

# Ograniczenie rozmiaru słownika uniwersalnych kodów gramatycznych

Łukasz Dębowski  
ldebowsk@ipipan.waw.pl

Instytut Podstaw Informatyki PAN

## Lingwistyczna motywacja problemu

Rozpatrujemy teksty w języku naturalnym (np. w j. polskim):

- **A** — liczba różnych słów w tekście,
- **N** — długość tekstu w słowach.

Wzór empiryczny („**prawo Guirauda-Herdana-Heapsa**”):

$$A \propto N^\alpha,$$

gdzie  $\alpha$  waha się między **0.5** a **1** w zależności od zbioru tekstów.

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

# Przełożenie na problem z teorii informacji

Słowa w tekście to w znacznej części powtarzające się napisy:

- liczba różnych słów  $\Rightarrow$  liczba **napisów powtarzających się**.

Modelowanie probabilistyczne tekstu wg hipotezy Hilberga (1990):

- tekst w języku naturalnym  $\Rightarrow$  **tekst wylosowany** z rozkładu stacjonarnego scharakteryzowanego parametrem globalnej zależności, którym jest **informacja wzajemna**.

## Problem

Ograniczyć w terminach **informacji wzajemnej** oczekiwaną liczbę **powtórzeń** w tekście wylosowanym z rozkładu stacjonarnego.

## Metoda

Wykorzystać problem **najmniejszej gramatyki** (CFG) dla tekstu.

- 1 Wstęp
- 2 Zasadniczy pomysł
- 3 Komplikacje dla kodowania uniwersalnego
- 4 Związki z modelowaniem j. naturalnego
- 5 Bibliografia

- 1 Wstęp
- 2 Zasadniczy pomysł
- 3 Komplikacje dla kodowania uniwersalnego
- 4 Związki z modelowaniem j. naturalnego
- 5 Bibliografia

# Gramatyka bezkontekstowa generująca jeden tekst

$$G = \left\{ \begin{array}{l} A_1 \rightarrow A_2 A_2 A_4 A_9 A_4 A_3 A_7 A_5 \\ A_2 \rightarrow A_6 A_9 A_6 A_3 \\ A_3 \rightarrow A_9 A_7 A_8, A_5 \\ A_4 \rightarrow \text{Jeszcze\_raz} \\ A_5 \rightarrow A_8 \text{\_nam\_} \\ A_6 \rightarrow \text{Sto\_lat} \\ A_7 \rightarrow \text{Niech} \\ A_8 \rightarrow \text{\_zyje} \\ A_9 \rightarrow \text{!\_} \end{array} \right\}$$

*Sto lat! Sto lat! Niech żyje, żyje nam.*

*Sto lat! Sto lat! Niech żyje, żyje nam.*

*Jeszcze raz! Jeszcze raz! Niech żyje, żyje nam.*

*Niech żyje nam.*

# Rozmiar słownika i długość gramatyki

Skrócony zapis gramatyki:

$$\mathbf{G} = \left\{ \begin{array}{l} \mathbf{A}_1 \rightarrow \alpha_1, \\ \mathbf{A}_2 \rightarrow \alpha_2, \\ \dots, \\ \mathbf{A}_n \rightarrow \alpha_n \end{array} \right\}, \quad \begin{array}{l} \forall[\mathbf{G}] := n \quad (\text{rozmiar słownika}), \\ |\mathbf{G}| := \sum_i |\alpha_i| \quad (\text{długość gramatyki}). \end{array}$$

Symbol  $\mathbf{A}_1$  ustalony jako symbol startowy.

**Najmniejsza gramatyka** dla tekstu to  $\mathbf{G}$  o najmniejszym  $|\mathbf{G}|$ , generująca ten tekst jako jedyną produkcję.

- Moses Charikar, Eric Lehman, ..., Abhi Shelat (2005),
- Wojciech Rytter (2003),
- John C. Kieffer, Enhui Yang (2000).

