

Hilberg Exponents: New Measures of Long Memory in the Process

Łukasz Dębowski
ldebowsk@ipipan.waw.pl



Institute of Computer Science
Polish Academy of Sciences

12th November 2015, Gifu

- 1 Introduction
- 2 Relating $\gamma_{\mathbf{P}}^{\pm}$ and $\delta_{\mathbf{P}}^{\pm}$
- 3 Relating $\delta_{\mathbf{P}}^{+}$ and $\delta_{\mathbf{R}}^{+}$
- 4 Evaluating $\gamma_{\mathbf{Q}}^{\pm}$ and $\delta_{\mathbf{Q}}^{\pm}$
- 5 Conclusion

Motivation: Hilberg's hypothesis

- According to a hypothesis by Hilberg (1990), the mutual information between two adjacent blocks of text in natural language grows like a power of the block length.
- This property differentiates natural language from \mathbf{k} -parameter sources, for which the mutual information is proportional to the logarithm of the block length.
- In 2011, we constructed processes, called Santa Fe processes, which feature the power-law growth of mutual information.
- In 2011, we also showed that Hilberg's hypothesis implies Herdan's law, some version of Zipf's law.
- In 2014, we showed experimentally that for a PPM-like code the estimates of mutual information grow as a power law for natural language and logarithmically for a \mathbf{k} -parameter source.

Preliminaries

- \mathbb{X} — a countable alphabet,
 $(\Omega, \mathcal{J}, \mathbf{Q})$ — probability space with $\Omega = \mathbb{X}^{\mathbb{Z}}$,
 $\mathbf{X}_k : \Omega \ni (\mathbf{x}_i)_{i \in \mathbb{Z}} \mapsto \mathbf{x}_k \in \mathbb{X}$ — random variables,
 \mathbf{Q} — a stationary measure (not necessarily ergodic),
 $\mathbf{X}_n^m = (\mathbf{X}_i)_{n \leq i \leq m}$ — blocks of symbols,
 $\mathbf{E}_{\mathbf{Q}} \mathbf{X}$ — expectation,
 $\mathbf{Var}_{\mathbf{Q}} \mathbf{X}$ — variance.
- \mathbf{P} — a code (incomplete measure), i.e., it satisfies $\mathbf{P}(\mathbf{x}_1^n) \geq 0$
 and the Kraft inequality $\sum_{\mathbf{x}_1^n} \mathbf{P}(\mathbf{x}_1^n) \leq 1$.
 For \mathbf{P} (or for \mathbf{Q}), we define the pointwise mutual information

$$I^{\mathbf{P}}(n) = -\log \mathbf{P}(\mathbf{X}_{-n+1}^0) - \log \mathbf{P}(\mathbf{X}_1^n) + \log \mathbf{P}(\mathbf{X}_{-n+1}^n).$$

In the formula **log** stands for the binary logarithm.

Hilberg exponents

Define the positive logarithm

$$\log^+ x = \begin{cases} \log(x + 1), & x \geq 0, \\ 0, & x < 0. \end{cases}$$

For a code \mathbf{P} we introduce

$$\begin{aligned} \gamma_P^+ &= \limsup_{n \rightarrow \infty} \frac{\log^+ I^P(n)}{\log n}, & \gamma_P^- &= \liminf_{n \rightarrow \infty} \frac{\log^+ I^P(n)}{\log n}, \\ \delta_P^+ &= \limsup_{n \rightarrow \infty} \frac{\log^+ E_Q I^P(n)}{\log n}, & \delta_P^- &= \liminf_{n \rightarrow \infty} \frac{\log^+ E_Q I^P(n)}{\log n}. \end{aligned}$$

We call these: γ_P^+ —the upper random Hilberg exponent, γ_P^- —the lower random Hilberg exponent, δ_P^+ —the upper expected Hilberg exponent, and δ_P^- —the lower expected Hilberg exponent.

Basic observations

By definition,

$$\gamma_P^+ \geq \gamma_P^- \geq 0,$$

$$\delta_P^+ \geq \delta_P^- \geq 0.$$

- For $\mathbf{P} = \mathbf{Q}$, Hilberg exponents quantify some sort of long-range non-Markovian dependence in the process.
- For an IID process or a hidden Markov process with a finite number of hidden states, $\mathbf{E}_Q \mathbf{I}^Q(\mathbf{n}) \leq \mathbf{D}$ so $\delta_Q^\pm = 0$.
- For a \mathbf{k} -parameter source, $\mathbf{E}_Q \mathbf{I}^Q(\mathbf{n}) \propto \mathbf{k} \log \mathbf{n}$, so $\delta_Q^\pm = 0$.
- If $\mathbf{E}_Q \mathbf{I}^Q(\mathbf{n}) \propto \mathbf{n}^\beta$ where $\beta \in [0, 1]$ then $\delta_Q^\pm = \beta$.
- There exist some non-Markovian but mixing sources, being a generalization of Santa Fe processes, for which $\delta_Q^\pm \in (0, 1)$.

