ab$ with a and $b \geq 2$, we have

$$L(ab) = L(a) + L(b) \leq \frac{3}{\log 2}(\log a) + \frac{3}{\log 2}(\log b) = \frac{3}{\log 2}(\log ab).$$

□

Remark 1. We do not know if the constant $3/\log 2$ in the above theorem is optimal. The proof makes one suspect that the quotient $L(n)/\log n$ may be large when $n = p_k$ is a prime such that there exists a sequence of primes $(p_j)_{j=1}^k$ with $p_{j+1} = 2p_j + 1$. For example, 89, 179, 359, 719, 1439, 2879 is such a sequence of prime numbers, and the maximum value of the quotient $L(n)/\log n$ that we know is

$$\frac{L(2879)}{\log 2879} = 3.766384578 \dots < 4.328085123 \dots = \frac{3}{\log 2}.$$

The main difference between the two functions $L(\cdot)$ and $\|\cdot\|$ is that $L(\cdot)$ is additive and $\|\cdot\|$ is not. For each pair of numbers n and m greater than 1 we have $L(mn) = L(m) + L(n)$. On the other hand there exist pairs n, m of numbers greater than 1 and such that $\|mn\| < \|m\| + \|n\|$. In such a case we shall say that $n \cdot m$ is a bad factorization.

In figure 1 we put a dot at each point (n, m) such that $n \cdot m$ is a bad factorization. The figure contains all the factors n and $m \leq 60$.

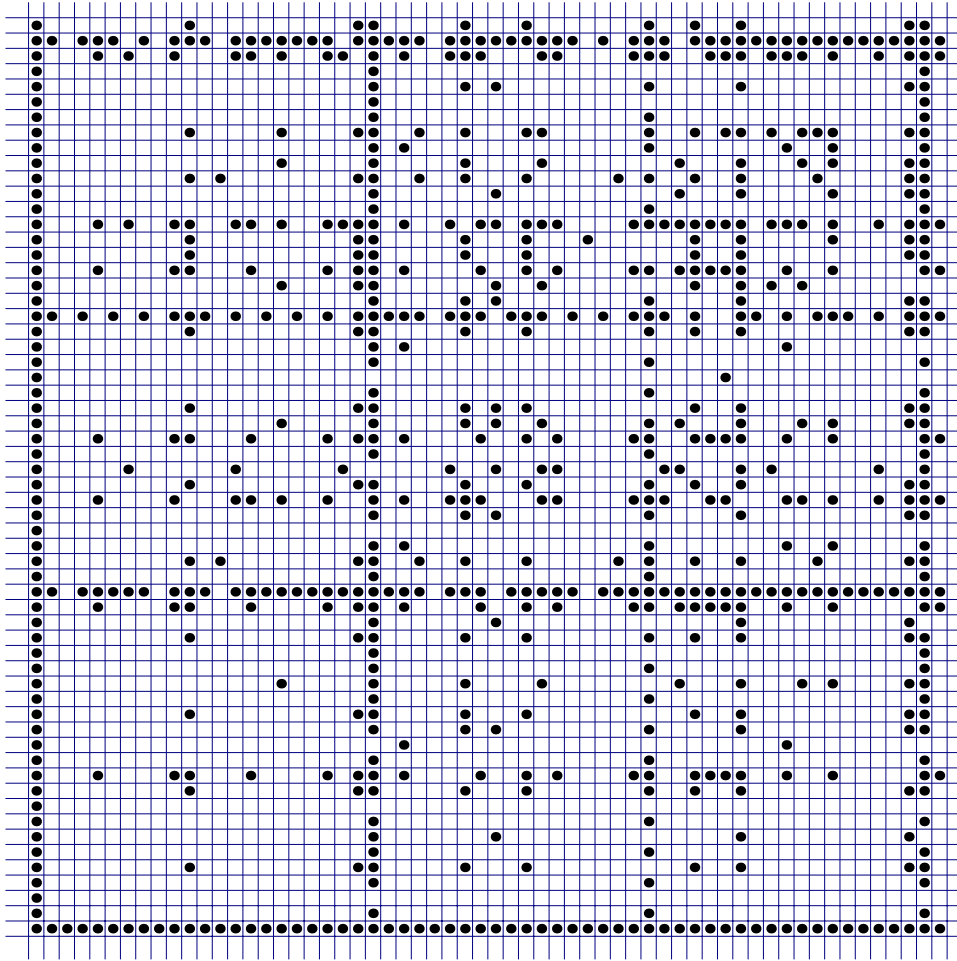


FIGURE 1. Bad Factors.

$1 \cdot n$ is always a bad factorization. In the figure we see some other surprising regularities. There are some conspicuous (vertical and horizontal) aligned points. Especially note the verticals at $n = 23, 41, 59$, which deserve an explanation.

These numbers, we may call them *bad factors*, appear to have great complexity. We define the *number with great complexity* n_k as the number n_k that is the less solution to $\|n\| = k$. The first values of this sequence are

1, 2, 3, 4, 5, 7, 10, 11, 17, 22, 23, 41, 47, 59,
 89, 107, 167, 179, 263, 347, 467, 683, 719, 1223,
 1438, 1439, 2879, 3767, 4283, 6299, 10079, 11807,
 15287, 21599, 33599, ...

This sequence appears in [8] with some errata. In this way we find the reference to Rawsthorne [7].

