

W analogii do języka naturalnego...
czyli twierdzenie o faktach i słowach
dla procesów stacjonarnych

Łukasz Dębowski
ldebowsk@ipipan.waw.pl



Institut Podstaw Informatyki
Polskiej Akademii Nauk

LVIII Szkoła Matematyki Poglądowej
24–28.08.2018

Punkt wyjścia

Teoria informacji:

- **Podstawowy problem:**

Jak zakodować dany napis (ciąg znaków) przy pomocy jak najmniejszej liczby cyfr binarnych (zer lub jedynek)?

- **Zastosowania:**

— kompresja i przesyłanie danych, teoretyczne podstawy informatyki, statystyki, a nawet matematyki jako takiej.

Statystyczne modelowanie języka naturalnego:

- **Podstawowy problem:**

Jakie przypisać prawdopodobieństwo dowolnym wypowiedziom w danym języku naturalnym (angielskim, polskim, chińskim)?

- **Zastosowania:**

— automatyczne rozpoznawanie mowy, klawiatury telefonów komórkowych, maszynowe tłumaczenie.

Temat wystąpienia

Twierdzenie o faktach i słowach (sformułowanie nieformalne)

Liczba **niezależnych** faktów opisywanych przez skończony tekst jest z grubsza mniejsza niż liczba **różnych** słów w tymże tekście.

W analogii do języka naturalnego, pojęcie faktu i słowa można zdefiniować dla dość dowolnego procesu stochastycznego.

Ku sformułowaniu formalnemu

Na dalszych slajdach zdefiniujemy:

- $\mathbb{U}(x_1^n)$ — zbiór faktów przewidywalnych z napisu x_1^n .
- $G_{\text{PPM}}(x_1^n)$ — rząd PPM napisu x_1^n .
- $V_{\text{PPM}}(x_1^n)$ — zbiór słów PPM napisu x_1^n .
- Wykładnik Hilberga: $\mathop{\text{hilb}}_{n \rightarrow \infty} n^\beta = \beta$.

Twierdzenie o faktach i słowach

Dla dowolnego procesu stacjonarnego $(X_i)_{i=1}^\infty$ o skończonym alfabecie zachodzi nierówność

$$\mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} \text{ card } \mathbb{U}(X_1^n) \leq \mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} [G_{\text{PPM}}(X_1^n) + \text{card } V_{\text{PPM}}(X_1^n)].$$

- 1 Wprowadzenie
- 2 Jak sformalizować pojęcie słowa?
- 3 Jak sformalizować pojęcie faktu?
- 4 Jak udowodnić twierdzenie o faktach i słowach?
- 5 Konkluzje

Jak zdefiniować słowo w dowolnym tekście?

- W przypadku tekstów w wielu językach naturalnych, słowo to ciąg liter od spacji do spacji:

Все люди рождаются свободными и равными в своем достоинстве и правах. Они наделены разумом и совестью и должны поступать в отношении друг друга в духе братства.

- Co począć z japońskim?

すべての人間は、生まれながらにして自由であり、かつ、尊厳と権利とについて平等である。人間は、理性と良心とを授けられており、互いに同胞の精神をもって行動しなければならない。

- Albo z ciągiem zer i jedynek?

$\Omega = 0.000000100000010000011000100001101000111\dots$

Podęście toporne, acz skuteczne

- **Notacja:** napis $x_j^k = x_j x_{j+1} \dots x_k$, gdzie x_i to symbole.
- Zbiór podstów długości m w napisie x_1^n to

$$V(m|x_1^n) := \{x_{t+1}^{t+m} : 0 \leq t \leq n - m\}.$$

- Niech $G_{opt}(x_1^n)$ będzie **optymalną długością** podstawa dla x_1^n .
- Zbiór **optymalnych słów** w napisie x_1^n to

$$V_{opt}(x_1^n) := V(G_{opt}(x_1^n)|x_1^n).$$

Jak wybrać sensowne $G_{opt}(x_1^n)$?

Rozkłady prawdopodobieństw PPM

Sugestii, jak wybrać sensowną optymalną długość podstawa, udziela konstrukcja rozkładów prawdopodobieństw PPM.

- Rozkłady PPM (Prediction by Partial Matching) to pewna rodzina rozkładów $\mathbf{PPM}_k(x_1^n)$, gdzie $k = -1, 0, 1, \dots$
- Każdy rozkład $\mathbf{PPM}_k(x_1^n)$ to rozkład procesu Markowa rzędu k estymowany sukcesywnie na danym napisie x_1^n .
- Dla każdego napisu x_1^n istnieje rząd k , dla którego prawdopodobieństwo $\mathbf{PPM}_k(x_1^n)$ jest największe.

