

# Maksymalne powtórzenia w tekście i hipotezy probabilistyczne

Łukasz Dębowski  
ldebowsk@ipipan.waw.pl



Instytut Podstaw Informatyki PAN  
Warszawa

2 (4) Warsztat badawczy i Perypatetyczna Konferencja  
„Modelowanie Systemów Poznawczych”  
Zakopane, 3-5.10.2014

- 1 Hipotezy probabilistyczne
- 2 Maksymalne powtórzenia
- 3 Dane empiryczne
- 4 Konkluzje

## Teksty — między losowością a determinizmem

*If a Martian scientist sitting before his radio in Mars accidentally received from Earth the broadcast of an extensive speech which he recorded perfectly through the perfection of Martian apparatus and studied at his leisure, what criteria would he have to determine whether the reception represented the effect of animate process on Earth, or merely the latest thunderstorm on Earth? It seems that the only criteria would be the arrangement of occurrences of the elements, and 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)

# Dwie hipotezy probabilistyczne

## Hipoteza skończonej energii:

Proces generowania tekstu jest procesem o skończonej energii.

Hipoteza skończonej energii jest spełniona, jeżeli teksty są zaszumione w pewien dość ogólny sposób.

## Mocna hipoteza Hilberga:

Proces generowania tekstu jest mocnym procesem Hilberga.

Zgodnie z hipotezą Hilberga tylko niewielka część możliwych tekstów podlega replikacji. ( $\implies$  ewolucja memetyczna?)

# Hipoteza skończonej energii

- $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  — proces stochastyczny — ciąg zmiennych losowych  $\mathbf{X}_i : \Omega \rightarrow \mathbb{X}$  na przestrzeni probabilistycznej  $(\Omega, \mathcal{J}, \mathbf{P})$ .
- $\mathbf{X}_k^l = (\mathbf{X}_k, \mathbf{X}_{k+1}, \dots, \mathbf{X}_l)$  — bloki zmiennych losowych.
- $|u|$  — długość napisu  $u$ .

## Definicja (Shields 1997)

Proces  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  nazywamy **procesem o skończonej energii**, jeżeli prawdopodobieństwo warunkowe bloków spełnia

$$\mathbf{P} \left( \mathbf{X}_{i+|w|+1}^{i+|wu|} = u \mid \mathbf{X}_{i+1}^{i+|w|} = w \right) \leq \mathbf{K}c^{|u|}$$

dla pewnych stałych  $0 < c < 1$  oraz  $\mathbf{K} > 0$ .

# Mocna hipoteza Hilberga

## Definicja

Dla zmiennej losowej  $\mathbf{X}$ , definiuje się **entropię topologiczną** jako logarytm liczby wartości o niezerowym prawdopodobieństwie,

$$H_{\text{top}}(\mathbf{X}) = \log \text{card} \{x : P(\mathbf{X} = x) > 0\} .$$

## Definicja

Proces  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  nazywamy **mocnym procesem Hilberga**, jeżeli

$$H_{\text{top}}(\mathbf{X}_{i+1}^{i+n}) \leq \mathbf{A}n^\beta$$

dla pewnych stałych  $0 < \beta < 1$  oraz  $\mathbf{A} > 0$ .

- 1 Hipotezy probabilistyczne
- 2 Maksymalne powtórzenia
- 3 Dane empiryczne
- 4 Konkluzje

# Maksymalne powtórzenie

## Definicja

Dla tekstu  $w$  definiuje się **maksymalne powtórzenie** jako

$$L(w) := \max \{ |s| : w = x_1 s y_1 = x_2 s y_2 \text{ i } x_1 \neq x_2 \},$$

gdzie  $s$ ,  $x_i$  oraz  $y_i$  są fragmentami tekstu  $w$ .



# Dwa twierdzenia

## Twierdzenie 1 (Shields 1997)

Dla procesu o skończonej energii  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  istnieje  $\mathbf{C} > \mathbf{0}$  takie, że

$$\mathbf{L}(\mathbf{X}_1^m) \leq \mathbf{C} \log m$$

zachodzi z prawdopodobieństwem  $\mathbf{1}$  dla dostatecznie dużych  $\mathbf{m}$ .

## Twierdzenie 2

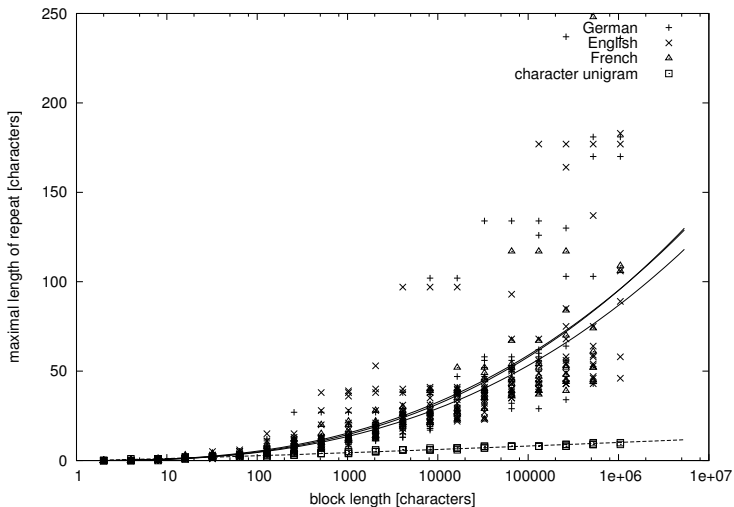
Dla mocnego procesu Hilberga  $(\mathbf{X}_i)_{i \in \mathbb{Z}}$  istnieje  $\mathbf{C} > \mathbf{0}$  takie, że

