The Vocabulary of Grammar-Based Codes and the Logical Consistency of Texts

Łukasz Dębowski Idebowsk@ipipan.waw.pl



Institute of Computer Science Polish Academy of Sciences



Consider texts in a natural language (such as English):

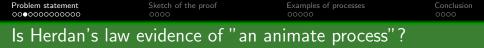
- V the number of different words in the text,
- **n** the length of the text.

We observe the relationship

$V\propto n^{eta},$

where β is between **0.5** a **1** depending on a text.

- Władysław Kuraszkiewicz, Józef Łukaszewicz (1951),
- Pierre Guiraud (1954),
- Gustav Herdan (1964),
- H. S. Heaps (1978).



If a Martian scientist sitting before his radio in Mars accidentally received from Earth the broadcast of an extensive speech [...], what criteria would he have to determine whether the reception represented the effect of animate process [...]? It seems that [...] the only clue to the animate origin would be this: the arrangement of the occurrences would be neither of rigidly fixed regularity such as frequently found in wave emissions of purely physical origin nor yet a completely random scattering of the same.

— George Kingsley Zipf (1965:187)



Zipf's and Herdan's law are observed if the letters and spaces in the text are obtained by pressing keys at random.

- Benoit B. Mandelbrot (1953),
- George A. Miller (1957).

Problem statement	Sketch of the proof	Examples of processes	Conclusion 0000
The new expla	anation of Herda	n's law	

We will prove a theorem which can be stated informally in this way, for $\beta \in (0, 1)$:

If a text of length **n** describes $\geq \mathbf{n}^{\beta}$ independent facts in a repetitive way then the text contains $\geq \mathbf{n}^{\beta} / \log \mathbf{n}$ distinct words.

For the formal statement, we shall adopt two postulates:

- Words are understood as nonterminal symbols in the shortest grammar-based encoding of the text.
- Texts are emitted by a finite-energy strongly nonergodic source.
- Facts are independent binary variables which can be predicted from the text in a shift-invariant way.

Problem statement	Sketch of the proof 0000	Examples of processes	Con 000
A context-free g	rammar that g	enerates one text	

nclusion

$$\left\{ \begin{array}{l} \mathsf{A}_1 \rightarrow \mathsf{A}_2\mathsf{A}_2\mathsf{A}_4\mathsf{A}_5\mathsf{dear_children}\mathsf{A}_5\mathsf{A}_3\mathsf{all}.\\ \mathsf{A}_2 \rightarrow \mathsf{A}_3\mathsf{you}\mathsf{A}_5\\ \mathsf{A}_3 \rightarrow \mathsf{A}_4_to_-\\ \mathsf{A}_4 \rightarrow \mathsf{Good_morning}\\ \mathsf{A}_5 \rightarrow \mathsf{,_} \end{array} \right.$$

Good morning to you, Good morning to you, Good morning, dear children, Good morning to all.

Problem statement	Sketch of the proof	Examples of processes	Conclusion
	0000	00000	0000
The vocabulary	size and grammar	-based codes	

The vocabulary size of a grammar:

$$\mathbb{V}[G] := n, \quad \text{jeżeli} \quad G = \left\{ \begin{array}{l} \mathsf{A}_1 \to \alpha_1, \\ \mathsf{A}_2 \to \alpha_2, \\ \dots, \\ \mathsf{A}_n \to \alpha_n \end{array} \right\}.$$

A grammar-based code is a function of form $C = B(\Gamma(\cdot))$, where

- a grammar transform $\Gamma : \mathbb{X}^+ \to \mathcal{G}$, for each string $\mathbf{w} \in \mathbb{X}^+$, returns a gramar $\Gamma(\mathbf{w})$ that generates this string.
- **2** a grammar encoder $B : \mathcal{G} \to \mathbb{X}^+$ codes the grammar as (another) string.
 - John C. Kieffer, En-hui Yang (2000),
 - Moses Charikar, Eric Lehman, ..., Abhi Shelat (2005).