# Przybliżenie algorytmicznej informacji wzajemnej

Minimalna transformacja gramatykowa:

$$\Gamma : \mathbf{w} \in \mathbb{X}^* \rightarrow \Gamma(\mathbf{w}),$$

gdzie  $\Gamma(\mathbf{w})$  to najmniejsza gramatyka dla  $\mathbf{w}$ .

$|\Gamma(\mathbf{w})|$  — przybliżenie **złożoności algorytmicznej** napisu  $\mathbf{w}$ .

Przybliżenie **algorytmicznej informacji wzajemnej** między  $\mathbf{u}$  i  $\mathbf{v}$ :

$$\underbrace{|\Gamma(\mathbf{u})| + |\Gamma(\mathbf{v})| - |\Gamma(\mathbf{uv})|}_{\text{nadwyżka długości gramatyki}}.$$



# Ograniczenie rozmiaru słownika (Dębowski 2006)

Dla **minimalnej transformacji gramatycznej** zachodzi nierówność:

$$0 \leq |\Gamma(\mathbf{u})| + |\Gamma(\mathbf{v})| - |\Gamma(\mathbf{w})| \leq \mathbb{V}[\Gamma(\mathbf{w})]\mathbb{L}(\mathbf{w}), \quad (1)$$

gdzie  $\mathbf{w} = \mathbf{uv}$  zaś

$$\mathbb{L}(\mathbf{w}) := \max_{s,x,y,z \in \mathbb{X}^* : \mathbf{w} = xsysz} |s|,$$

to długość najdłuższego **powtórzenia** w napisie  $\mathbf{w}$ .

W dalszej części interesować nas będą uogólnienia nierówności (1) w przypadku kompresji **procesu stochastycznego**.

- 1 Wstęp
- 2 Zasadniczy pomysł
- 3 Komplikacje dla kodowania uniwersalnego**
- 4 Związki z modelowaniem j. naturalnego
- 5 Bibliografia

# Stacjonarny proces stochastyczny

Ciąg zmiennych losowych:

$$(\dots, \mathbf{X}_{-1}, \mathbf{X}_0, \mathbf{X}_1, \mathbf{X}_2, \dots), \quad \mathbf{X}_i : \Omega \rightarrow \mathbb{X}.$$

Rozkład prawdopodobieństwa niezmienniczy względem przesunięć:

$$\mathbb{P}(\mathbf{X}_{i+1} = \mathbf{x}_1, \mathbf{X}_{i+2} = \mathbf{x}_2, \dots, \mathbf{X}_{i+n} = \mathbf{x}_n) = \mathbb{P}(\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_n),$$

gdzie  $\mathbb{P} : \mathbb{X}^* \rightarrow \mathbb{R}$ .

# Kody jednoznacznie dekodowalne

Kod  $\mathbf{C} : \mathbb{X}^* \rightarrow \mathbb{X}^*$  nazywa się **jednoznacznie dekodowalnym**,  
gdy jego **rozszerzenie** (na skończone konkatenacje)  
 $\mathbf{C}^* : (\mathbb{X}^*)^* \ni (\mathbf{u}_1, \dots, \mathbf{u}_k) \mapsto \mathbf{C}(\mathbf{u}_1)\dots\mathbf{C}(\mathbf{u}_k) \in \mathbb{X}^*$   
jest przekształceniem wzajemnie jednoznacznym.

Kod jest **j.d.**, jeżeli jest **bezprzedrostkowy**, tzn. żaden napis z jego obrazu nie jest przedrostkiem innego napisu z obrazu.