$$\mathbf{L}(\mathbf{X}_1^m) \geq \mathbf{C}(\log m)^{1/\beta}$$

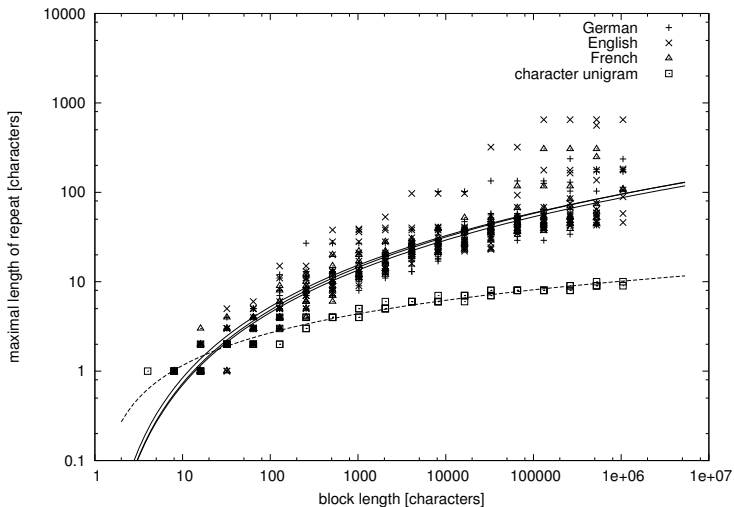
zachodzi z prawdopodobieństwem  $\mathbf{1}$ .

- 1 Hipotezy probabilistyczne
- 2 Maksymalne powtórzenia
- 3 Dane empiryczne
- 4 Konkluzje

# 35 tekstów w 3 językach + tekst unigramowy



# To samo w skali podwójnie logarytmicznej



# Parametry modelu

$$L(X_1^m) \approx A(\log m)^\alpha$$

klasa tekstów	<b>A</b>	<b><math>\alpha</math></b>
j. niemiecki	<b>0.076 ± 0.011</b>	<b>2.71 ± 0.07</b>
j. angielski	<b>0.093 ± 0.012</b>	<b>2.64 ± 0.06</b>
j. francuski	<b>0.074 ± 0.009</b>	<b>2.69 ± 0.06</b>
unigramowe	<b>0.42 ± 0.03</b>	<b>1.21 ± 0.03</b>

# Przykłady maksymalnych powtórzeń

tekst	maksymalny powtórzony napis
Thomas Mann, Buddenbrooks	“Mit raschen Schritten, die Arme ausgebreitet und den Kopf zur Seite geneigt, in der Haltung eines Mannes, welcher sagen will: Hier bin ich! Töte mich, wenn du willst!”
Jules Verne, Vingt mille lieues sous les mers	“à la partie supérieure de la coque du «Nautilus», et”
Jonathan Swift, Gulliver's Travels	“of meat and drink sufficient for the support of 1724”
tekst unigramowy	“e t tloeu ”

- 1 Hipotezy probabilistyczne
- 2 Maksymalne powtórzenia
- 3 Dane empiryczne
- 4 Konkluzje

# Wnioski

- 1 W tekstach w języku naturalnym maksymalne powtórzenie rośnie znacznie szybciej niż w tekstach unigramowych.
- 2 Obserwacja ta falsyfikuje hipotezę skończonej energii.
- 3 A zatem teksty nie są zaszumione.
- 4 Obserwowany wzrost maksymalnego powtórzenia może być natomiast konsekwencją mocnej hipotezy Hilberga.
- 5 Zgodnie z hipotezą Hilberga tylko niewielka część możliwych tekstów podlega replikacji. ( $\implies$  ewolucja memetyczna?)

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



## Post scriptum: Losowość a dowolność

*El universo (que otros llaman la Biblioteca) se compone de un número indefinido, y tal vez infinito, de galerías hexagonales, con vastos pozos de ventilación en el medio, cercados por barandas bajísimas. Desde cualquier hexágono se ven los pisos inferiores y superiores: interminablemente. La distribución de las galerías es invariable. Veinte anaqueles, a cinco largos anaqueles por lado, cubren todos los lados menos dos; su altura, que es la de los pisos, excede apenas la de un bibliotecario normal. Una de las caras libres da a un angosto zaguán, que desemboca en otra galería, idéntica a la primera y a todas. ...*

— *La Biblioteca de Babel*, Jorge Luis Borges