Further simple observations

Inequalities

$$\gamma_P^- \leq \gamma_P^+ \leq 1,$$

$$\delta_P^- \leq \delta_P^+ \leq 1$$

hold in the following cases:

- 1 For $\mathbf{P} = \mathbf{Q}$ — by the Shannon-McMillan-Breiman theorem and by stationarity.
- 2 For \mathbf{P} being universal almost surely and in expectation, that is if for $\mathbf{k}(\mathbf{n})$ and $\mathbf{l}(\mathbf{n})$ being nondecreasing functions of \mathbf{n} , where $\mathbf{k}(\mathbf{n}) + \mathbf{l}(\mathbf{n}) \rightarrow \infty$ we have

$$\lim_{n \rightarrow \infty} \frac{1}{\mathbf{k}(\mathbf{n}) + \mathbf{l}(\mathbf{n}) + 1} \left[-\log \mathbf{P}(\mathbf{X}_{-\mathbf{k}(\mathbf{n})}^{\mathbf{l}(\mathbf{n})}) \right] = h_{\mathbf{Q}},$$

where $h_{\mathbf{Q}}$ is the entropy rate of measure \mathbf{Q} , and if

$$\lim_{n \rightarrow \infty} \frac{1}{\mathbf{k}(\mathbf{n}) + \mathbf{l}(\mathbf{n}) + 1} \mathbf{E}_{\mathbf{Q}} \left[-\log \mathbf{P}(\mathbf{X}_{-\mathbf{k}(\mathbf{n})}^{\mathbf{l}(\mathbf{n})}) \right] = \mathbf{E}_{\mathbf{Q}} h_{\mathbf{Q}}.$$

A research program

There are many Hilberg exponents, for different measures and for different codes. Seeking for some order, we may look for results of three kinds:

- ① For a fixed code \mathbf{P} and a measure \mathbf{Q} , we relate the random exponents $\gamma_{\mathbf{P}}^{\pm}$ and the expected exponents $\delta_{\mathbf{P}}^{\pm}$.
- ② For two codes \mathbf{P} and \mathbf{R} , we relate the exponents of a fixed kind, say $\delta_{\mathbf{P}}^{\pm}$ and $\delta_{\mathbf{R}}^{\pm}$ for some measure \mathbf{Q} .
- ③ For a fixed code \mathbf{P} and a measure \mathbf{Q} , we directly evaluate exponents $\gamma_{\mathbf{P}}^{\pm}$ and $\delta_{\mathbf{P}}^{\pm}$.

In the following we will present some results of these three sorts.

- 1 Introduction
- 2 Relating $\gamma_{\mathbf{P}}^{\pm}$ and $\delta_{\mathbf{P}}^{\pm}$
- 3 Relating $\delta_{\mathbf{P}}^{+}$ and $\delta_{\mathbf{R}}^{+}$
- 4 Evaluating $\gamma_{\mathbf{Q}}^{\pm}$ and $\delta_{\mathbf{Q}}^{\pm}$
- 5 Conclusion

“Second-order” Shannon-McMillan-Breiman theorem

- The original idea of the SMB theorem was to relate the asymptotic growth of pointwise and expected entropies for an ergodic process \mathbf{Q} with $\mathbf{P} = \mathbf{Q}$.
- In contrast, relating the random Hilberg exponents γ_Q^\pm and the expected Hilberg exponents δ_Q^\pm means relating the speed of growth of the pointwise and expected mutual informations, which is a subtler effect than the SMB theorem.
- Thus relating γ_Q^\pm and δ_Q^\pm could be called a “second-order” analogue of the SMB theorem.

Our main result

For a code \mathbf{P} with exponent $\delta_P^- > 0$, let us introduce

$$\epsilon_P = \limsup_{n \rightarrow \infty} \frac{\log^+ [\text{Var}_Q I^P(n) / E_Q I^P(n)]}{\log n}.$$

Theorem

For an ergodic measure \mathbf{Q} over a finite alphabet, random Hilberg exponents γ_Q^\pm are almost surely constant. Moreover, we have \mathbf{Q} -almost surely

$$\delta_Q^+ \geq \gamma_Q^+ \geq \delta_Q^+ - \epsilon_Q,$$

$$\delta_Q^- \geq \gamma_Q^- \geq \delta_Q^- - \epsilon_Q,$$

where the left inequalities hold without restrictions, whereas the right inequalities hold for $\delta_Q^- > 0$.