Admissibly minimal codes				
Problem statement	Sketch of the proof	Examples of processes	Conclusion 0000	

Let $X = \{0, 1, ..., D - 1\}$. A grammar transform Γ and the code $B(\Gamma(\cdot))$ are called admissibly minimal if

 ${\color{black} 9} \ |B(\Gamma(w))| \leq |B(G)| \ \text{for each grammar } G \ \text{that generates } w,$

- ${f Q}$ the encoder has the form ${f B}({f G})={f B}^*_{\sf S}({f B}_{\sf N}({f G})),$
- $\textcircled{\textbf{3}} \quad \textbf{B}_{N} \text{ encodes the grammar}$

$$\mathbf{G} = \{\mathbf{A}_1 \rightarrow \alpha_1, \mathbf{A}_2 \rightarrow \alpha_2, ..., \mathbf{A}_n \rightarrow \alpha_n\}$$

as a string of integers

 $\mathsf{B}_{\mathsf{N}}(\mathsf{G}) := \mathsf{F}_1^*(\alpha_1)\mathsf{D}\mathsf{F}_2^*(\alpha_2)\mathsf{D}...\mathsf{D}\mathsf{F}_{\mathsf{n}}^*(\alpha_{\mathsf{n}})(\mathsf{D}+1),$

using $F_i(x) := x$ for $x \in \mathbb{X}$ and $F_i(A_j) := D + 1 + j - i$, **9** $B_S : \{0\} \cup \mathbb{N} \to \mathbb{X}^+$ is an injection, the set $B_S(\{0\} \cup \mathbb{N})$ is prefix-free, $|B_S(\cdot)|$ is non-decreasing and

 $\limsup_{n\to\infty}|B_S(n)|/\log_Dn=1.$

Problem statement	Sketch of the proof 0000	Examples of processes	Conclusion 0000	
Two classes of stochastic processes				

Let $(X_i)_{i\in\mathbb{Z}}$ be a stochastic process on the space $(\Omega, \mathfrak{J}, \mathsf{P})$, where $X_i : \Omega \to \mathbb{X}$ for a countable alphabet \mathbb{X} . Denote blocks as $X_{m:n} := (X_i)_{m \leq k \leq n}$.

The proces $(X_i)_{i\in\mathbb{Z}}$ is called strongly nonergodic if there exist variables $(Z_k)_{k\in\mathbb{N}} \sim IID$, $P(Z_k = 0) = P(Z_k = 1) = \frac{1}{2}$, and functions $s_k : \mathbb{X}^* \to \{0, 1\}$, $k \in \mathbb{N}$, such that

$$\lim_{n\to\infty}\mathsf{P}(\mathsf{s}_k(\mathsf{X}_{t+1:t+n})=\mathsf{Z}_k)=1,\qquad \forall t\in\mathbb{Z}.$$

 $\mathbf{Y} = \sum_{\mathbf{k} \in \mathbb{N}} 2^{-\mathbf{k}} \mathbf{Z}_{\mathbf{k}}$ is measurable against the shift-invariant σ -field.

The process $(X_i)_{i \in \mathbb{Z}}$ is called a finite-energy process if

$$\mathsf{P}(\mathsf{X}_{t+1:t+m}=\mathsf{u}|\mathsf{X}_{t-n:t}=\mathsf{w})\leq\mathsf{K}\mathsf{c}^m,\qquad\forall t\in\mathbb{Z}.$$

Problem statement	Sketch of the proof	Examples of processes 00000	Conclusion 0000
The main result			

$$\mathsf{U}_{\delta}(\mathsf{n}) := \left\{\mathsf{k} \in \mathbb{N} : \mathsf{P}\left(\mathsf{s}_{\mathsf{k}}(\mathsf{X}_{1:\mathsf{n}}) = \mathsf{Z}_{\mathsf{k}}\right) \geq \delta\right\}.$$

Theorem 1

Let $(X_i)_{i\in\mathbb{Z}}$ be a stationary finite-energy strongly nonergodic process over a finite alphabet X. Suppose that

$$\liminf_{n\to\infty}\frac{\operatorname{card}\mathsf{U}_{\delta}(\mathsf{n})}{\mathsf{n}^{\beta}}>0$$

for some $\beta \in (0,1)$ and $\delta \in (\frac{1}{2},1)$. Then

$$\limsup_{n\to\infty}\mathbb{E}\,\left(\frac{\mathbb{V}[\Gamma(X_{1:n})]}{n^\beta(\log n)^{-1}}\right)^p>0,\quad p>1,$$

for any admissibly minimal grammar transform $\mathbf{\Gamma}$.

