

A Normal(?) Sequence

Aresh Pourkavoos

May 23, 2022

This idea is about OEIS entry [A330731](#), which I submitted, and most of the information here is also available there. Nonetheless, I wanted to explore the sequence in a bit more detail with less of a wait time before publication, hence the separate source.

The sequence is infinite and binary, i.e. its entries are all either 0 or 1. I designed it to be *normal*, meaning that as the number of terms increases, every subsequence of the same length appears with equal frequency. For example, a randomly chosen subsequence of length 5 has a 1 in 32 chance of being 01001. At the moment, I still have not proven this property, but experimentally, this seems to be the case. However, I hope to prove that it is normal eventually, and in fact, I believe that the frequency of a given substring converges to its final value *faster* than a “random” sequence.

The sequence is defined by the following procedure to create the next digit, starting from the empty string:

- List all of the tails, or suffixes, of the sequence so far in order from longest to shortest. For example, the tails of 011 are 011, 11, 1, and the empty string.
- For each tail, count the number of times it occurs elsewhere in the sequence (i.e. other than the end), and note how many times 0 or 1 appears afterward. For the example above, 011 and 11 do not appear anywhere other than the end, but 1 appears in one other place (the middle), where it is followed by a 1.
- The *longest* tail where one bit appears afterward more often than the other decides which bit comes next, namely the bit which appears *less* often. Since 1 is followed by 0 zero times and by 1 once, the next bit is 0. This is why the occurrences of the empty string (and the bits that come after them) were not counted in the previous step: these numbers do not matter for the result. If they did, though, 0 would still be added since the empty string is followed by 0 once and by 1 twice.
- In the edge case where all of these comparisons result in a tie, 0 is the next bit. For example, in 010011, the tails are 010011, 10011, 0011, 011, 11, 1, and the empty string. All but the last two don't occur elsewhere, 1 is followed by 0 once and by 1 once, and the empty string is followed by both three times. Hence the next bit is 0 by default.

The naive algorithm to generate the sequence per the definition takes $O(n^3)$ time for the first n terms. This is because to generate each new term, the frequencies of $O(n)$ different tails are checked (from the whole sequence so far down to a relatively short tail), and each counting requires a pass over all terms, which takes $O(n)$ time. In other words, each term takes $O(n^2)$ time to generate. However, with the right data structures, I was able to reduce the time per term to $O(n)$, bringing the time for the first n terms to $O(n^2)$.

The following page contains a C program that prints the first 8192 terms of the sequence.

```

1  #include <stdio.h>
2  #define N_TERMS 8192
3  // Stores generated terms
4  int a[N_TERMS];
5  // b[j-1] is the number of bits before (not including) a[n-j]
6  // which match the tail of the first n entries
7  int b[N_TERMS];
8  // c induces a linked list structure on b
9  // with an extra node to make it cyclic
10 // Indices are offset by 1 since the extra node is in front
11 int c[N_TERMS];
12 int cEnd = 0;
13 int main() {
14     FILE *bfile = fopen("b330731.txt", "w+");
15     for (int n = 0; n < N_TERMS; n++) {
16         // Append new digit to list from previous loop
17         // or (n = 0) initialize c
18         c[n] = 0;
19         c[cEnd] = n;
20         cEnd = n;
21         // Find new digit by iterating over b
22         // using the indices given in c
23         int newD = 0;
24         int freq0 = 0;
25         int freq1 = 0;
26         int prevTail = -1;
27         int j = c[0];
28         while (j != 0) {
29             int currTail = b[j-1];
30             if (currTail != prevTail) {
31                 // All tails of a given length have been accumulated in freqs,
32                 // so they need to be compared to decide whether to continue
33                 if (freq1 != freq0) {
34                     break;
35                 }
36                 freq0 = 0;
37                 freq1 = 0;
38             }
39             // Use digit that comes after current tail to adjust freqs
40             if (a[n-j] == 0) {
41                 freq0++;
42             } else {
43                 freq1++;
44             }
45             prevTail = currTail;
46             j = c[j];
47         }
48         // 0 is chosen by default (if freq1 == freq0)
49         if (freq1 < freq0) {
50             newD = 1;
51         }
52         // Update matching tail lengths
53         j = 0;
54         for (int numVisited = 0; numVisited < n; numVisited++) {

```

```

55     int k = c[j];
56     if (a[n-k] == newD) {
57         // Matching tail grows by 1, position in list is unaffected
58         b[k-1]++;
59         j = k;
60     } else {
61         // Matching tail resets to 0, moved to back of list
62         b[k-1] = 0;
63         c[j] = c[k];
64         if (cEnd == k) {
65             cEnd = j;
66         }
67         c[k] = 0;
68         c[cEnd] = k;
69         cEnd = k;
70     }
71 }
72 a[n] = newD;
73 b[n] = 0;
74 printf("%d", newD);
75 fprintf(bfile, "%d_%d\n", n, newD);
76 }
77 printf("\n");
78 fclose(bfile);
79 return 0;
80 }

```

A line-by-line breakdown:

- The arrays a , b and c declared on lines 3-5 are responsible for the algorithm's $O(n)$ space complexity. Since the program generates a fixed number of terms here, the arrays are initialized to their full length. However, they could also be implemented as unbounded arrays and have an element appended to them at the beginning of each iteration of the main loop.
- The $O(n^2)$ time complexity comes from the nested for loops. Each outer loop on line 8 generates a single term, and the loops on lines 17 and 30 take $O(n)$ steps each.
- For all integers $n \in [0, \text{N_TERMS}]$, after n iterations of the outer loop (i.e. at all points right after n is incremented and the next loop is about to begin):
 - $a[i] = A330731(i)$ for all $i \in [0, n)$, i.e. a stores the first n terms of A330731.
 - $b[i]$ stores the length of the longest subsequence ending at $a[n - i - 1]$ which matches the

For all arrays, the terms beyond the first n are not initialized and thus undefined.