4. MEAN VALUES.

There is another proof of $\|n\| \leq 3 \log n / \log 2$. We observe that if we write n in binary $n = \sum_{j=0}^{k-1} \varepsilon_j 2^j + 2^k$ we have a means to express n :

$$n = \varepsilon_0 + 2(\varepsilon_1 + 2(\varepsilon_2 + \cdots + 2(\varepsilon_{k-2} + 2(\varepsilon_{k-1} + 2)) \cdots)).$$

If we substitute each 2 by 1+1 and observe that each ε_j is equal to 0 or 1, we have an expression for n that uses at most $2k + k$ ones, and where k is determined by $2^k \leq n < 2^{k+1}$. It follows that $\|n\| \leq 3 \log n / \log 2$.

The above reasoning proves that the function $L_2(n) = 2k + \varepsilon_0 + \varepsilon_1 + \cdots + \varepsilon_{k-1}$ is another upper bound for $\|n\|$. The relation between $L_2(n)$ and $L(n)$ is not very simple. Amongst the first 1000 numbers we generally have $L(n) \leq L_2(n)$ but this inequality has exceptions. The first one is $L_2(161) = 16 < 17 = L(161)$. In this range the difference is small.

The function $L_2(n)$ allows us to obtain information about $\|\cdot\|$. Consider the numbers n that in binary take the form $1\varepsilon_{k-1} \cdots \varepsilon_0$, i. e. numbers than in binary have $k + 1$ digits. By the above expression we have

$$\|n\| \leq 2k + \varepsilon_0 + \cdots + \varepsilon_{k-1}.$$

We may suppose that the ε_k are independent random variables with mean 1/2. The inequality of Chernoff (see [2] or [1] for a simple exposition) says that

$$\mathbb{P}\left(\left|\sum \varepsilon_j - k/2\right| < x\sqrt{k}\right) \geq 1 - 2e^{-2x^2}.$$

It follows that $\mathbb{P}(\|n\| \leq 2k + k/2 + x\sqrt{k}) \geq 1 - 2e^{-2x^2}$, and taking $x = \sqrt{\log k}$ we get

$$\mathbb{P}\left(\|n\| > 5k/2 + \sqrt{k \log k}\right) \leq 2k^{-2}.$$

Hence between the 2^k values of n with $2^k \leq n < 2^{k+1}$ at most $(2/k^2)2^k$ satisfy $\|n\| > 5k/2 + \sqrt{k \log k}$. The other ones, most of them, satisfy

$$\|n\| \leq \frac{5k}{2} + \sqrt{k \log k} = \frac{5 \log n}{2 \log 2} + O(\sqrt{\log n \log \log n}).$$

Therefore, for almost all large values of n we have

$$\|n\| \leq \frac{5 \log n}{2 \log 2} + O(\sqrt{\log n \log \log n}).$$

The upper bound $L(n)$ is very good for small values of n . For example for the first 220 values of n , $L(n) = \|n\|$, except for the values in the following table:

n	$\ n\ $	$L(n)$	n	$\ n\ $	$L(n)$	n	$\ n\ $	$L(n)$
46	12	13	115	15	16	164	15	16
47	13	14	118	15	16	165	15	16
55	12	13	121	15	16	166	16	17
82	13	14	138	15	16	167	17	18
83	14	15	139	16	17	184	16	17
92	14	15	141	16	17	188	17	18
94	15	16	145	15	16	217	16	17
110	14	15	161	16	17	220	16	17

In these cases the bound $L_2(n)$ is equal or greater than $L(n)$, except for the case $n = 161$.

The two functions $L(n)$ and $\|n\|$ coincide in 771 values of n in the range $1 \leq n \leq 1000$, the difference being equal to 1 for the 229 other values in this range with a few exceptions.

5. PARTICULAR VALUES.

5.1. Numbers with small complexity. A good lower bound for $\|n\|$ is obtained from the knowledge of the largest number we may write with m ones. That is, given m , which is the largest natural number N with $\|N\| = m$. The answer roughly is that we must group the m ones in groups of three and multiply them. To show this we define the concept of *extremal expression*. Let \mathbf{M}_m be an expression with $\|\mathbf{M}_m\| = m$ (that is \mathbf{M}_m is formed with m symbols \mathbf{x} and the operations of sum and product), and such that its value $v(\mathbf{M}_m)$ is the maximum of all the expression with m ones, i. e.

$$N = v(\mathbf{M}_m) = \sup_{\|\mathbf{A}\|=m} v(\mathbf{A}).$$

We say that such an expression \mathbf{M}_m is extremal.

In the above situation $\|N\| = m$. In fact, since $N = v(M_m)$ and $\|M_m\| = m$, we have $\|N\| \leq m$. Assume, by contradiction, that $\|N\| < m$. Then there will exist an expression B such that $v(B) = N$ and $\|B\| = \|N\| < m$. Let d be such that $m = d + \|B\|$. We may construct an expression C such that $C = B + x + \dots + x$ and such that $\|C\| = \|B\| + d = m$ and $v(C) = v(B) + d > N$. This contradicts the definition of M_m .

It is easy to see that the following expressions are extremal

$$\begin{aligned} M_1 &= x, & M_2 &= (x + x), & M_3 &= (x + (x+x)), \\ M_4 &= (x+x)(x+x), & M_5 &= (x+(x+x))(x+x), \dots \end{aligned}$$

We see that given m the extremal expression M_m is not unique. For example for $m = 4$ the expression $M_4 = (x+(x+(x+x)))$ is another possibility.

We shall use here a not very precise notation. For example, we shall write $M_3^a M_2$ to denote any expression having this form, not defining how the product is constructed from its factors. So, M_3^4 denotes any of the expressions $((M_3 M_3)(M_3 M_3))$, $(M_3(M_3(M_3 M_3)))$ or any other form of grouping the factors.