Rozkłady prawdopodobieństw PPM — definicja

Częstość podstawa w napisie to

$$N(\mathbf{w}_1^k | x_1^n) := \sum_{i=1}^{n-k+1} \mathbf{1} \{x_i^{i+k-1} = \mathbf{w}_1^k\},$$

gdzie $x_j^k := \lambda$ oraz $\sum_{i=j}^k f(i) := \mathbf{0}$ dla $k < j$.

Dla $x_i \in \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_D\}$ definiujemy

$$\text{PPM}_k(x_i | x_1^{i-1}) := \begin{cases} \frac{1}{D}, & k = -1, \\ \frac{N(x_{i-k}^i | x_1^{i-1}) + 1}{N(x_{i-k}^{i-1} | x_1^{i-2}) + D}, & k \geq 0, \end{cases}$$

$$\text{PPM}_k(x_1^n) := \prod_{i=1}^n \text{PPM}_k(x_i | x_1^{i-1}).$$

Maksymalne powtórzenie, rząd PPM i słownik PPM

- Długość **maksymalnego powtórzenia** w napisie x_1^n to

$$L(x_1^n) := \max \left\{ k : N(w_1^k | x_1^n) \geq 2 \text{ dla pewnego } w_1^k \right\}.$$

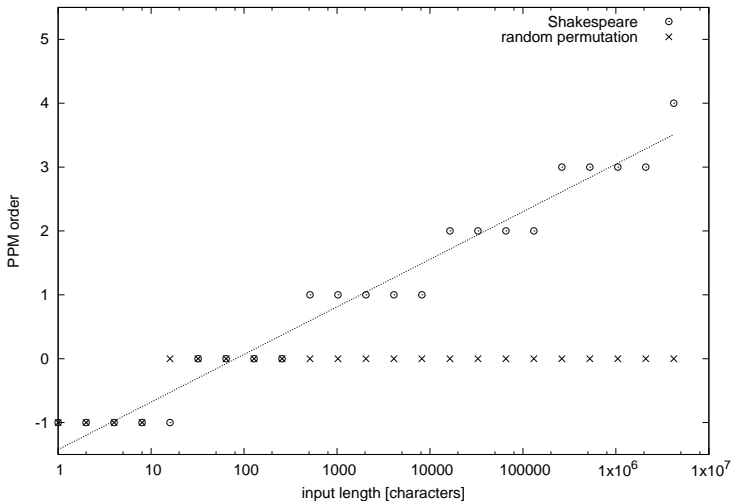
- Zauważmy, że $\text{PPM}_k(x_1^n) = D^{-n}$ dla $k > L(x_1^n)$.
- **Rząd PPM** $G_{\text{PPM}}(x_1^n)$ to najmniejsza liczba G taka, że

$$\text{PPM}_G(x_1^n) \geq \text{PPM}_k(x_1^n) \text{ dla każdego } k \geq -1.$$

- Mamy $G_{\text{PPM}}(x_1^n) \leq L(x_1^n) \leq n$.
- **Zbiór słów PPM** w napisie x_1^n to

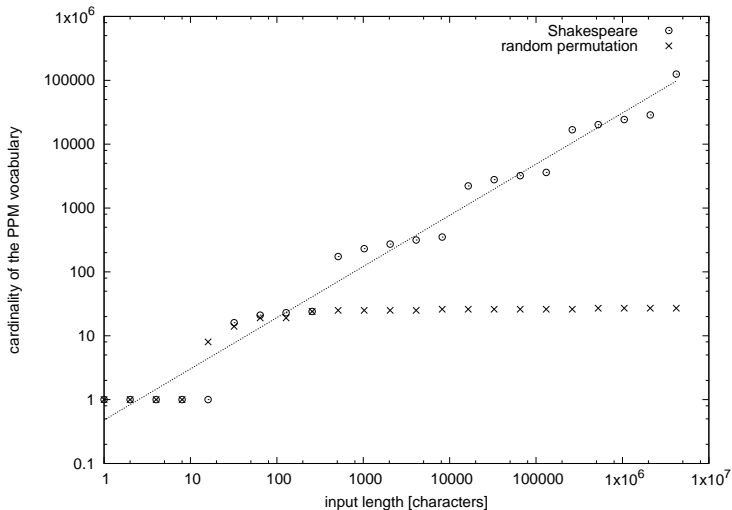
$$V_{\text{PPM}}(x_1^n) := V(G_{\text{PPM}}(x_1^n) | x_1^n).$$

