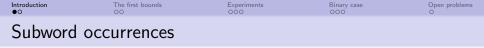
Introduction	The first bounds	Experiments	Binary case	Open problems
OO	OO	000	000	O

Maximal number of subword occurrences in a word

Wenjie Fang LIGM, Université Gustave Eiffel arXiv:2406.02971

AofA 2024, 17 June 2024, University of Bath

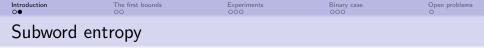


A word: $w = (w_1, \ldots, w_\ell)$ with w_i in a finite alphabet \mathcal{A} .

Two notions of patterns: **subword** (scattered) and factor (consecutive) Example: 01 occurs 5 times as subword in 011001, but twice as factor Counting pattern occurrences: harder for subwords, easier for factors

Occurrence of u in w: a subset of positions in w that gives uocc(w, u) or $\binom{w}{u}$: number of occurrences of u as subword of wFlajolet, Szpankowski, Vallée (2006): normal limit law and large deviation of occ(w, u) for fixed u and $w \sim \text{Unif}(A^n)$, $n \to \infty$.

Quite some research in many directions! Difficulty from self-correlation.

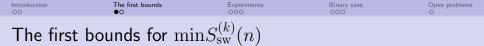


Given $w \in A^*$, what are its most frequent subwords? Related to data-mining for finding patterns appearing frequently. Surprisingly difficult! Complexity unknown.

$$\begin{split} \max(w) &:= \max_u \operatorname{occ}(w, u) : \text{ maximal number of subword occurrences} \\ & \mathsf{Subword entropy:} \ S_{\mathrm{sw}}(w) := \log_2 \max(w) : \\ & \mathsf{Easy to maximize:} \ \max(0^n) = \binom{n}{\lfloor n/2 \rfloor}, \ S_{\mathrm{sw}}(0^n) = n + O(\log_2 n). \end{split}$$

Minimal subword entropy for $|\mathcal{A}| = k$, length *n*:

$$\min S_{\mathrm{sw}}^{(k)}(n) := \min_{u \in \mathcal{A}^n} S_{\mathrm{sw}}(w).$$



Trivial upper bound (from 0^n): for some constant c,

$$\min S_{sw}^{(k)}(n) \le n - \frac{1}{2}\log_2 n + c.$$

Easy lower bound: for some constant c',

$$\min S_{\rm sw}^{(k)}(n) \ge \log_2(1+k^{-1})n - \frac{1}{2}\log_2 n + c'.$$

Reasoning: For any fixed w, take random word u of length αn . Then

$$S_{sw}(w) \ge \log_2 \mathbb{E}[\operatorname{occ}(w, u)] = \log_2 \left(\binom{n}{\alpha n} k^{-\alpha n} \right)$$

Maximized at $\alpha = (k+1)^{-1}$. Holds for all w.

Introduction	The first bounds	Experiments	Binary case	Open problems
OO	O●	000	000	O
Super-additi	vity			

Proposition (Super-additivity of $minS_{sw}$)

Given $k \ge 2$, for $n, m \ge 1$,

$$\min S_{sw}^{(k)}(n+m) \ge \min S_{sw}^{(k)}(n) + \min S_{sw}^{(k)}(m).$$

Not difficult, but a little twist!

Lemma (Fekete's lemma)

For (g_n) super-additive, when $n \to +\infty$, then g_n/n either tends to $+\infty$, or converges to some limit L.

Corollary

The minimal subword entropy per letter $\min S_{sw}^{(k)}(n)/n$ has a limit L_k :

$$\log_2(1+k^{-1}) \le L_k \le 1.$$

Better bounds?

Introduction	The first bounds	Experiments	Binary case	Open problems
00	00	OO	000	O
D'		1 .		

Binary words with minimal entropy

When no idea, brute force!

Very hard... Start with the binary case.

n	Words \boldsymbol{w} with min. subword entropy	$\max \operatorname{occ}(w)$	Symmetry
1	0	1	Р
2	01	1	А
3	001	2	
4	0110	2	Р
5	01110	3	Р
6	011001	5	А
$\overline{7}$	0110001	6	
8	01110001	9	А
9	011000110	16	Р
10	0110001110	22	
11	01110001110	33	Р
12	011000111001	52	А
13	0111001001110	72	Р
14	01100010111001	108	А
15	011000101110001	162	

Introduction	The first bounds	Experiments	Binary case	Open problems
00	00	000	000	O
Binary word	s with minimal	entropy (cont'd)	

Interesting, some more!

n	Words \boldsymbol{w} with min. subword entropy	$\maxocc(w)$	Symmetry
16	0111000101110001	252	А
17	01100011111000110	390	Р
18	011100100101110001	588	
19	0110001011101000110	900	Р
	0110001110110001110		
20	01110001011011000110	1320	
21	011100011011010001110	2049	
22	0110001110101000111001	2958	А
23	01110001011011010001110	4473	Р
24	011000111010101000111001	6979	А
25	0111000101101101000111001	10602	
26	01110001011011001000111001	15962	
27	011100010101110101000111001	24150	
28	01100011110100100101111000110	36450	
	0111000101110101000101110001		А
29	01100011101010001010111000110	53671	Р
30	011000111001100010101111000110	83862	

Introduction	The first bounds	Experiments	Binary case	Open problems
00	OO	000	000	O
Binary word	s with minimal	entropy ((cont'd 2)	

Confusing... A last push!

n	Words w with min. subword entropy	$\maxocc(w)$	Symmetry
31	0110001110101000101011110001110	127998	
32	01100011101010001010111010001110	189131	
33	011000111101010001011011010001110	288900	
34	0110001110101000101011101001001110	442386	
35	01110001011011001000110111001001110	681966	

The last line took 11 days on a single core.