Proposition 7. *Let $M_2 = (x + x)$, $M_3 = (x + (x+x))$ and $M_4 = (x+x)(x+x)$. For $n > 1$, the expressions*

$$M_n = \begin{cases} M_3^k & \text{if } n = 3k, \\ M_3^{k-1} M_4 & \text{if } n = 3k + 1, \\ M_3^k M_2 & \text{if } n = 3k + 2, \end{cases}$$

are extremal.

Proof. We may check the proposition for $n = 2, 3$ and 4 directly.

Assume the assertion for all $s < n$ and try to prove it for $n \geq 5$. Certainly there is one extremal expression K with $\|K\| = n$. Then there are two expressions A and B such that $K = (A + B)$ or $K = (AB)$. A and B are extremal expressions because K is extremal. We may replace A and B by extremal expressions of the same complexity and value and the resulting expression K' will be also extremal. Hence, without loss of generality, we may assume, using the induction hypothesis, that A and B are of the form given in the Proposition or $A = x$ and B is as in the Proposition.

The case $K = (A+B)$ it is only possible if $v(A)$ or $v(B) = 1$, because, in other cases, the expression (AB) contradicts the extremality of K . But $K = (x + M_3^k)$, $K = (x + M_3^{k-1} M_4)$, or $K = (x + M_3^k M_2)$ are impossible with $n \geq 5$. Because these expressions are clearly not extremal. (Compare with $M_3^{k-1} M_4$, $M_3^k M_2$ or M_3^{k+1} respectively).

Therefore $K = (AB)$ where A and B are like those in the Proposition. Some of the combinations are not possible: for example $A = M_3^k M_2$ and $B = M_3^{j-1} M_4$ are not possible since $M_3^{k+j-1} M_4 M_2$ is improved by M_3^{k+j+1} and K will not be extremal. A case by case analysis proves that K is one of the three forms in the Proposition. \square

Corollary 8. *For $a = 0, 1$, or 2 and $b \in \mathbb{N}$ we have:*

$$\|2^a 3^b\| = 2a + 3b, \quad a = 0, 1, 2.$$

All natural numbers $n > 1$ may be written in a unique way as $n = 2a + 3b$ with $a = 0, 1$ or 2 . In this case $2^a 3^b$ is the greatest number m with $\|m\| = n$. Hence $m > 2^a 3^b$ implies $\|m\| > 2a + 3b$.

We define g by

$$g(n) = \begin{cases} 3a & \text{if } n \in [3^a, 3^a + 3^{a-1}), \\ 3a + 1 & \text{if } n \in [3^a + 3^{a-1}, 2 \cdot 3^a), \\ 3a + 2 & \text{if } n \in [2 \cdot 3^a, 3^{a+1}), \end{cases}$$

we then have $g(n) \leq \|n\|$ for each n .

Corollary 9. *For any $n \geq 2$ we have*

$$3 \frac{\log n}{\log 3} \leq \|n\| \leq L(n) \leq 3 \frac{\log n}{\log 2}.$$

Proof. We only need to prove the first inequality. If $n = 3^a$, we see directly that the inequality is true. If $x \in (3^a, 3^a + 3^{a-1}]$, we have $\|x\| \geq 3a + 1$. Then

$$\|x\| \geq \|3^a\| + 1 = 3a + 1 \geq 3 \frac{\log(4 \cdot 3^{a-1})}{\log 3} \geq 3 \frac{\log x}{\log 3}.$$

Analogously for $x \in (4 \cdot 3^{a-1}, 2 \cdot 3^a]$ we have

$$\|x\| \geq \|4 \cdot 3^{a-1}\| + 1 \geq 3a + 2 \geq 3 \frac{\log(2 \cdot 3^a)}{\log 3}.$$

Finally for $x \in (2 \cdot 3^a, 3^{a+1}]$, we only need to check that

$$\|x\| \geq \|2 \cdot 3^a\| + 1 = 3a + 3 \geq 3 \frac{\log(3^{a+1})}{\log 3}.$$

\square

6. THE PROBLEM $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ AND THE COMPLEXITY OF THE
NATURAL NUMBERS.

6.1. **Idea of the problem $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$.** Before explaining the problem we must describe the classes \mathbf{P} and \mathbf{NP} . Consider a finite alphabet A , and let A^* be the set of *words*, that is, the set of finite sequences of elements of A .

We call *language* a subset $S \subset A^*$. We say that S is in the class \mathbf{P} if there is an algorithm T and a polynomial $p(t)$ such that with a word x as input, T gives an output $T(x)$, such that $T(x) = 1$ if $x \in S$ and $T(x) = 0$ if $x \notin S$. Also T gives the output $T(x)$ in a time bounded by $p(|x|)$ (here $|x|$ denotes the length of the word x). We then say that T is a polynomial algorithm. In a few words we may say that \mathbf{P} is the class of languages recognizable in polynomial time. It is important to notice that this concept is very stable with respect to the diverse definitions of what is an algorithm, how we compute the “time” that the algorithm T takes to give the output, or even if we consider the same language in a different alphabet (as when we consider a set of natural numbers written in different basis). In other words, the concept does not change if we give proper definitions of these concepts.

The class \mathbf{NP} consists of the languages recognizable by non deterministic polynomial algorithms. That is $S \subset A^*$ is in \mathbf{NP} if there exists an algorithm T and a polynomial $p(x)$ such that for each $x \in S$ there is $y \in A^*$ with $|y| \leq p(|x|)$ and such that with the input (x, y) the algorithm gives the output $T(x, y) = 1$ in time bounded by $p(|x|)$. On the other hand if $x \notin S$ we have $T(x, y) = 0$ for all y with $|y| \leq p(|x|)$.