Problem statement	Sketch of the proof	Examples of processes	Conclusion 0000
The first asso	ciated result		

Denote the mutual information between \mathbf{n} -blocks

$$\mathsf{E}(\mathsf{n}) := \mathsf{I}(\mathsf{X}_{1:\mathsf{n}};\mathsf{X}_{\mathsf{n}+1:2\mathsf{n}}) = \mathbb{E} \log \frac{\mathsf{P}(\mathsf{X}_{1:2\mathsf{n}})}{\mathsf{P}(\mathsf{X}_{1:\mathsf{n}})\mathsf{P}(\mathsf{X}_{\mathsf{n}+1:2\mathsf{n}})}.$$

Theorem 2

Let $(X_i)_{i\in\mathbb{Z}}$ be a stationary strongly nonergodic process over a finite alphabet X. Suppose that

$$\liminf_{\mathsf{n}\to\infty}\frac{\mathsf{card}\,\mathsf{U}_\delta(\mathsf{n})}{\mathsf{n}^\beta}>0$$

for some $eta \in (0,1)$ and $\delta \in (rac{1}{2},1)$. Then

$$\limsup_{\mathsf{n}\to\infty}\frac{\mathsf{E}(\mathsf{n})}{\mathsf{n}^\beta}>0.$$

Problem statement	Sketch of the proof	Examples of processes	Conclusion
00000000000000	0000		0000
The second asso	ciated result		

Theorem 3

Let $(X_i)_{i\in\mathbb{Z}}$ be a stationary finite-energy process over a finite alphabet $\mathbb{X}.$ Suppose that

$$\liminf_{n\to\infty}\frac{\mathsf{E}(n)}{n^\beta}>0$$

for some $eta \in (0,1)$. Then

$$\limsup_{n\to\infty}\mathbb{E}\,\left(\frac{\mathbb{V}[\Gamma(\mathsf{X}_{1:n})]}{n^\beta(\log n)^{-1}}\right)^p>0,\quad p>1,$$

for any admissibly minimal grammar transform $\mathbf{\Gamma}$.

 Problem statement
 Sketch of the proof
 Examples of processes
 Conclusion

 000000000000
 0000
 0000
 0000
 0000

Mutual information for natural language

Basing on Shannon's (1950) estimates of conditional entropy of printed English, Hilberg (1990) conjectured that mutual information between two **n**-blocks drawn from natural language satisfies

$\mathsf{E}(\mathsf{n}) \asymp \mathsf{n}^{\beta}, \quad \beta \approx 1/2.$

- Theorem 2 a rational motivation of Hilberg's conjecture
- Theorem 3 Hilberg's conjecture implies Herdan's law

Besides that, Theorem 1 indicates

- why Herdan's law may be observed for the same text translated into different languages,
- why certain variation of the exponent in Herdan's law may be expected depending on a text.

000000000000000000000000000000000000000	usion

1 Problem statement

- 2 Sketch of the proof
- 3 Examples of processes
- 4 Conclusion

Entropy, pseudoentropy, and code length

Denote the entropy of the \mathbf{n} -block and entropy rate as:

$$H(n):=H(X_{1:n})=-\mathbb{E}\,\log \mathsf{P}(X_{1:n}),\quad h:=\lim_{n\to\infty}H(n)/n.$$

Define also "pseudoentropy"

$$\mathsf{H}^{\mathsf{U}}(\mathsf{n}):=\mathsf{hn}+[\log 2-\eta(\delta)]\cdot\mathsf{card}\:\mathsf{U}_{\delta}(\mathsf{n}).$$

Let $\mathbb{X}=\{0,1,...,D-1\}$ and $\mathsf{C}=\mathsf{B}(\Gamma(\cdot))$ be an admissibly minimal codes. Put its expected length

$$\mathsf{H}^\mathsf{C}(\mathsf{n}) := \mathbb{E} \, \left| \mathsf{C}(\mathsf{X}_{1:\mathsf{n}}) \right| \log \mathsf{D}.$$

We have inequality

$$H^C(u) \geq H(n) \geq H^U(n)$$

and equality of rates

$$\lim_{n\to\infty} H^{\mathsf{C}}(n)/n = \lim_{n\to\infty} H(n)/n = \lim_{n\to\infty} H^{\mathsf{U}}(n)/n = h.$$



Let a function $f:\mathbb{N}\to\mathbb{R}$ satisfy $\lim_k f(k)/k=0$ and $f(n)\geq 0$ for all but finitely many n.Then we have $2f(n)-f(2n)\geq 0$ for infinitely many n.