Dla dowolnego **j.d.** kodu  $\mathbf{C}$  i każdego rozkładu p-stwa  $\mathbb{P}$  zachodzi

$$\sum_{\mathbf{u} \in \mathbb{X}^{*n}} \mathbb{P}(\mathbf{u}) |\mathbf{C}(\mathbf{u})| \geq \mathbf{H}(n) := - \sum_{\mathbf{u} \in \mathbb{X}^n} \mathbb{P}(\mathbf{u}) \log_{|\mathbb{X}|} \mathbb{P}(\mathbf{u}) \quad (\text{entropia}).$$

# Entropia a kodowanie arytmetyczne

Dla każdego dyskretnego rozkładu  $p$ -stwa  $\mathbb{P} : \mathbb{X}^n \rightarrow \mathbb{R}$  istnieje **bezbudrostkowy** kod  $\mathbf{C} : \mathbb{X}^n \rightarrow \mathbb{X}^*$  taki, że

$$|\mathbf{C}(\mathbf{u})| = \left\lceil -\log_{|\mathbb{X}|} \mathbb{P}(\mathbf{u}) \right\rceil + 1.$$

Dla tegoż kodu mamy

$$0 \leq \sum_{\mathbf{u} \in \mathbb{X}^n} \mathbb{P}(\mathbf{u}) |\mathbf{C}(\mathbf{u})| - \mathbf{H}(\mathbf{n}) \leq 2.$$

Kod **j.d.** nazywa się **uniwersalnym**, gdy asymptotycznie optymalnie kompresuje on **dowolny** proces stacjonarny.

# Kody uniwersalne

Ponieważ  $\mathbf{H}(\mathbf{n})/n$  jest nieujemne i maleje z  $\mathbf{n}$ , istnieje granica

$$\mathbf{h} = \lim_{n \rightarrow \infty} \mathbf{H}(\mathbf{n})/n \quad (\text{intensywność entropii}).$$

Mówimy, że kod **j.d.**  $\mathbf{C}$  jest **uniwersalny**, gdy spełnia

$$\lim_{n \rightarrow \infty} \sum_{\mathbf{u} \in \mathbb{X}^n} \mathbb{P}(\mathbf{u}) |\mathbf{C}(\mathbf{u})| / n = \mathbf{h}$$

dla dowolnego rozkładu stacjonarnego  $\mathbb{P}$ .

Przykładem kodu **uniwersalnego** jest kod Lempel-Ziva.

# Informacja wzajemna (Dębowski 2006)

$$E(n) = 2H(n) - H(2n) \quad (\text{entropia nadwyżkowa}),$$
$$E^C(n) = \sum_{u,v \in \mathbb{X}^n} \mathbb{P}(uv) \underbrace{[|C(u)| + |C(v)| - |C(uv)|]}_{\text{nadwyżka długości kodu}}.$$

Jeżeli  $C$  jest **uniwersalny**, to

$$\limsup_{n \rightarrow \infty} [E^C(n) - E(n)] \geq 0.$$

Jeśli  $C$  i  $C'$  są **uniwersalne** oraz  $|C(\cdot)| \geq |C'(\cdot)|$ , to

$$\limsup_{n \rightarrow \infty} [E^C(n) - E^{C'}(n)] \geq 0.$$

# Transformacja gramatyczna a kod uniwersalny

Nie istnieje kod uniwersalny o długości równej długości **najmniejszej gramatyki** dla napisu  $u$ , gdyż ta ma długość

$$O(|u| / \log |u|).$$

## Intuicja

Nie można zakodować nieskończonej liczby symboli nieterminalnych za pomocą pojedynczych symboli (lub słów ustalonej długości) nad skończonym alfabetem.

Mimo to istnieje kod uniwersalny,

- który koduje gramatykę **minimalną w uogólnionym sensie**,
- i którego nadwyżka długości jest mniejsza od iloczynu **rozmiaru słownika i długości najdłuższego powtórzenia**.



# Konstrukcja kodu

$\mathbf{B} : \mathcal{G} \rightarrow \mathbb{X}^+$  nazywam **lokalnym koderem gramatycznym**, gdy

$$\mathbf{B}(\mathbf{G}) = \mathbf{B}_S(\mathbf{B}_N(\mathbf{G})),$$

gdzie  $\mathbf{B}_N : \mathcal{G} \rightarrow (\{0\} \cup \mathbb{N})^*$  koduje **gramatykę**