We say that in this case T is a non-deterministic algorithm since to obtain $x \in S$ we must first choose y . If we know which y to take this process is fast, but if we do not know y , we may try each possible y , but this will need a time $\geq |A|^{p(|x|)}$ which in practice is impossible.

Again the class \mathbf{NP} is very stable with respect to possible changes in the definitions. Also many practical problems are in this class.

It is easy to check that $\mathbf{P} \subset \mathbf{NP}$. The question is whether these two classes are the same. To understand a bit more of the difficulty observe the following.

Our experience as mathematicians teaches us that to understand a proof, or better to check the correctness of a proof is a task of type \mathbf{P} . That is the time needed is proportional to the length of the proof.

On the other hand to determine if a conjecture x is a Theorem we need first to write the proof y and then apply the above procedure to check the correctness of the pair (x, y) . The set of Theorems is not in the class \mathbf{NP} since as we know the length of the proof $|y|$ is not

bounded by the length of the theorem x , that is $|y| \preceq p(|x|)$. But for each polynomial $p(t)$, the following set is in **NP**

$$\mathcal{T}_p = \{x : x \text{ is a theorem with a proof of length } \leq p(|x|)\}.$$

Maybe someone finds these definitions rather vague, but the formal logic allows one to make things precise.

If $\mathbf{P} = \mathbf{NP}$ and the proof were sufficiently constructive (technically, that we can find a polynomial algorithm for an **NP**-complete problem), then there would exist a polynomial algorithm that would allow not only decide if $x \in \mathcal{T}_p$, but also to find in this case a proof for x in polynomial time. The mathematicians would not be needed any more.

When one recalls the achievements of the 20th century: proof of Fermat's theorem, classification of finite simple groups, pointwise convergence of Fourier series of function in L^p , Riemann's hypothesis for algebraic varieties over fields of characteristic p , independence of continuum hypothesis, and many more, one gets the impression that there exists an algorithm to decide $x \in \mathcal{T}_p$, by searching directly for a proof, not by trial and error. This algorithm consists in taking promising students, give them the possibility to travel and speak with specialists on the topic in question, let them try to solve analogous questions, study the solution of related problems, and so on ...

7. CONNECTION OF THE COMPLEXITY OF NATURAL NUMBERS AND THE PROBLEM $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$.

Consider the assertion $\|4787\| = 28$. We may decompose it in two parts. The first, $\|4787\| \leq 28$, has a very easy proof

$$(\star) \quad 4787 = 2 + 3(2 + 3^2)(1 + 2^4 3^2).$$

The other part of the assertion $\|4787\| \geq 28$, has a much more laborious proof. Just now I do not know any other way than computing the values of $\|n\|$ for all $n \leq 4787$, a task that, on my personal computer, took several hours.

Of course this does not imply that it is easy to find proof as in (\star) .

Consider the sets

$$A = \{(n, c) \in \mathbb{N}^2 : \|n\| \leq c\}, \quad B = \{(n, c) \in \mathbb{N}^2 : \|n\| > c\}.$$

The fact, as we have remarked, that if $(n, c) \in A$, then there is a relatively short proof of it, shows us that A is in the class **NP**.

Roughly, a set A is in **NP**, if to prove that $x \in A$ an exhaustive search is required, which in principle is exponential in the size of x , but once the proof has been found, it is easily recognized (in polynomial time with respect to the size of x). Complete information may be found in

the book [3]. These problems bring to mind the one of finding a needle in a haystack. Once we have found the needle there is no doubt that the task is done, but at first it appears unreachable since the straw is so similar to the needle that we do not see any other means than search methodically.

The core of the problem $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ is whether in situations where there exists a short proof, there is always a direct path to find it. If $\mathbf{P} = \mathbf{NP}$, then there is always a direct path to the proof without hesitations. At first sight this appears a wild assumption, but the rigorous proof of $\mathbf{P} \neq \mathbf{NP}$ eludes us still after twenty seven years of study.

Recently Microsoft has funded an investigation center and has contracted Michael Friedman, (Fields medal in 1986). Friedman has the intention of trying to solve the question $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$. Microsoft will invest 2.6 million dollars each year in this program.

It appears that $\mathbf{P} = \mathbf{NP}$ is false, but not all is so simple. Sometimes tasks that appear to need an exhaustive search have been proved simple. We shall give an example.

Let $\mathcal{C} \subset \mathbb{N}$ the set of composite numbers. At first sight it appears that the only means to proof that n is composite is to divide n by each number $m \leq \sqrt{n}$ and check if some remainder equals 0. The size of n is of the order of the number of digits needed to write it, i.e. of the order $\log n$. The number of needed checks maybe $\sqrt{n} = e^{(\log n)/2}$, which grows exponentially with $\log n$. And if really n is composite there is a short proof: to exhibit a proper divisor d of n . That is \mathcal{C} is in the class \mathbf{NP} .

But it is not so difficult to decide whether n is composite. If n is prime and b is prime with n we have $b^{n-1} \equiv 1 \pmod{n}$. An idea somewhat more elaborate, let n be a prime and $n - 1 = 2^{st}$, in the sequence of the rests of $b^t, b^{2t}, \dots, b^{2^{st}} \pmod{n}$ the last different from 1 must be -1 . In the other case it is sure that n is composite. This is the famous Miller-Rabin test. It is known that if the generalized Riemann hypothesis is true, then if n is composite, the test of Miller-Rabin is not satisfied for some $b < 2(\log n)^2$. Hence, under the mentioned hypothesis, we have a fast algorithm (polynomial) to decide whether n is composite: to do the test of Miller-Rabin for all $b < 2(\log n)^2$.