The equalities and inequalities on the previous slide yield

$$\begin{split} \liminf_{n \to \infty} \frac{\text{card } U_{\delta}(n)}{n^{\beta}} > 0 & \Longrightarrow & \limsup_{n \to \infty} \frac{\mathsf{E}^{\mathsf{C}}(n)}{n^{\beta}} > 0, \quad (\text{Th. 1}) \\ \liminf_{n \to \infty} \frac{\text{card } U_{\delta}(n)}{n^{\beta}} > 0 & \Longrightarrow & \limsup_{n \to \infty} \frac{\mathsf{E}(n)}{n^{\beta}} > 0, \quad (\text{Th. 2}) \\ \liminf_{n \to \infty} \frac{\mathsf{E}(n)}{n^{\beta}} > 0 & \Longrightarrow & \limsup_{n \to \infty} \frac{\mathsf{E}^{\mathsf{C}}(n)}{n^{\beta}} > 0. \quad (\text{Th. 3}) \end{split}$$

for E(n) = 2H(n) - H(2n) and $E^{C}(n) = 2H^{C}(n) - H^{C}(2n)$.



$$\mathsf{E}^{\mathsf{C}}(\mathsf{n}) = \mathbb{E}\left[|\mathsf{C}(\mathsf{X}_{1:\mathsf{n}})| + |\mathsf{C}(\mathsf{X}_{\mathsf{n}+1:2\mathsf{n}})| - |\mathsf{C}(\mathsf{X}_{1:2\mathsf{n}})|\right]\log\mathsf{D}.$$

For an admissibly minimal code $C = B(\Gamma(\cdot))$, we have

$$\mathsf{C}(\mathsf{u})|+|\mathsf{C}(\mathsf{v})|-|\mathsf{C}(\mathsf{w})|\leq\mathsf{W}_0\mathbb{V}[\mathsf{\Gamma}(\mathsf{w})](1+\mathbb{L}(\mathsf{w})),$$

where w = uv for $u, v \in \mathbb{X}^+$, $\mathbb{L}(w)$ denotes the maximal length of a repeat in w, and $W_0 = |B_S(D+1)|$.

For a finite-energy process
$$(X_i)_{i \in \mathbb{Z}}$$
,
$$\sup_{n \in \mathbb{N}} \mathbb{E} \left(\frac{\mathbb{L}(X_{1:n})}{\log n} \right)^q < \infty, \quad q > 0.$$

Problem statement	Sketch of the proof	Examples of processes	Conclusion
		00000	

1 Problem statement

- 2 Sketch of the proof
- 3 Examples of processes

4 Conclusion

Problem statement	Sketch of the proof	Examples of processes	Conclusion
	0000	0●000	0000
The binary excha	angeable process		

Consider a family of binary IID processes

$$\mathsf{P}(\mathsf{X}_{1:n} = \mathsf{x}_{1:n} | \theta) = \prod_{i=1}^{n} \theta^{\mathsf{x}_i} (1-\theta)^{1-\mathsf{x}_i}.$$

Construct such a process $(X_i)_{i\in\mathbb{Z}}$ that

$$P(X_{1:n} = x_{1:n}) = \int_0^1 P(X_{1:n} = x_{1:n}|\theta)\pi(\theta)d\theta$$

for a prior $\pi(\theta) > 0$. For $\mathbf{Y} = \lim_{n \to \infty} n^{-1} \sum_{i=1}^{n} X_i$ we have $\mathsf{P}(\mathbf{Y} \le \mathbf{y}) = \int_{0}^{\mathbf{y}} \pi(\theta) d\theta.$