Rząd PPM na wykresie



$$G_{\text{PPM}}(x_1^n) \approx -1.423 + 0.323 \ln n$$

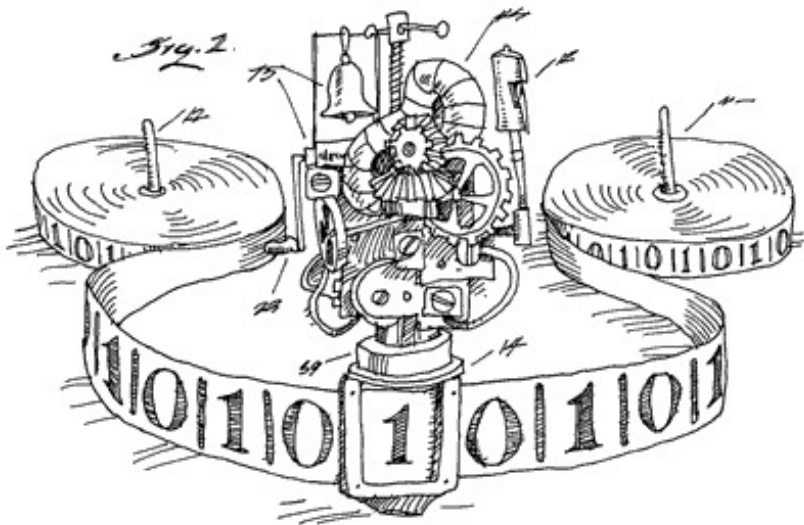
Liczba różnych słów PPM na wykresie



$$\text{card } V_{\text{PPM}}(x_1^n) \approx 1.300n^{0.801}$$

- 1 Wprowadzenie
- 2 Jak sformalizować pojęcie słowa?
- 3 Jak sformalizować pojęcie faktu?
- 4 Jak udowodnić twierdzenie o faktach i słowach?
- 5 Konkluzje

Maszyna Turinga (z dzwonkiem)



Algorytmiczna teoria informacji

- Mamy ustaloną bezprzedrostkową maszynę Turinga.
- **Złożoność Kolmogorowa** $\mathbb{H}(x_1^n)$ to długość najkrótszego programu generującego napis x_1^n :

$$\mathbb{H}(x_1^n) := \min \{ |p| : S(p) = x_1^n \},$$

gdzie $S(p)$ to wynik programu p .

(Funkcja $x_1^n \mapsto \mathbb{H}(x_1^n)$ nie jest obliczalna.)

- nieskończony ciąg $(x_i)_{i=1}^{\infty} = (x_1, x_2, x_3, \dots)$ jest nazywany **algorytmicznie losowym**, gdy istnieje stała $c > 0$ taka, że

$$\mathbb{H}(x_1^n) \geq n - c.$$

(Zachodzi to, gdy najkrótszy program ma postać **print** x_1^n .)

- **Ciąg rzutów uczciwą monetą** jest algorytmicznie losowy pnp.

Prawdopodobieństwo stopu

- **Prawdopodobieństwo stopu** Ω to liczba

$$\Omega = \sum_{p:S(p)\neq \perp} 2^{-|p|} \in (0, 1).$$

- Zdefiniujmy rozwinięcie binarne $(\Omega_k)_{k=1}^{\infty} = (\Omega_1, \Omega_2, \Omega_3, \dots)$, gdzie $\Omega_k \in \{0, 1\}$ oraz $\sum_{k=1}^{\infty} 2^{-k} \Omega_k = \Omega$.
- Ciąg $(\Omega_k)_{k=1}^{\infty}$ jest **ciągami algorytmicznie losowym**.
- Ponadto, gdybyśmy znali napis Ω_1^n , moglibyśmy odpowiedzieć na pytanie, które **stwierdzenia matematyczne** długości mniejszej od n są prawdziwe, a które nie.

Prawdopodobieństwo stopu Ω jest „kamieniem filozoficznym”.
Cyfry Ω_k są niezależnymi faktami matematycznymi.