Another incentive to pose the problem $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ is the existence of \mathbf{NP} -complete problems. That is sets $B \subset \mathbb{N}$ such that B is in the class \mathbf{NP} and, for which from $B \in \mathbf{P}$ it follows that $\mathbf{P} = \mathbf{NP}$.

From Euclid's times, mathematicians have had a clear concept of algorithm. Turing gives a further step and by an effort of introspection

gives us a precise definition. Turing's mental image is that of a mathematician, notebook in hand, computing. By abstracting the procedure Turing created the idea of a modern computer. Starting from Turing's definitions it is possible to quantify the time a computer will spend on a given task and so to give a precise definition of the classes **P** and **NP**.

The first connection of the complexity of the natural numbers with the problem $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ is the fact that $\mathbf{P} = \mathbf{NP}$ implies the existence of a fast algorithm to compute $\|n\|$. There will be constants C and $k \in \mathbb{N}$ and an algorithm that will compute $\|n\|$ in time $\leq C(\log n)^k$.

8. COMPLEXITY OF BOOLEAN FUNCTIONS.

There is another connection, this time structural, between the complexity of natural numbers and the problem $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$. To explain this connection we must define a related concept, that of the complexity of a boolean function.

The set $\{0, 1\}$ is a field when we consider the composition laws sum and product mod 2. For each number n let \mathcal{F}_n be the set of functions $f: \{0, 1\}^n \rightarrow \{0, 1\}$. The set \mathcal{F}_n is a ring if we take sum and product with respect to the field in the image $\{0, 1\}$.

For example consider the constant functions **1**, **0** and the components π_j defined by $\pi_j(\mathbf{x}) = \pi_j(x_1, x_2, \dots, x_n) = x_j$.

The ring \mathcal{F}_n is generated by these functions, i. e. we may write any function $f \in \mathcal{F}_n$ as a polynomial of the above functions. To see this given $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n) \in \{0, 1\}^n$, we define the function $f_\varepsilon = \prod_j (\delta_j + \pi_j)$, where, for each j , $\delta_j = 1 + \varepsilon_j$. Then $f_\varepsilon(\mathbf{x}) = 0$, except for $\mathbf{x} = \varepsilon$. Hence, any function g may be written

$$g = \sum_{\varepsilon \in S} f_\varepsilon,$$

where S is the set of ε such that $g(\varepsilon) = 1$.

As in the case of the natural numbers, we may define the complexity of the elements of \mathcal{F}_n . It will be the greatest function $f \mapsto \|f\|$ such that

$$\|0\| = \|1\| = 0; \quad \|\pi_j\| = 1; \quad \|f+g\| \leq \|f\| + \|g\|; \quad \|fg\| \leq \|f\| + \|g\|.$$

For any $\theta \in (0, 1)$, most of the elements of \mathcal{F}_n have complexity $\geq 2^{\theta n}$. The proof of this result is done by counting how many elements have complexity k , say a_k . It is easy to see that $a_0 = 2$, $a_1 = 2n$. From f and g with $\|f\| = j$ and $\|g\| = k - j$ we get, at most, four elements

with complexity $\leq k$. They are $f + g$, fg , $1 + f + g$, $1 + fg$. With these observations we get

$$a_k \leq 4(a_1 a_{k-1} + a_2 a_{k-2} + \cdots + a_{k-1} a_1).$$

It follows that $a_k \leq A_k$, where A_k is defined by

$$A_0 = 2; \quad A_1 = 2n; \quad A_k = 4 \sum_{j=1}^{k-1} A_j A_{k-j}.$$

From this definition we get

$$\sum_{k=0}^{\infty} A_k x^k = \frac{17 - \sqrt{1 - 32nx}}{8}; \quad A_k = \frac{1}{2(2k-2)} \binom{2k-2}{k} (8n)^k.$$

Hence

$$a_k \leq A_k \sim \frac{2^{5k}}{8\sqrt{2\pi}k^{3/2}} n^k.$$

Therefore for x large

$$\sum_{k=0}^x A_k \leq c \sum_{k=0}^x (32n)^k \leq c'(32n)^x \leq Ae^{Bx \log n},$$

hence if $x < 2^{\theta n}$, with $0 < \theta < 1$, we get

$$\sum_{k=0}^x A_k \ll \text{card}(\mathcal{F}_n) = 2^{2^n},$$

proving our assertion.

Each construction of $f(x_1, \dots, x_n)$ as a polynomial allows one to prove an assertion of type $\|f\| \leq a$. But from the polynomial expression we may get something more practical: a circuit that allows to compute $f(x_1, \dots, x_n)$ starting from the inputs x_j .

As in the case of natural numbers, it is difficult to prove inequalities of type $\|f\| > a$. In fact the situation is surprising: we have seen that in the set of functions with n variables, the complexity is usually larger than $2^{\theta n}$. Hence one would expect to have an easy task in defining a sequence of functions (f_n) , where f_n depends on n variables and such that $\|f_n\| > 2^{\theta n}$. On the contrary it has only been achieved that $\|f_n\| > p(n)$, where p is a polynomial of small degree (see [9], [5]). The problem here is not to prove that there exist sequences with $\|f_n\| > 2^{\theta n}$, which, as we have seen is easy, but to define explicitly a concrete sequence of functions for which this is so. When we speak of “define explicitly” we refer to a technical concept that needs some explanation. We must exclude easy solutions as: *let f_n the first function of n variables with maximum complexity.* We say that (f_n) is given explicitly if there is

an algorithm that computes the value of $f_n(x_1, \dots, x_n)$ in a reasonable time.