Process $(X_i)_{i\in\mathbb{Z}}$ is strongly nonergodic, because Y has a continuous distribution. However, block $X_{1:n}$ is conditionally independent from $X_{n+1:2n}$ given the sum $S_n := \sum_{i=1}^n X_i$. Thus

$$\mathsf{E}(n) = \mathsf{I}(\mathsf{X}_{1:n};\mathsf{X}_{n+1:2n}) = \mathsf{I}(\mathsf{S}_n;\mathsf{X}_{n+1:2n}) \leq \mathsf{H}(\mathsf{S}_n) \leq \log(n+1).$$

The unesses which I increased at Canta Fallwatitute				
000000000000	0000	00000	0000	
Problem statement	Sketch of the proof	Examples of processes	Conclusion	

The process which I invented at Santa Fe Institute

Let a process $(X_i)_{i \in \mathbb{Z}}$ have the form

$$X_i := (K_i, Z_{K_i}),$$

where $(K_i)_{i\in\mathbb{Z}}$ and $(Z_k)_{k\in\mathbb{N}}$ are independent IID processes,

$$\begin{split} \mathsf{P}(\mathsf{K}_{\mathsf{i}} = \mathsf{k}) &= \mathsf{k}^{-1/\beta} / \zeta(\beta^{-1}), \qquad \beta \in (0,1), \\ \mathsf{P}(\mathsf{Z}_{\mathsf{k}} = \mathsf{z}) &= \frac{1}{2}, \qquad \qquad \mathsf{z} \in \{0,1\}. \end{split}$$

A linguistic interpretation

Process $(X_i)_{i\in\mathbb{Z}}$ is a sequence of random statements consistently describing the state of an "earlier drawn" random object $(Z_k)_{k\in\mathbb{N}}$. $X_i = (k, z)$ asserts that the k-th bit of $(Z_k)_{k\in\mathbb{N}}$ has value $Z_k = z$.

- We have card $U_{\delta}(n) \geq An^{\beta}$.
- Unfortunately, the alphabet $\mathbb{X}=\mathbb{N}\times\{0,1\}$ is infinite.

Problem statement	0000	00000	0000
Stationary (va	riable-length) co	oding of this process	

A function $f:\mathbb{X}\to\mathbb{Y}^+$ is extended to a function $f^{\mathbb{Z}}:\mathbb{X}^{\mathbb{Z}}\to\mathbb{Y}^{\mathbb{Z}},$

$$f^{\mathbb{Z}}((x_i)_{i\in\mathbb{Z}}):=...f(x_{-1})f(x_0).f(x_1)f(x_2)..., \qquad x_i\in\mathbb{X}.$$

For a measure u on $(\mathbb{Y}^{\mathbb{Z}},\mathfrak{Y}^{\mathbb{Z}})$ we define its stationary mean $ar{
u}$ as

$$\overline{\nu}(\mathsf{A}) = \lim_{\mathsf{n} \to \infty} \frac{1}{\mathsf{n}} \sum_{\mathsf{i}=0}^{\mathsf{n}-1} \nu \circ \mathsf{T}^{-\mathsf{i}}(\mathsf{A}),$$

where $T((y_i)_{i\in\mathbb{Z}}):=(y_{i+1})_{i\in\mathbb{Z}}$ is the shift.

Theorem 4

Let $\mu = P((X_i)_{i \in \mathbb{Z}} \in \cdot)$ for the process from the previous slide. Put $\mathbb{Y} = \{0, 1, 2\}$ and f(k, z) := b(k)z2, where $1b(k) \in \{0, 1\}^+$ is the binary expansion of k. A process with measure $\mu \circ (f^{\mathbb{Z}})^{-1}$ satisfies the hypothesis of Th. 1 for $\beta > 0.78$.

