

Praktyczne i teoretyczne problemy statystycznego modelowania języka naturalnego

Część III: Twierdzenie o faktach i słowach

Łukasz Dębowski
ldebowsk@ipipan.waw.pl



Instytut Podstaw Informatyki
Polskiej Akademii Nauk

XLIII Konferencja „Statystyka Matematyczna”
Będlewo, 04–08.12.2017

Statystyczne prawa językowe

Teksty w języku naturalnym spełniają **przybliżone** prawa ilościowe:

- 1 **Prawo Zipfa:** częstość słowa jest odwrotnie proporcjonalna do rangi słowa.
- 2 **Prawo Heapsa:** liczba różnych słów w tekście rośnie potęgowo z długością tekstu.
- 3 **Intensywność entropii Shannona:** jest rzędu 1 bita na literę.
- 4 **Hipoteza Hilberga:** entropia warunkowa litery maleje potęgowo z długością kontekstu.
- 5 **Prawo kodu PPM:** liczba różnych „słów” wykrywanych przez algorytm PPM w tekście rośnie potęgowo z długością tekstu.
- 6 **Prawo maksymalnego powtórzenia:** długość maksymalnego powtórzenia rośnie jak sześcian logarytmu długości tekstu.

Czy można coś wywnioskować o języku jako procesie stochastycznym na podstawie tych obserwacji/hipotez?

Pytania matematyka

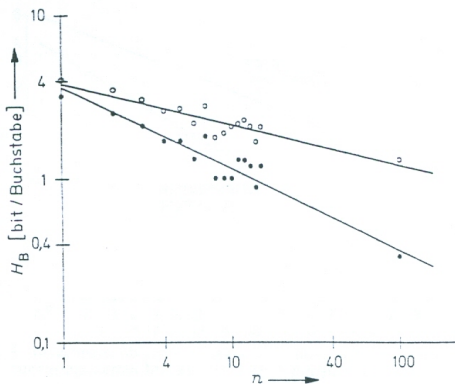
- 1 Czy istnieje idealny probabilistyczny model języka?
- 2 Czy model ten może być modelem Markowa?
- 3 Czy model ten może być ukrytym modelem Markowa?
- 4 Czy model ten jest ergodyczny?
- 5 Czy model ten jest stacjonarny?
- 6 Czy model ten jest asymptotycznie średnio stacjonarny?
- 7 Czy model ten jest kodem uniwersalnym?
- 8 Czy model ten jest efektywnie obliczalny?

- 1 Punkt wyjścia
- 2 Słowa PPM i informacja wzajemna
- 3 Fakty i informacja wzajemna
- 4 Konkluzje
- 5 Dodatek

Hipoteza Hilberga dla entropii Shannona (1990)

$$H(X) := \mathbb{E} [-P(X)]$$

$$H(X|Y) := \mathbb{E} [-P(X|Y)]$$



$$H(X_n|X_1^{n-1}) \approx Bn^{\beta-1} + h, \quad \beta \approx 1/2, \quad n \leq 100$$

Prawa potęgowe dla entropii i informacji wzajemnej

Informacja wzajemna dla procesu stacjonarnego:

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

Hipoteza Hilberga:

$$\begin{aligned} H(X_n | X_1^{n-1}) &\approx Bn^{\beta-1} + h \\ &\Downarrow \\ H(X_1^n) &= \sum_{i=1}^n H(X_i | X_1^{i-1}) \approx B'n^\beta + hn \\ &\Downarrow \\ I(X_1^n; X_{n+1}^{2n}) &\approx B''n^\beta \end{aligned}$$

Wykładnik Hilberga

Wykładnik Hilberga definiujemy jako

$$\mathbf{hilb} s(n) := \limsup_{n \rightarrow \infty} \frac{\log^+ s(2^n)}{\log 2^n}, \quad \log^+ x = \begin{cases} \log(x+1), & x \geq 0, \\ 0, & x < 0. \end{cases}$$

Na przykład $\mathbf{hilb} n^\beta = \beta$.