The problem $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ induces one to consider a special sequence of boolean functions. Let a be a natural number and consider $n = \binom{a}{2}$ the number of pairs. Our variables will be

$$x_{12}, x_{13}, x_{23}, x_{14}, x_{24}, x_{34}, \dots, x_{1a}, x_{2a}, \dots, x_{a-1a}.$$

In this way, each set of values of these variables $\in \{0, 1\}^n$ may be seen as a graph with a vertices and where $x_{jk} = 1$ if and only if the vertices j and k are connected by an edge of the graph. For each $b \leq a$ let $f_b^a(x_{12}, \dots, x_{a-1a})$ be the function that is equal 1 if and only if there is a set of b vertices such that all of them are connected in the graph.

It is plausible that $\|f_b^a\| \geq \binom{a}{b}$, since to compute the value of f_b^a in a given graph we need to check each set of b vertices. It can be shown that, if this is so, then $\mathbf{P} \neq \mathbf{NP}$. In this way to prove $\|f_b^a\| \geq \binom{a}{b}$ is, I think, the most promising path to solve the $\mathbf{P} \stackrel{?}{=} \mathbf{NP}$ question.

In the case of the complexity of natural numbers, an analogous question is the following, posed by Guy [4]:

Problem. *Is there a sequence of natural numbers (a_n) such that*

$$(1) \quad \lim_{n \rightarrow \infty} \frac{\|a_n\|}{\log a_n} > \frac{3}{\log 3}?$$

A good candidate is the sequence 2^n . All computed values satisfy $\|2^n\| = 2n$. Selfridge asks (see [4]) whether there exists any n with $\|2^n\| < 2n$.

If for some n and k we would have $2^n = 3^k$, (which is clearly impossible), the second expression would give us $\|2^n\| < 2n$. Of course the advantage would be greater for big n than for small n . Although the above is impossible, maybe another type of equality would yield $\|2^n\| < 2n$. For example, if for some n , 2^n written in base 3 has small digits. Again, this is unlikely but not impossible. Also, there may exist another type of expression of 2^n . The question here is whether a number of the form

$$(1 + 1)(1 + 1) \cdots (1 + 1),$$

may be written in some way with fewer 1's. We have almost a trivial example $4 = (1 + 1)(1 + 1) = 1 + 1 + 1 + 1$. Here we have the same number of 1's so that I call it an almost-example. Maybe there are non-trivial almost-examples, for example

$$2^{27} = 1 + (1 + 2 \cdot 3)(1 + 2^3 \cdot 3^2)(1 + 2^9 \cdot 3^3(1 + 2 \cdot 3^2)).$$

If we replace each 2 by 1+1 and each 3 by 1+1+1 we get an expression for 2^{27} with 57 ones, in which the multiplicative structure of 2^{27} is not used.

The above equality proves that $\|2^{27} - 1\| \leq 56$. In spite of an intense search I have not found an $n > 2$ such that $\|2^n - 1\| < 2n - 1$, but I think this may happen.

The evidence appears to be in favor of the existence of a sequence that satisfies (1). For example, we may look at figure 2. There we have put a little disk with center at each point $(n, \|n\|)$ with $1 \leq n \leq 2000$ and also we have drawn the smooth curves that bound $\|n\|$, i. e. $3(\log t)/\log 3$ and $3(\log t)/\log 2$, and also the curve $5 \log t/2 \log 2$. The points overlap and we see some lines parallel to the x -axis. We see that the upper bound appears to be bad and that apparently $\|n\| \leq 5 \log n/2 \log 2$, whereas in reality we have only proved that this inequality is true for almost all $n \in \mathbb{N}$.

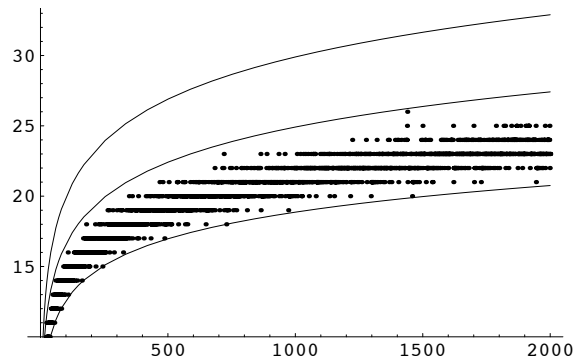


FIGURE 2. Graph of $\|n\|$.

But this figure says nothing about the limit $\lim \|n\|/\log n$, in which we are interested in. We only see that for the first 2000 values of n this sequence is bounded by the limits $5/2 \log 2$ and $3/\log 3$.

9. CONJECTURES

I have computed, using Proposition 1, the complexity of the first 200 000 natural numbers. Looking at these numbers, one sees many regularities. We will call them conjectures about the behavior of the function $\|\cdot\|$, although I have not much confidence in that they persist for larger numbers.

These conjectures were derived from tables such as this one

3	6	9	12	15	18	21	24
10	100	1000	10000	100000	1000000	10000000	100000000
	22	220	2200	22000	220000	2200000	22000000
	21	210	2101	21010	210100	2101000	21010000
		202	2100	21000	210000	2100000	21000000
		201	2020	20200	202000	2020000	20200000
		<u>122</u>	2010	20100	201000	2010000	20100000
			2002	20020	200222	2002220	20022200
			2001	20010	200200	2002000	20020000
			<u>1221</u>	20002	200100	2001000	20010000
			1220	20001	200020	2000200	20002000
			1212	<u>12221</u>	200010	2000100	20001000
			1211	12210	200002	2000020	20000200
			1201	12200	200001	2000010	20000100
			1122	12122	<u>122210</u>	2000002	20000020
			1121	12120	122100	2000001	20000010
			1112	12111	122000	<u>1222100</u>	20000002
				12110	121220	1221000	<u>20000001</u>
				12102	121200	1220000	<u>12221000</u>
				12101	121121	1212200	12210000
				12012	121110	1212000	12200000
				12010	121100	1211210	12122000
				12001	121022	1211100	12121201
				11221	121020	1211000	12120000

In this table we have written in columns the numbers with complexity $3n$ ($n = 1, 2, \dots, 8$), written in base 3 and in decreasing order.