Naïve complexity: $O(4^n n^2)$. A lot of optimizations needed.

Observations

- For larger *n*, symmetry runs out.
- Average run length 1.6-2, mostly 1, 2, 3, but length 4 and 5 exist.
- Growth rate slightly larger than 1.5 given by lower bound of L_2 .

Idea: Find words like them, but analyzable.

Introduction	The first bounds	Experiments	Binary case	Open problems
00	00	000	●00	O
			4 a	

Three families inspired by experiments

Average run length slightly less than 2. Most runs have length 1, 2, 3.

Candidates: $(01)^m$, $(0011)^m$, $(000111)^m$.

Proposition

The following words has a most frequent subword of the form

- $(01)^m$: subword $(01)^r$;
- $(0011)^m$: subword $(01)^r$;
- $(000111)^m$: subword $(0011)^r$.

With local analysis in subword.

Key result for analysis, as most frequent subwords are hard to compute! Experimentally, periodic words have periodic most frequent subwords. But no proof!

Introduction	The first bounds	Experiments	Binary case	Open problems
00	00	000	O●O	O

Generating functions of periodic subword occurrences

Occurrence generating function: $f_{w,u}(x,y) = \sum_{m,r>0} \operatorname{occ}(w^m, u^r) x^m y^r$

Proposition

$$f_{01,01} = \frac{1-x}{(1-x)^2 - xy},$$

$$f_{0011,01} = \frac{1-x}{(1-x)^2 - 4xy},$$

$$f_{000111,0011} = \frac{(1-x)^3}{(1-x)^4 - 9x(1+2x)^2y}.$$

 $\max(w^m) = \max_r [x^m y^r] f_{w,u}$ for these families.

In fact a universal and effective result!

Theorem

For any words $w, v \in \mathcal{A}^*$, the g.f. $f_{w,v}(x, y)$ is rational in x, y.

Problem is that we don't know the most frequent subwords...

Introduction	The first bounds	Experiments	Binary case	Open problems
00	OO	000	00●	O
۸ I.		т		

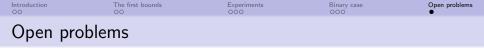
Asymptotics and bounds on L_2

Proposition			
Word w	Subword	Max at	$S_{\rm sw}(w)$
$(01)^m$	$(01)^r$	$r = \frac{m}{\sqrt{5}}$	$m \log_2 \frac{3+\sqrt{5}}{2} + \frac{\log_2 m}{2} + O(1)$
$(0011)^m$	$(01)^{r}$	$r = \frac{m}{\sqrt{2}}$	$m\log_2(3+2\sqrt{2}) + \frac{\log_2 m}{2} + O(1)$
$(000111)^m$	$(0011)^r$	$r = \alpha m$	$m\gamma - \frac{\log_2 m}{2} + O(1)$
Here, $\alpha \approx 0.66$ is the pos. sol. of $457\alpha^4 - 246\alpha^2 + 72\alpha - 27 = 0$, and			
$\gamma = \alpha \log_2 9 + 2\alpha \log_2 \frac{1+2\zeta}{(1-\zeta)^2} - (1-\alpha) \log_2 \zeta,$			
ζ	$T = \frac{1 - 9\alpha}{1 - 9\alpha}$	$\frac{+\sqrt{73\alpha^2}-}{4+4\alpha}$	$\frac{18\alpha + 9}{2}$.

The last needs (automated) ACSV or saddle-point on large powers.

Upper bounds of L_2 : 0.694..., 0.636..., 0.654....

We have $0.585... = \log_2(3/2) \le L_2 \le \frac{1}{2}\log_2(1+\sqrt{2}) = 0.636...$

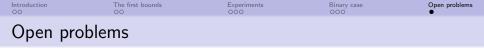


- Value of L_2 ? Value of other L_k ? Better bounds?
- Does periodic word have a quasi-periodic most frequent subword?
- Any structure on words almost realizing $\min S^{(k)}_{sw}(n)$?

Difficult "minimal of maximal" structure, chaos in experimental data

A lot of unknowns, even intuitive ones!

- Is $\min S_{sw}^{(k)}(n)/n$ ultimately increasing?
- For any w, is every most frequent subword is of length $\leq |w|/2?$



- Value of L_2 ? Value of other L_k ? Better bounds?
- Does periodic word have a quasi-periodic most frequent subword?
- Any structure on words almost realizing $\min S^{(k)}_{sw}(n)$?

Difficult "minimal of maximal" structure, chaos in experimental data

A lot of unknowns, even intuitive ones!

- Is $\min S_{sw}^{(k)}(n)/n$ ultimately increasing?
- For any w, is every most frequent subword is of length $\leq |w|/2?$

Thank you for your attention!