Twierdzenie

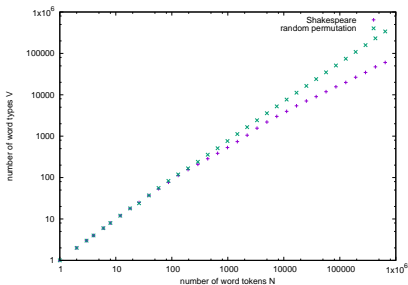
Niech $\mathfrak{J}(n) := 2\mathfrak{G}(n) - \mathfrak{G}(2n)$. Jeżeli $\lim_{n \rightarrow \infty} \mathfrak{G}(n)/n = g$, to

$$\mathbf{hilb} (\mathfrak{G}(n) - ng) \leq \mathbf{hilb} \mathfrak{J}(n),$$

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

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

Prawo Heapsa (liczba różnych słów ortograficznych)



Kuraszkiewicz i Łukaszewicz 1951; Herdan 1964; Heaps 1978:

$$V \propto N^\gamma, \quad \gamma < 1$$

V — liczba różnych słów w tekście (typów/types).

N — liczba wszystkich słów tekście (okazów/tokens).

- 1 Punkt wyjścia
- 2 Słowa PPM i informacja wzajemna
- 3 Fakty i informacja wzajemna
- 4 Konkluzje
- 5 Dodatek

Kod PPM (Prediction by Partial Matching)

Definiujemy

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

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

$$\text{PPM}(x_1^n) := \frac{6}{\pi^2} \sum_{k=-1}^{\infty} \frac{\text{PPM}_k(x_1^n)}{(k+2)^2}.$$

Wielkość $\text{PPM}(x_1^n)$ nazywa się **p-stwem PPM** napisu x_1^n .

Zauważmy, że $\text{PPM}_k(x_1^n) = D^{-n}$ dla $k > L(x_1^n)$.

Uniwersalność p-stwa PPM

Entropia bloku: $H(X_1^n) = \mathbb{E} [-\log P(X_1^n)]$

Intensywność entropii: $h = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} [-\log P(X_1^n)]$

Twierdzenie

P-stwo PPM jest p-stwem **uniwersalnym**, tzn. zachodzi

$$\mathbb{E} [-\log \text{PPM}(X_1^n)] \geq \mathbb{E} [-\log P(X_1^n)]$$
$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} [-\log \text{PPM}(X_1^n)] = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} [-\log P(X_1^n)]$$

dla dowolnego procesu stacjonarnego $(X_i)_{i=1}^{\infty}$ o skończ. alfabecie.

Rząd PPM i słownik PPM

- **Rząd PPM** $G_{\text{PPM}}(x_1^n)$ to najmniejsza liczba G taka, że
– $\log \text{PPM}_G(x_1^n) \leq -\log \text{PPM}_k(x_1^n)$ dla każdego $k \geq -1$.

- Zbiór wszystkich podstów długości m w napisie x_1^n to

$$V(m|x_1^n) := \{y_1^m : x_{t+1}^{t+m} = y_1^m \text{ dla pewnego } 0 \leq t \leq n - m\}.$$

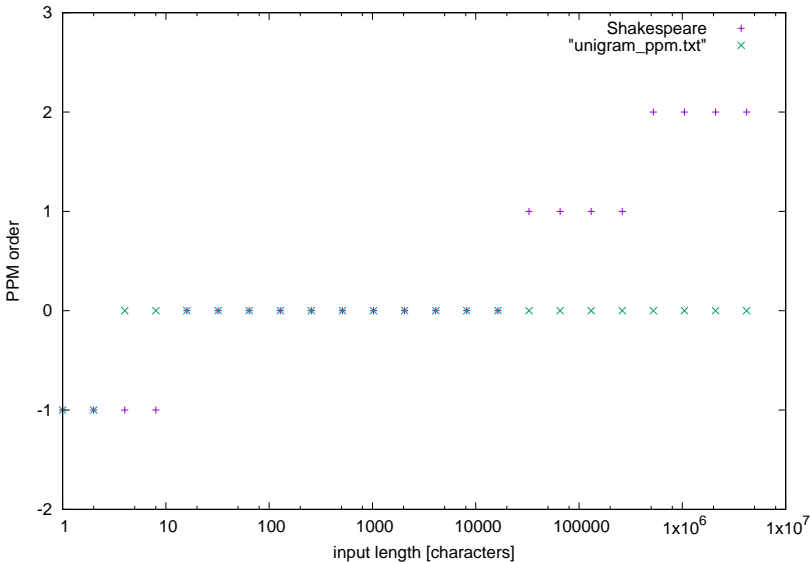
- Zbiór różnych **słów PPM** w napisie X_1^n to