The first observation: $\|3n\| = 3 + \|n\|$ is wrong. $\|107\| = 16$ and $\|321\| = \|1 + 2^6 5\| = 18$. But the following conjectures seem to be true:

Conjecture 1. *For each natural number n , there is an integer $a \geq 0$ such that $\|3^j n\| = 3(j - a) + \|3^a n\|$ for each natural number $j \geq a$.*

Let us define the set $A = \{n \in \mathbb{N} : \|3^j n\| = 3j + \|n\| \text{ for all } j\}$.

Conjecture 2. *For each pair of natural numbers p and q , there exists $a \geq 0$ such that, for $j \geq a$, we have $\|p(q3^j + 1)\| = 3j + 1 + \|p\| + \|q\|$.*

The main observation in the above table is that the greatest numbers with complexity $3n$ are those natural numbers contained in the

sequence $(3^n a_n)$, where a_n is given by

$$1, \frac{2(3+1)}{3^2}, \frac{2^6}{3^4}, \frac{2 \cdot 3 + 1}{3^2}, \frac{2(3^2+1)}{3^3}, \frac{2 \cdot 3^2 + 1}{3^3}, \frac{2^9}{3^6}, \\ \frac{2(3^3+1)}{3^4}, \frac{2 \cdot 3^3 + 1}{3^4}, \dots, \frac{2(3^k+1)}{3^{k+1}}, \frac{2 \cdot 3^k + 1}{3^{k+1}}, \dots$$

Conjecture 3. *There exist three transfinite sequences $(a_\alpha)_{\alpha < \xi}$, $(b_\alpha)_{\alpha < \xi}$, $(c_\alpha)_{\alpha < \xi}$ of rational numbers, such that the (greatest) numbers of complexity $3n$ (respectively $3n+1$, $3n+2$) are the (first) natural numbers contained in the sequence $(3^n a_\alpha)$, (resp. $(3^n b_\alpha)$, $(3^n c_\alpha)$).*

ξ is an infinite numerable ordinal such that $\omega\xi = \xi$.

These sequences start in the following way:

$$(a_\alpha), \quad 1, \frac{8}{9}, \frac{64}{81}, \frac{7}{9}, \frac{20}{27}, \dots \rightarrow \frac{2}{3}, \frac{160}{243}, \frac{52}{81}, \dots \rightarrow \frac{16}{27}, \frac{1280}{2187}, \frac{140}{243}, \dots \rightarrow \frac{5}{9} \dots \\ (b_\alpha), \quad \frac{4}{3}, \frac{32}{27}, \frac{10}{9}, \frac{256}{243}, \frac{28}{27}, \dots \rightarrow 1, \frac{80}{81}, \frac{26}{27}, \dots \rightarrow \frac{8}{9}, \frac{640}{729}, \frac{70}{81}, \dots \rightarrow \frac{64}{81} \dots \\ (c_\alpha), \quad 2, \frac{16}{9}, \frac{5}{3}, \frac{128}{81}, \frac{14}{9}, \dots \rightarrow \frac{4}{3}, \frac{320}{243}, \frac{35}{27}, \dots \rightarrow \frac{32}{27}, \frac{95}{81}, \frac{2560}{2187}, \dots \rightarrow \frac{10}{9} \dots$$

where the dots indicate infinite sequences, and where the indicated limits are not terms of the sequences.

Conjecture 4. *The three sequences are decreasing. The denominators of each term a_α , b_α or c_α are powers of 3.*

Conjecture 5. *The numbers of the sequence (a_α) are the numbers of the set*

$$\left\{ \frac{n}{3^{\lfloor n/3 \rfloor}} : \|n\| \equiv 0 \pmod{3}, \quad \text{and} \quad n \in A \right\},$$

ordered decreasingly.

Conjecture 6. *The numbers of the sequence (b_α) are the numbers of the set*

$$\left\{ \frac{n}{3^{(\lfloor n-1 \rfloor)/3}} : \|n\| \equiv 1 \pmod{3}, \quad \text{and} \quad n \in A \right\},$$

ordered decreasingly.

Conjecture 7. *The numbers of the sequence (c_α) are the numbers of the set*

$$\left\{ \frac{n}{3^{(\lfloor n-2 \rfloor)/3}} : \|n\| \equiv 2 \pmod{3}, \quad \text{and} \quad n \in A \right\},$$

ordered decreasingly.

The following conjectures are more doubtful. They are only based on a few cases.

Conjecture 8. *For all ordinals $\beta < \xi$ we have*

$$\lim_{n \rightarrow \infty} a_{\omega\beta+n} = c_\beta/3, \quad \lim_{n \rightarrow \infty} b_{\omega\beta+n} = a_\beta, \quad \lim_{n \rightarrow \infty} c_{\omega\beta+n} = b_\beta.$$

This is the basis of the assertion about the value of ξ , which appears to be at least $\xi = \omega^\omega$, since this is the least solution of $\omega\xi = \xi$.

The following assertions, along with conjecture 8, allow to predict, with some accuracy, the values of the transfinite sequences.