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

jako **ciąg liczb**  $\mathbf{B}_N(\mathbf{G}) = \mathbf{F}_1(\alpha_1)\mathbf{D}\mathbf{F}_2(\alpha_2)\mathbf{D}\dots\mathbf{D}\mathbf{F}_n(\alpha_n)(\mathbf{D} + 1)$ ,

$$\mathbf{F}_i(\mathbf{x}) := \mathbf{x}, \quad \mathbf{x} \in \mathbb{X} = \{0, 1, \dots, \mathbf{D} - 1\},$$

$$\mathbf{F}_i(\mathbf{A}_j) := \mathbf{D} + 1 + j - i, \quad \mathbf{F}_i(\beta\gamma) = \mathbf{F}_i(\beta)\mathbf{F}_i(\gamma),$$

zaś  $\mathbf{B}_S : (\{0\} \cup \mathbb{N})^* \rightarrow \mathbb{X}^*$  spełnia  $\mathbf{B}_S(\beta\gamma) = \mathbf{B}_S(\beta)\mathbf{B}_S(\gamma)$ .

# Uogólnienie minimalnej transformacji gramatykowej

Dla kodera gramatykowego  $\mathbf{B}$ , transformację gramatykową  $\Gamma$  oraz kod  $\mathbf{B}(\Gamma(\cdot))$  nazywam **B-minimalnymi**, jeżeli

$$|\mathbf{B}(\Gamma(\mathbf{w}))| \leq |\mathbf{B}(\mathbf{G})|$$

dla każdej gramatyki  $\mathbf{G}$  generującej słowo  $\mathbf{w}$ .

## Twierdzenie (Dębowski 2007)

Niech  $\mathbf{B}_S$  będzie koderem **bezprzedrostkowym** zaś  $|\mathbf{B}_S(\cdot)|$  będzie funkcją rosnącą i spełnia

$$\limsup_{n \rightarrow \infty} |\mathbf{B}_S(\mathbf{n})| / \log_D n = 1.$$

Jeżeli kod  $\mathbf{C}$  jest  $\mathbf{B}_S(\mathbf{B}_N(\cdot))$ -minimalny, to jest **uniwersalny**.

# Główne twierdzenie

## Twierdzenie (Dębowski 2007)

Dla lokalnego kodera gramatykowego  $\mathbf{B} = \mathbf{B}_S(\mathbf{B}_N(\cdot))$ , oznaczmy

$$\mathbf{W}_m := \max_{0 \leq n \leq D+2+m} |\mathbf{B}_S(n)|.$$

Jeżeli kod  $\mathbf{C} = \mathbf{B}(\Gamma(\cdot))$  jest  $\mathbf{B}$ -minimalny, to

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

gdzie  $\mathbf{w} = \mathbf{uv}$  dla  $\mathbf{u}, \mathbf{v} \in \mathbb{X}^+$ .

(— Tudzież kod ten ma kilka innych przyjemnych własności.)

- 1 Wstęp
- 2 Zasadniczy pomysł
- 3 Komplikacje dla kodowania uniwersalnego
- 4 Związki z modelowaniem j. naturalnego**
- 5 Bibliografia

# Entropia języka naturalnego

Spróbujmy wyobrazić sobie proces tworzenia tekstów w **języku naturalnym** jako proces stacjonarny  $(\mathbf{X}_k)_{k \in \mathbb{Z}}$ :

- $\mathbf{X}_i$  — pojedyncze **litery** tekstu ( $\mathbb{X}$  skończone),
- $(\mathbf{X}_i)_{m \leq i \leq n}$  — konkretne **teksty** tworzone przez ludzi.

W oparciu o pomiary entropii warunkowej (Shannon 1950) sformułowano hipotezę (Hilberg 1990), że proces ten ma entropię nadwyżkową  $\mathbf{E}(\mathbf{n}) \asymp \sqrt{\mathbf{n}}$ .

- W. Hilberg, (1990). *Frequenz*, 44:243–248.
- Ł. Dębowski, (2006). *On Hilberg's law and its links with Guiraud's law*. *Journal of Quantitative Linguistics*, 13:81–109.