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

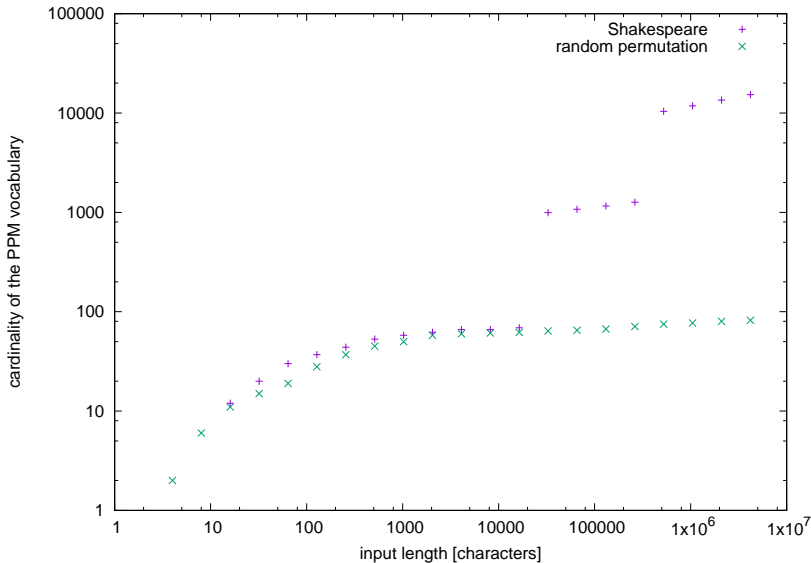
- Ogólnie zachodzi nierówność

$$\text{card } V_{\text{PPM}}(x_1^n) \leq \min \left\{ D^{G_{\text{PPM}}(x_1^n)}, n - G_{\text{PPM}}(x_1^n) + 1 \right\}.$$

Rząd PPM na wykresie



Moc słownika PPM na wykresie



Wykładnik Hilberga dla mocy słownika PPM

Dla procesów Markowa k -tego rzędu nad skończonym alfabetem

$$\lim_{n \rightarrow \infty} G_{\text{PPM}}(X_1^n) \leq k, \text{ czyli } \lim_{n \rightarrow \infty} \text{hilb card } V_{\text{PPM}}(X_1^n) = 0.$$

Dla języka naturalnego mamy prawdopodobnie

$$\lim_{n \rightarrow \infty} \text{hilb card } V_{\text{PPM}}(X_1^n) = \beta, \quad \beta \in (0, 1).$$

Czyli język naturalny **nie jest procesem Markowa**.

W poprzednim wykładzie wykazaliśmy mocniejszy fakt,
że język nie jest ukrytym procesem Markowa.

Twierdzenie o słowach i informacji wzajemnej

Entropia krzyżowa PPM:

$$H_{\text{PPM}}(\mathbf{X}_1^n) = \mathbb{E} \left[-\log \text{PPM}(\mathbf{X}_1^n) \right]$$

Informacja wzajemna PPM dla procesu stacjonarnego:

$$\begin{aligned} I_{\text{PPM}}(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) &:= H_{\text{PPM}}(\mathbf{X}_1^n) + H_{\text{PPM}}(\mathbf{X}_{n+1}^{2n}) - H_{\text{PPM}}(\mathbf{X}_1^{2n}) \\ &= 2H_{\text{PPM}}(\mathbf{X}_1^n) - H_{\text{PPM}}(\mathbf{X}_1^{2n}) \end{aligned}$$