Problem statement	Sketch of the proof	Examples of processes 0000●	Conclusion 0000
A mixing process			

Let a process $(X_i)_{i \in \mathbb{Z}}$ have the form

$$X_i := (K_i, Z_{i,K_i}),$$

where $(K_i)_{i\in\mathbb{Z}}$ and $(Z_{ik})_{i\in\mathbb{Z}}$, $k\in\mathbb{N}$, are independent,

$$\mathsf{P}(\mathsf{K}_{\mathsf{i}}=\mathsf{k})=\mathsf{k}^{-1/\beta}/\zeta(\beta^{-1}), \qquad (\mathsf{K}_{\mathsf{i}})_{\mathsf{i}\in\mathbb{Z}}\sim\mathsf{IID},$$

whereas $(Z_{ik})_{i \in \mathbb{Z}}$ are Markov chains with

$$\begin{split} \mathsf{P}(\mathsf{Z}_{ik}=z) &= \frac{1}{2},\\ \mathsf{P}(\mathsf{Z}_{ik}=z|\mathsf{Z}_{i-1,k}=z) &= 1-p_k. \end{split}$$

Object $(Z_{ik})_{k \in \mathbb{N}}$ described by text $(X_i)_{i \in \mathbb{Z}}$ is a function of time i.

- We have $\liminf_{n\to\infty}\mathsf{E}(n)/n^\beta>0$ for $p_k\leq\mathsf{P}(\mathsf{K}_i=\mathsf{k}).$
- The stationary coding of this process is an ergodic process and also satisfies $\liminf_{n\to\infty} E(n)/n^{\beta} > 0$.

Problem statement	Sketch of the proof	Examples of processes	Conclusion
			0000

1 Problem statement

- 2 Sketch of the proof
- 3 Examples of processes



Problem statement	Sketch of the proof 0000	Examples of processes	Conclusion 0●00
Can we check wh	nich explanation is	better?	

Monkey-typing explanation

Zipf's and Herdan's law are observed if the letters and spaces in the text are obtained by pressing keys at random.

VS.

New explanation

If a text of length **n** describes $\geq \mathbf{n}^{\beta}$ independent facts in a repetitive way then the text contains $\geq \mathbf{n}^{\beta} / \log \mathbf{n}$ distinct words.

Can we estimate mutual information well?				
Problem statement	Sketch of the proof	Examples of processes	Conclusion 0000	

- Can we strengthen Theorems 1, 2, and 3?
 - Consider asymptotically mean stationary (AMS) processes.
 - Infer almost sure growth of vocabulary.
 - $\bullet~ \mbox{Replace}~ \mbox{lim} \mbox{sup}_{n \to \infty} \mbox{ with } \mbox{lim} \mbox{inf}_{n \to \infty}.$
- Let C(u) be the shortest program that generates u.

Then $E^{C}(n)$ is the algorithmic information between blocks.

- Let $(\omega_k)_{k\in\mathbb{N}}$ be an algorithmically random real in (0, 1). Mind that $\mathsf{E}(\mathsf{n}) = 0$ but $\mathsf{E}^{\mathsf{C}}(\mathsf{n}) \asymp \mathsf{n}^{\beta}$ for $\mathsf{X}_i := (\mathsf{K}_i, \omega_{\mathsf{K}_i})$.
- Can we use some universal codes to distinguish between some AMS sources with little vs. large $E^{C}(n)$?
- Can we use vocabulary of grammar-based codes to distinguish between some AMS sources with little vs. large E^C(n)?
- O there exist admissibly minimal codes that are computable in polynomial time? (Or sufficiently similar codes?)
 - Let $(X_i)_{i \in \mathbb{Z}}$ be a binary IID process. Then $\mathbb{V}[\Gamma(X_{1:n}] = \Omega\left(\sqrt{\frac{hn}{\log n}}\right)$ for irredicible grammar transforms.

Problem statement	Sketch of the proof	Examples of processes	Conclusion
Mv work			

- Ł. Dębowski, (2012). Mixing, Ergodic, and Nonergodic Processes with Rapidly Growing Information between Blocks. IEEE Transactions on Information Theory, vol. 58, pp. 3392-3401.
- Ł. Dębowski, (2011). On the Vocabulary of Grammar-Based Codes and the Logical Consistency of Texts. IEEE Transactions on Information Theory, vol. 57, pp. 4589–4599.
- Ł. Dębowski, (2010). Variable-Length Coding of Two-Sided Asymptotically Mean Stationary Measures. Journal of Theoretical Probability, 23:237–256.
- Ł. Dębowski, (2007). Menzerath's law for the smallest grammars. In: P. Grzybek, R. Koehler, eds., Exact Methods in the Study of Language and Text. Mouton de Gruyter. (77–85)

www.ipipan.waw.pl/~ldebowsk