# Hipoteza

Jeżeli tekst długości  $N$  liter opisuje w sposób **niesprzeczny**  
 $N^\beta$  **losowych niezależnych faktów** o świecie,  
a proces generowania tekstu jest dostatecznie losowy,  
to tekst ten zawiera co najmniej  $N^\beta / \log N$  **różnych słów**  
(rozumianych jako symbole nieterminalne najkrótszej gramatyki).

# Pewien wyidealizowany proces stochastyczny

Rozważmy proces  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  o zmiennych losowych

$$\mathbf{X}_i := (\mathbf{K}_i, \mathbf{Z}_{\mathbf{K}_i}),$$

gdzie niezależne procesy IID  $(\mathbf{K}_i)_{i \in \mathbb{Z}}$  oraz  $(\mathbf{Z}_k)_{k \in \mathbb{N}}$  spełniają  $\mathbf{P}(\mathbf{K}_i = \mathbf{k}) > 0$  dla  $\mathbf{k} \in \mathbb{N}$  oraz  $\mathbf{P}(\mathbf{Z}_k = \mathbf{z}) = \frac{1}{2}$  dla  $\mathbf{z} \in \{0, 1\}$ .

## Interpretacja quasi-lingwistyczna

Proces  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  jest ciągiem losowych **stwierdzeń**  $\mathbf{X}_i$  **niesprzecznie opisujących** stan „wcześniej” wylosowanego obiektu  $(\mathbf{Z}_k)_{k \in \mathbb{N}}$  przyjmującego **nieprzeliczalnie** wiele wartości.

Stwierdzenie  $\mathbf{X}_i = (\mathbf{k}, \mathbf{z})$  orzeka, że  $\mathbf{k}$ -ty bit obiektu  $(\mathbf{Z}_k)_{k \in \mathbb{N}}$  ma wartość  $\mathbf{z}$ , w taki sposób, że można ustalić zarówno  $\mathbf{k}$  jak  $\mathbf{z}$ .

Dla stwierdzeń  $\mathbf{X}_i = (\mathbf{k}, \mathbf{z})$  i  $\mathbf{X}_j = (\mathbf{k}', \mathbf{z}')$  nie wiadomo, które bity opiszą i jakie wartości im przypiszą, ale jeżeli  $\mathbf{k} = \mathbf{k}'$ , to  $\mathbf{z} = \mathbf{z}'$ .

## Dolne ograniczenie entropii nadwyżkowej

$E(n)$  można ograniczyć przez liczbę bitów procesu  $(Z_k)_{k \in \mathbb{N}}$ , które można przewidzieć dostatecznie dobrze na podstawie  $X_{1:n}$ :

$$H_\delta^U(n) := |\{k \in \mathbb{N} : P(f_{nk}((X_i)_{1 \leq i \leq n}) = Z_k) \geq \delta\}|.$$

Twierdzenie (Dębowski 2007)

Dla  $H^I(n) := I((X_i)_{1 \leq i \leq n}; (Z_k)_{k \in \mathbb{N}})$  mamy  $\lim_n H^I(n)/n = 0$  oraz

$$E(n) \geq 2H^I(n) - H^I(2n),$$

$$H^I(n) \geq H_\delta^U(n) \cdot \frac{\log 2 + \delta \log \delta + (1-\delta) \log(1-\delta)}{\log |\mathbb{X}|}.$$

Zatem jeżeli  $H_\delta^U(n) \geq B \cdot n^\beta$  (np. dla  $P(K_i = k) \sim k^{-1/\beta}$ ), to

$$\limsup_{n \rightarrow \infty} [E(n) - B' \cdot n^\beta] \geq 0.$$



# Procesy nieprzeliczalnego opisu

Twierdzenie z poprzedniego slajdu obowiązuje też dla **p.n.o.**

Proces stacjonarny  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  nazywam **p.n.o.**, jeżeli istnieją