Twierdzenie (informacja wzajemna i słowa PPM)

Dla dowolnego procesu stacjonarnego $(\mathbf{X}_i)_{i=1}^{\infty}$ o skończ. alfabecie

$$\begin{aligned} \text{hilb}_{n \rightarrow \infty} I(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) &\leq \text{hilb}_{n \rightarrow \infty} I_{\text{PPM}}(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) \\ &\leq \text{hilb}_{n \rightarrow \infty} \mathbb{E} \left[G_{\text{PPM}}(\mathbf{X}_1^n) + \text{card } V_{\text{PPM}}(\mathbf{X}_1^n) \right]. \end{aligned}$$

- 1 Punkt wyjścia
- 2 Słowa PPM i informacja wzajemna
- 3 **Fakty i informacja wzajemna**
- 4 Konkluzje
- 5 Dodatek

Twierdzenie ergodyczne i procesy ergodyczne

- 1 Rozpatrzmy proces dyskretny $(X_i)_{i=1}^{\infty} = (X_1, X_2, X_3, \dots)$.
- 2 Dla napisu $\mathbf{w} = (x_1, \dots, x_n)$ określmy zmienną losową

$$Y_i^{\mathbf{w}} := \begin{cases} 1 & \text{jeżeli } X_i = x_1, \dots, X_{i+n-1} = x_n, \\ 0 & \text{inaczej.} \end{cases}$$

- 3 Proces $(X_i)_{i=1}^{\infty}$ nazywa się **stacjonarnym**, gdy wartości oczekiwane $\mathbb{E} Y_i^{\mathbf{w}}$ nie zależą od i dla wszystkich \mathbf{w} .

Twierdzenie (ergodyczne)

Dla dowolnego procesu **stacjonarnego** $(X_i)_{i=1}^{\infty}$, istnieją granice

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n Y_i^{\mathbf{w}} = Y^{\mathbf{w}} \text{ prawie na pewno.}$$

- 4 Proces $(X_i)_{i=1}^{\infty}$ nazywa się **ergodycznym**, gdy granice $Y^{\mathbf{w}}$ są stałe dla każdego \mathbf{w} .

Przykłady ergodycznych procesów stacjonarnych

- 1 Proces $(X_i)_{i=1}^{\infty}$ nazywa się **IID**, jeżeli

$$P(X_1 = x_1, \dots, X_n = x_n) = \pi(x_1) \dots \pi(x_n).$$

— Procesy IID są ergodyczne.

- 2 Proces $(X_i)_{i=1}^{\infty}$ nazywa się **Markowa**, jeżeli

$$P(X_1 = x_1, \dots, X_n = x_n) = \pi(x_1) p(x_2|x_1) \dots p(x_n|x_{n-1}).$$

— Proces Markowa jest ergodyczny, jeśli $p(x_i|x_{i-1}) > c > 0$.

- 3 Proces $(Y_i)_{i=1}^{\infty}$ nazywa się **ukrytym Markowa**, jeżeli $Y_i = f(X_i)$ dla pewnego procesu Markowa $(X_i)_{i=1}^{\infty}$.

— Ukryty proces Markowa jest ergodyczny, jeśli jego proces Markowa jest ergodyczny.

Czy język naturalny jest nieergodyczny?

- 1 Proces jest **ergodyczny**, jeżeli częstości napisów w dostatecznie długiej próbie zbiegają do wartości stałych.
- 2 Wyobraźmy sobie teraz, że sięgamy po tekst będący losowo wybraną książką z biblioteki.
- 3 Spróbujmy policzyć częstości **słowa kluczowego**, np. słowa *bijekcja* czy słowa *skamielina*.
- 4 Spodziewamy się, że częstości **słów kluczowych** są zmiennymi losowymi o wartościach zależnych od **tematu** losowego tekstu.
- 5 Ponieważ słowa kluczowe są pewnymi napisami, to proces stochastyczny modelujący język naturalny powinien być **nieergodyczny**.