Proces Santa Fe

- 1 Rozpatrzmy przestrzeń probabilistyczną $(\mathbb{J}, \mathcal{J}, P)$.
- 2 Niech $(K_i)_{i=1}^{\infty}$ będzie ciągiem niezależnych zmiennych losowych o wartościach w liczbach naturalnych, $K_i : \mathbb{J} \rightarrow \mathbb{N} = \{1, 2, 3, \dots\}$, i o rozkładzie

$$P(K_i = k) = \frac{k^{-\alpha}}{\zeta(\alpha)}, \quad \alpha > 1,$$

gdzie $\zeta(\alpha) = \sum_{k=1}^{\infty} k^{-\alpha}$.

- 3 **Proces Santa Fe** to ciąg zmiennych $(X_i)_{i=1}^{\infty}$ złożonych z par

$$X_i = (K_i, \Omega_{K_i}),$$

gdzie Ω_k to cyfry prawdopodobieństwa stopu.

Od procesu Santa Fe do prawdopodobieństwa stopu

- Zdefiniujemy funkcję:

$$g(k, x_1^n) = \begin{cases} 0, & \text{jeśli } \forall_{i \in \{1, \dots, n\}} (x_i = (k, z) \implies x_i = (k, 0)), \\ 1, & \text{jeśli } \forall_{i \in \{1, \dots, n\}} (x_i = (k, z) \implies x_i = (k, 1)), \\ 2, & \text{inaczej.} \end{cases}$$

- Funkcja ta jest obliczalna i dla procesu Santa Fe spełnia

$$\lim_{n \rightarrow \infty} g(k, X_{i+1}^{i+n}) = \Omega_k \text{ prawie na pewno.}$$

- Określmy zbiór **niezależnych faktów** przewidywalnych z x_1^n :

$$\mathbb{U}(x_1^n) := \{l \in \mathbb{N} : g(k, x_1^n) = \Omega_k \text{ dla wszystkich } k \leq l\}.$$

- Dla procesu Santa Fe mamy wzrost potęgowy

$$\limsup_{n \rightarrow \infty} \frac{\text{card } \mathbb{U}(X_1^n)}{n^{1/\alpha}} \in (0, \infty) \text{ prawie na pewno.}$$

Zbiór przewidywalnych faktów — w ogólności

- Niech $(z_k)_{k=1}^{\infty} = (z_1, z_2, z_3, \dots)$ będzie ustalonym ciągiem algorytmicznie losowym, niekoniecznie równym rozwinięciu binarnemu prawdopodobieństwa stopu.
- Symbole z_k nazywać będziemy **niezależnymi faktami**.
- Niech $g(k, x_1^n)$ będzie ustaloną funkcją obliczalną.
- Niech x_1^n będzie dowolnym napisem.
- Określmy zbiór **niezależnych faktów** przewidywalnych z x_1^n :

$$\mathbb{U}(x_1^n) := \{l \in \mathbb{N} : g(k, x_1^n) = z_k \text{ dla wszystkich } k \leq l\}.$$

- 1 Wprowadzenie
- 2 Jak sformalizować pojęcie słowa?
- 3 Jak sformalizować pojęcie faktu?
- 4 Jak udowodnić twierdzenie o faktach i słowach?**
- 5 Konkluzje

Wykładnik Hilberga

- Wykładnik Hilberga definiujemy jako

$$\mathbf{hilb} s(n) := \limsup_{n \rightarrow \infty} \frac{\log^+ s(n)}{\log n}, \quad \log^+ x = \begin{cases} \log(x+1), & x \geq 0, \\ 0, & x < 0. \end{cases}$$

Mamy $\mathbf{hilb} n^\beta = \beta$ dla $\beta \geq 0$.

- Dla procesu Santa Fe mamy prawie na pewno

$$\mathbf{hilb} \text{ card } \cup(X_1^n) = 1/\alpha \in (0, 1).$$

- Dla języka naturalnego mamy prawdopodobnie

$$\mathbf{hilb} \text{ card } V_{\text{PPM}}(x_1^n) \approx 0.8.$$

Twierdzenie o faktach i słowach

Na poprzednich slajdach zdefiniowaliśmy:

- $\mathbb{U}(x_1^n)$ — zbiór faktów przewidywalnych z napisu x_1^n .
- $G_{\text{PPM}}(x_1^n)$ — rząd PPM napisu x_1^n .
- $V_{\text{PPM}}(x_1^n)$ — zbiór słów PPM napisu x_1^n .
- Wykładnik Hilberga: $\mathop{\text{hilb}}_{n \rightarrow \infty} n^\beta = \beta$.