- zmienne  $(\mathbf{Z}_k)_{k \in \mathbb{N}} \sim \text{IID}$ ,  $\mathbf{P}(\mathbf{Z}_k = \mathbf{z}) = \frac{1}{2}$  dla  $\mathbf{z} \in \{\mathbf{0}, \mathbf{1}\}$ ,
- i funkcje  $\mathbf{f}_{nk} : \mathbb{X}^n \rightarrow \{\mathbf{0}, \mathbf{1}\}$

takie, że

- $\mathbf{P}(\mathbf{f}_{nk}((\mathbf{X}_i)_{j+1 \leq i \leq j+n}) = \mathbf{Z}_k)$  nie zależy od  $\mathbf{j} \in \mathbb{Z}$ ,
- $\lim_n \mathbf{P}(\mathbf{f}_{nk}((\mathbf{X}_i)_{j+1 \leq i \leq j+n}) = \mathbf{Z}_k) = 1$ .

Są to procesy nieergodyczne.

# Procesy o skończonej energii

Proces stacjonarny  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  nazywa się **p. o s.e.**, jeżeli

$$\mathbb{P}(\mathbf{u}\mathbf{v})/\mathbb{P}(\mathbf{u}) \leq K \exp[-c|\mathbf{v}|], \quad \mathbf{u}, \mathbf{v} \in \mathbb{X}^+.$$

Przykładem **p. o s.e.** jest ciąg zmiennych

$$\mathbf{X}_i \equiv \mathbf{Y}_i + \mathbf{U}_i \pmod{D},$$

gdzie proces  $(\mathbf{Y}_i)_{i \in \mathbb{Z}}$  jest ergodyczny zaś  $(\mathbf{U}_i)_{i \in \mathbb{Z}} \sim \text{IID}$ .

**Twierdzenie (Shields 1997)**

Jeżeli  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  jest **p. o s.e.**, to długość najdłuższego powtórzenia jest ograniczona nierównością

$$\liminf_{n \rightarrow \infty} [K' \log n - \mathbb{L}((\mathbf{X}_i)_{1 \leq i \leq n})] \geq 0 \quad \text{prawie na pewno.}$$

# Składając wszystko razem

## Twierdzenie

Jeżeli kod  $\mathbf{C} = \mathbf{B}(\Gamma(\cdot))$  jest  $\mathbf{B}$ -minimalny dla lokalnego kodera  $\mathbf{B}$ , zaś  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  jest p. o s.e., to

$$\limsup_{n \rightarrow \infty} \left[ \sum_{\mathbf{w} \in \mathbb{X}^n} \mathbb{P}(\mathbf{w}) \mathbb{V}[\Gamma(\mathbf{w})] - \frac{\mathbf{E}(n)}{K' \log n} \right] \geq 0.$$

Jeżeli  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  jest p.n.o. takim, że

$$|\{ \mathbf{k} \in \mathbb{N} : \mathbb{P}(\mathbf{f}_{nk}((\mathbf{X}_i)_{1 \leq i \leq n}) = \mathbf{Z}_k) \geq \delta \}| \geq \mathbf{B} \cdot n^\beta,$$

to

$$\limsup_{n \rightarrow \infty} \left[ \mathbf{E}(n) - \mathbf{B}' \cdot n^\beta \right] \geq 0.$$

Niestety  $\limsup_n (a_n + b_n) \geq \limsup_n a_n + \liminf_n b_n$ .

- 1 Wstęp
- 2 Zasadniczy pomysł
- 3 Komplikacje dla kodowania uniwersalnego
- 4 Związki z modelowaniem j. naturalnego
- 5 Bibliografia**

## Moje prace

- Ł. Dębowski, (2007). *On vocabulary size of grammar-based codes*. <http://xxx.lanl.gov/abs/cs.IT/0701047>.  
(ISIT 2007, Nicea)
- Ł. Dębowski, (2006). *Ergodic decomposition of excess entropy and conditional mutual information*. Prace IPI PAN nr 993.
- Ł. Dębowski, (2006). *On Hilberg's law and its links with Guiraud's law*. Journal of Quantitative Linguistics, 13:81–109.

[www.ipipan.waw.pl/~ldebowsk](http://www.ipipan.waw.pl/~ldebowsk)