Zliczając **słowa kluczowe**, możemy rozpoznać **temat** tekstu.

Procesy nieergodyczne — inna perspektywa

Intuicja: Proces jest nieergodyczny \iff istnieją \geq dwa tematy.

Twierdzenie

Proces $(X_i)_{i=1}^{\infty}$ jest **nieergodyczny** wtedy i tylko wtedy, gdy istnieje funkcja $f(x_1, \dots, x_n)$ ciągów symboli i binarna zmienna losowa Z takie, że $0 < P(Z = 0) < 1$ oraz

$$\lim_{n \rightarrow \infty} P(f(X_{t+1}, \dots, X_{t+n}) = Z) = 1 \quad (1)$$

dla dowolnej pozycji t .

Definicja

Zmienną losową Z spełniającą (1) nazywam **losowym faktem**.

Zatem proces jest **nieergodyczny**, jeśli istnieje \geq jeden **losowy fakt**.
Losowy fakt informuje, na który z **dwóch tematów** jest tekst.

Proces Santa Fe — przykład procesu nieergodycznego

- ❶ Niech $(Z_k)_{k=1}^{\infty}$ będzie procesem IID o $Z_k \in \{0, 1\}$ oraz

$$P(Z_k = 0) = P(Z_k = 1) = 1/2.$$

- ❷ Niech $(K_i)_{i=1}^{\infty}$ będzie procesem IID o $K_i \in \{1, 2, 3, \dots\}$ oraz

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

- ❸ **Proces Santa Fe** to proces $(X_i)_{i=1}^{\infty}$, gdzie

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

- ❹ Proces Santa Fe jest nieergodyczny, gdyż **wszystkie** Z_k są probabilistycznie niezależnymi **faktami losowymi**.

Mocna nieergodyczność

Intuicja: Proces Santa Fe process jest **mocno nieergodyczny**, gdyż istnieje **nieskończenie** wiele probab. niezależnych **faktów losowych**.

Definicja

Proces $(X_i)_{i=1}^{\infty}$ nazywa się **mocno nieergodycznym**, gdy istnieją funkcje $f_k(x_1, \dots, x_n)$ ciągów symboli i binarny proces IID $(Z_k)_{k=1}^{\infty}$ takie, że $P(Z_k = 0) = P(Z_k = 1) = 1/2$ oraz

$$\lim_{n \rightarrow \infty} P(f_k(X_{t+1}, \dots, X_{t+n}) = Z_k) = 1$$

dla dowolnej pozycji t i dowolnego $k = 1, 2, 3, \dots$

(liczba tematów) $\approx 2^{(\text{liczba niezależnych faktów losowych})}$

mocna nieergodyczność \iff **kontinuum tematów**

Potęgowy wzrost liczby przewidywalnych faktów

Zbiór niezależnych **losowych faktów** przewidywalnych z \mathbf{X}_1^n to

$$U(\mathbf{X}_1^n) := \{l \in \{1, 2, \dots\} : f_k(\mathbf{X}_1^n) = Z_k \text{ for all } k \leq l\}.$$

Dla procesów Santa Fe mamy dokładne równości

$$\lim_{n \rightarrow \infty} \mathbb{E} \text{ card } U(\mathbf{X}_1^n) = \lim_{n \rightarrow \infty} I(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) = \beta,$$

gdzie $\beta = 1/\alpha \in (0, 1)$.

W ogólnym przypadku wykładniki Hilberga mogą być różne.