Conjecture 9. *The numbers of the sequence $b_{\omega\beta+n}$ that converges to $a_\beta = b/3^a$ (with $\|b\| = 3a$) are numbers from the sequences*

$$\frac{p(q3^j + 1)}{3^{a+j}}, \quad \text{where } b = pq, \text{ and, } \|p(q3^j + 1)\| = 3a + 3j + 1,$$

and those sporadic terms of the sequence $2^{3j+2}/3^{2j+1}$ contained between $\sup_{\gamma < \beta} a_\gamma$ and a_β .

Conjecture 10. *The numbers of the sequence $c_{\omega\beta+n}$ that converges to $b_\beta = b/3^a$ (with $\|b\| = 3a + 1$) are numbers from the sequences*

$$\frac{p(q3^j + 1)}{3^{a+j}}, \quad \text{where } b = pq, \text{ and, } \|p(q3^j + 1)\| = 3a + 3j + 2,$$

and those sporadic terms of the sequence $2^{3j+1}/3^{2j}$ contained between $\sup_{\gamma < \beta} b_\gamma$ and b_β .

Conjecture 11. *The numbers of the sequence $a_{\omega\beta+n}$ that converges to $c_\beta/3 = b/3^a$ (with $\|b\| = 3a - 1$) are numbers from the sequences*

$$\frac{p(q3^j + 1)}{3^{a+j}}, \quad \text{where } b = pq, \text{ and, } \|p(q3^j + 1)\| = 3a + 3j,$$

and those sporadic terms of the sequence $2^{3j}/3^{2j}$ contained between $\sup_{\gamma < \beta} \frac{1}{3}c_\gamma$ and $\frac{1}{3}c_\beta$.

In Conjecture 9, 10 and 11 we observe that some terms come from subsequent sequences. For example, the term $c_\omega = 320/243$ is the term corresponding to $j = 0$ of the sequence $2^6(4 \cdot 3^j + 1)/3^{j+5}$, that converges to $b_3 = 256/243$.

The above conjectures allow one to predict, for example, the 200 largest numbers with complexity 30.

The numbers with complexity 14 divided by 81, are

$$\begin{array}{cccc}
 c_0 = \frac{162}{81}, & c_1 = \frac{144}{81}, & c_2 = \frac{135}{81}, & c_3 = \frac{128}{81}, \\
 c_4 = \frac{126}{81}, & c_5 = \frac{120}{81}, & c_6 = \frac{117}{81}, & c_7 = \frac{114}{81}, \\
 c_9 = \frac{112}{81}, & c_{10} = \frac{111}{81}, & c_{11} = \frac{110}{81}, & c_{13} = \frac{109}{81}, \\
 c_{\omega+1} = \frac{105}{81}, & c_{\omega+2} = \frac{104}{81}, & c_{\omega+3} = \frac{102}{81}, & c_{\omega+6} = \frac{100}{81}, \\
 c_{\omega+8} = \frac{99}{81}, & c_{\omega+10} = \frac{98}{81}, & c_{\omega+14} = \frac{97}{81}, & c_{\omega 2} = \frac{95}{81}, \\
 c_{\omega 2+3} = \frac{93}{81}, & c_{\omega 2+5} = \frac{92}{81}, & c_{\omega 2+8} = \frac{91}{81}, & c_{\omega 3+4} = \frac{88}{81}, \\
 c_{\omega 3+7} = \frac{87}{81}, & c_{\omega 3+15} = \frac{86}{81}, & c_{\omega 4+2} = \frac{85}{81}, & c_{\omega 5+1} = \frac{83}{81}, \\
 c_{\omega^2+\omega+2} = \frac{79}{81}, & c_{\omega^2+\omega 2+3} = \frac{77}{81}, & & \\
 & & \frac{71}{81}, & \frac{69}{81}, \quad \frac{67}{81}, \quad \frac{59}{81},
 \end{array}$$

For the last four numbers I do not have enough data to know the corresponding ordinal.

REFERENCES

- [1] N. ALON & J. H. SPENCER, *The probabilistic method*, John Wiley and Sons, New York, 1992.
- [2] H. CHERNOFF, *A measure of the asymptotic efficiency for test of a hypothesis based on the sum of observations*, *Annals of Mathematical Statistics*, **23** (1952) 493–509.
- [3] M. R. GAREY & D. S. JOHNSON, *Computers and Intractability, a guide to the theory of NP-completeness*, W. H. Freeman and Co., 1979
- [4] R. K. GUY, *What is the least number of ones needed to represent n using only $+$ and \times (and parentheses)?*, *American Mathematical Monthly* **93** (1986) 189–190.
- [5] J. HASTAD, *The Shrinkage exponent of the Morgan formulas is 2*, *Siam J. Comput.* **27** (1998) 48–64.
- [6] K. MAHLER & J. POPKEN, *On a maximum problem in arithmetic*, (Dutch), *Nieuw Arch. Wiskunde* (3) **1** (1953) 1–15.
- [7] D. A. RAWSTHORNE, *How many 1's are needed?*, *Fibonacci Quart.* **27** (1989) 14–17.
- [8] N. J. A. SLOANE & S. PLOUFFE, *The Encyclopedia of Integer Sequences*, Academic Press, London, 1995.
<http://oeis.org>.

- [9] U. ZWICK, *A $4n$ lower bound on the combinatorial complexity of certain symmetric boolean functions over the basis of unate dyadic boolean functions*, Siam J. Comput. **20** (1991) 499–505.

FACULTAD DE MATEMÁTICAS, UNIVERSIDAD DE SEVILLA,
APDO. 1160, 41080-SEVILLA, SPAIN
E-mail address: `arias@us.es`