Twierdzenie o faktach i słowach

Dla dowolnego procesu stacjonarnego $(X_i)_{i=1}^\infty$ o skończonym alfabecie zachodzi nierówność

$$\mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} \text{ card } \mathbb{U}(X_1^n) \leq \mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} [G_{\text{PPM}}(X_1^n) + \text{card } V_{\text{PPM}}(X_1^n)].$$

Twierdzenie o wykładnikach Hilberga

Twierdzenie

Niech $J(n) := 2G(n) - G(2n)$. Jeżeli $\lim_{n \rightarrow \infty} G(n)/n = g$, to

$$\mathop{\text{hilb}}_{n \rightarrow \infty} [G(n) - ng] \leq \mathop{\text{hilb}}_{n \rightarrow \infty} J(n),$$

gdzie równość zachodzi, gdy $J(n) \geq 0$.

- Entropia Shannona $H(X_1^n) := \mathbb{E} [-\log P(X_1^n)]$.
- Intensywność entropii Shannona $h := \lim_{n \rightarrow \infty} H(X_1^n)/n$.
- Informacja wzajemna Shannona

$$\begin{aligned} I(X_1^n; X_{n+1}^{2n}) &:= H(X_1^n) + H(X_{n+1}^{2n}) - H(X_1^{2n}) \\ &= 2H(X_1^n) - H(X_1^{2n}) \geq 0. \end{aligned}$$

- Stąd $\mathop{\text{hilb}}_{n \rightarrow \infty} [H(X_1^n) - nh] = \mathop{\text{hilb}}_{n \rightarrow \infty} I(X_1^n; X_{n+1}^{2n})$.

Zarys dowodu twierdzenia o faktach i słowach

- Niech $\mathbb{H}_{\text{PPM}}(x_1^n) := -\log \left[\sum_{k=-1}^{\infty} \text{PPM}_k(x_1^n) / (k+2)^2 \right]$.
- Mamy równość intensywności entropii

$$h = \lim_{n \rightarrow \infty} \frac{\mathbb{E} \mathbb{H}(X_1^n)}{n} = \lim_{n \rightarrow \infty} \frac{\mathbb{E} \mathbb{H}_{\text{PPM}}(X_1^n)}{n}.$$

- Oznaczmy informację wzajemną PPM

$$\mathbb{I}_{\text{PPM}}(x_1^n; x_{n+1}^{2n}) := \mathbb{H}_{\text{PPM}}(x_1^n) + \mathbb{H}_{\text{PPM}}(x_{n+1}^{2n}) - \mathbb{H}_{\text{PPM}}(x_1^{2n}).$$

- Twierdzenie o faktach i słowach jest konsekwencją nierówności

$$\begin{aligned} \liminf_{n \rightarrow \infty} \mathbb{E} \text{card } \mathcal{U}(X_1^n) &\leq \liminf_{n \rightarrow \infty} [\mathbb{E} \mathbb{H}(X_1^n) - hn] \\ &\leq \liminf_{n \rightarrow \infty} [\mathbb{E} \mathbb{H}_{\text{PPM}}(X_1^n) - hn] \leq \liminf_{n \rightarrow \infty} \mathbb{E} \mathbb{I}_{\text{PPM}}(X_1^n; X_{n+1}^{2n}) \\ &\leq \liminf_{n \rightarrow \infty} \mathbb{E} [\mathbf{G}_{\text{PPM}}(X_1^n) + \text{card } \mathcal{V}_{\text{PPM}}(X_1^n)]. \end{aligned}$$

- 1 Wprowadzenie
- 2 Jak sformalizować pojęcie słowa?
- 3 Jak sformalizować pojęcie faktu?
- 4 Jak udowodnić twierdzenie o faktach i słowach?
- 5 Konkluzje

Konkluzje

- Przedstawiliśmy twierdzenie o faktach i słowach, które orzeka, że liczba **niezależnych** faktów opisywanych przez skończony tekst jest mniejsza niż liczba **różnych** słów w tymże tekście.
- Dla tekstów w języku naturalnym, liczba **różnych** słów zdaje się rosnać potęgowo z długością tekstu.
- Czy zatem liczba **niezależnych** faktów opisywanych przez teksty w języku naturalnym również rośnie potęgowo z długością tekstu?

Literatura

- 1 G. J. Chaitin, (1975). *A theory of program size formally identical to information theory*. Journal of the ACM, 22:329–340.
- 2 T. M. Cover and J. A. Thomas, (2006). *Elements of Information Theory, 2nd ed.* New York: John Wiley.
- 3 Ł. Dębowski, (2018). *Is Natural Language a Perigraphic Process? The Theorem about Facts and Words Revisited*. Entropy, 20(2):85.