Twierdzenie o faktach i słowach

Twierdzenie (fakty i informacja wzajemna)

Dla **dowolnego mocno nieergodycznego** procesu stacjonarnego $(X_i)_{i=1}^{\infty}$ o skończ. alfabecie

$$\mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} \text{ card } U(X_1^n) \leq \mathop{\text{hilb}}_{n \rightarrow \infty} I(X_1^n; X_{n+1}^{2n}).$$

Twierdzenie (informacja wzajemna i słowa PPM)

Dla **dowolnego** procesu stacjonarnego $(X_i)_{i=1}^{\infty}$ o skończ. alfabecie

$$\begin{aligned} \mathop{\text{hilb}}_{n \rightarrow \infty} I(X_1^n; X_{n+1}^{2n}) &\leq \mathop{\text{hilb}}_{n \rightarrow \infty} I_{\text{PPM}}(X_1^n; X_{n+1}^{2n}) \\ &\leq \mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} \left[G_{\text{PPM}}(X_1^n) + \text{card } V_{\text{PPM}}(X_1^n) \right]. \end{aligned}$$

- 1 Punkt wyjścia
- 2 Słowa PPM i informacja wzajemna
- 3 Fakty i informacja wzajemna
- 4 Konkluzje**
- 5 Dodatek

Konkluzje

Czy język naturalny jest **mocno nieergodyczny**?

Nie możemy tego wykluczyć, gdyż:

- 1 Liczba różnych słów PPM w tekstach zdaje się rosnąć potęgowo z długością tekstu.
- 2 W świetle twierdzenia o faktach i słowach liczba niezależnych faktów przewidywalnych na podstawie tekstu może także rosnąć potęgowo z długością tekstu.

Problem otwarty

Udowodnić, że zachodzi prawie na pewno

$$\liminf_{n \rightarrow \infty} \text{card } U(X_1^n) \leq \liminf_{n \rightarrow \infty} [\text{G}_{\text{PPM}}(X_1^n) + \text{card } V_{\text{PPM}}(X_1^n)] .$$

- 1 Punkt wyjścia
- 2 Słowa PPM i informacja wzajemna
- 3 Fakty i informacja wzajemna
- 4 Konkluzje
- 5 **Dodatek**

Podążanie w duchu złożoności Kolmogorowa

- 1 Modyfikacje definicji mocnej nieergodyczności:

$(Z_k)_{k=1}^{\infty} \longrightarrow (z_k)_{k=1}^{\infty}$ — ciąg **algorytmicznie losowy**

$(f_k)_{k=1}^{\infty} \longleftarrow$ funkcja **efektywnie obliczalna**

- 2 Modyfikacja zbioru przewidywalnych faktów:

$$U_a(X_1^n) := \{I \in \{1, 2, \dots\} : f_k(X_1^n) = z_k \text{ for all } k \leq I\}.$$

- 3 Dla zmodyfikowanego procesu Santa Fe postaci $X_i = (K_i, z_{K_i})$, gdzie K_i jak poprzednio, mamy

$$\lim_{n \rightarrow \infty} \mathbb{E} \text{card } U_a(X_1^n) = \beta,$$

gdzie $\beta = 1/\alpha \in (0, 1)$.

- 4 Informację wzajemną $I(X_1^n; X_{n+1}^{2n})$ zastępujemy przez oczekiwaną **algorytmiczną informację wzajemną** $I_a(X_1^n; X_{n+1}^{2n})$

Algorytmiczne twierdzenie o faktach i słowach

Twierdzenie (fakty i informacja wzajemna)

Dla **dowolnego** procesu stacjonarnego $(X_i)_{i=1}^{\infty}$ o skończ. alfabecie

$$\mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} \text{ card } U_a(X_1^n) \leq \mathop{\text{hilb}}_{n \rightarrow \infty} I_a(X_1^n; X_{n+1}^{2n}).$$

Twierdzenie (informacja wzajemna i słowa PPM)

Dla **dowolnego** procesu stacjonarnego $(X_i)_{i=1}^{\infty}$ o skończ. alfabecie

$$\begin{aligned} \mathop{\text{hilb}}_{n \rightarrow \infty} I_a(X_1^n; X_{n+1}^{2n}) &\leq \mathop{\text{hilb}}_{n \rightarrow \infty} I_{\text{PPM}}(X_1^n; X_{n+1}^{2n}) \\ &\leq \mathop{\text{hilb}}_{n \rightarrow \infty} \mathbb{E} \left[G_{\text{PPM}}(X_1^n) + \text{card } V_{\text{PPM}}(X_1^n) \right]. \end{aligned}$$