The fundamental idea of the proof

Our theorem can be demonstrated without invoking the ergodic theorem. Instead, we use an auxiliary “Kolmogorov code”

$$S(\mathbf{x}_1^n) = 2^{-K(\mathbf{x}_1^n|\mathbf{F})},$$

where $K(\mathbf{x}_1^n|\mathbf{F})$ is the prefix-free Kolmogorov complexity of a string \mathbf{x}_1^n given an object \mathbf{F} on an additional infinite tape. The object \mathbf{F} can be another string or, here, a definition of measure \mathbf{Q} .

The first auxiliary result

Theorem

Consider Kolmogorov code \mathbf{S} and an ergodic \mathbf{Q} over a finite alphabet \mathbb{X} . Exponents $\gamma_{\mathbf{S}}^-$ and $\gamma_{\mathbf{S}}^+$ are \mathbf{Q} -almost surely constant.

The idea of the proof:

- 1 $|\mathbf{K}(\mathbf{x}_1^n | \mathbf{F}) - \mathbf{K}(\mathbf{x}_{t+1}^{t+n} | \mathbf{F})| \leq \mathbf{C}t.$
- 2 Hence $\gamma_{\mathbf{S}}^\pm$ are shift invariant.
- 3 Hence $\gamma_{\mathbf{S}}^\pm$ are constant on ergodic sources.

Second auxiliary result (via Borel-Cantelli lemma)

Consider Kolmogorov code \mathbf{S} . For $I^{\mathbf{S}}(n) + \mathbf{B} \geq 1$, define

$$\zeta_{\mathbf{S}}^+ = \limsup_{n \rightarrow \infty} \frac{\log^+ [\mathbf{E}_{\mathbf{Q}}(I^{\mathbf{S}}(n) + \mathbf{B})^{-1}]^{-1}}{\log n},$$

$$\zeta_{\mathbf{S}}^- = \liminf_{n \rightarrow \infty} \frac{\log^+ [\mathbf{E}_{\mathbf{Q}}(I^{\mathbf{S}}(n) + \mathbf{B})^{-1}]^{-1}}{\log n}.$$

These will be called inverse expected Hilberg exponents.

Theorem

Consider Kolmogorov code \mathbf{S} and a stationary measure \mathbf{Q} . Then:

- 1 $\delta_{\mathbf{S}}^+ \geq \gamma_{\mathbf{S}}^+$ \mathbf{Q} -almost surely and $\text{ess sup}_{\mathbf{Q}} \gamma_{\mathbf{S}}^+ \geq \zeta_{\mathbf{S}}^+$.
- 2 $\delta_{\mathbf{S}}^- \geq \text{ess inf}_{\mathbf{Q}} \gamma_{\mathbf{S}}^-$ and $\gamma_{\mathbf{S}}^- \geq \zeta_{\mathbf{S}}^-$ \mathbf{Q} -almost surely.

A corollary

Corollary

For an ergodic measure \mathbf{Q} over a finite alphabet, equalities $\gamma_S^+ = \mathbf{ess\,sup}_Q \gamma_S^+$ and $\gamma_S^- = \mathbf{ess\,inf}_Q \gamma_S^-$ hold \mathbf{Q} -almost surely. Hence, \mathbf{Q} -almost surely we have

$$\delta_S^+ \geq \gamma_S^+ \geq \zeta_S^+,$$

$$\delta_S^- \geq \gamma_S^- \geq \zeta_S^-.$$

Third auxiliary result

Theorem

Consider Kolmogorov code \mathbf{S} with $\mathbf{F} = \mathbf{Q}$, where \mathbf{Q} is a stationary measure. Then:

- ① $\delta_S^- = \delta_Q^-$ and $\delta_S^+ = \delta_Q^+$.
- ② $\gamma_S^- = \gamma_Q^-$ and $\gamma_S^+ = \gamma_Q^+$ \mathbf{Q} -almost surely.

(By Shannon-Fano coding and Barron's inequality.)

Fourth, the last auxiliary result

For a code \mathbf{P} with exponent $\delta_P^- > 0$, let us introduce

$$\epsilon_P = \limsup_{n \rightarrow \infty} \frac{\log^+ [\text{Var}_Q I^P(n) / E_Q I^P(n)]}{\log n}.$$

Theorem

Consider Kolmogorov code \mathbf{S} and a stationary measure \mathbf{Q} . If $\delta_S^- > 0$ then $\zeta_S^+ \geq \delta_S^+ - \epsilon_S$ and $\zeta_S^- \geq \delta_S^- - \epsilon_S$.

(By Markov inequality.)

Theorem

Consider Kolmogorov code \mathbf{S} with $\mathbf{F} = \mathbf{Q}$, where \mathbf{Q} is a stationary measure. If $\delta_Q^- > 0$ then $\epsilon_S = \epsilon_Q$.

(By Shannon-Fano coding and Barron's inequality.)

Resuming, our main result

Theorem

For an ergodic measure \mathbf{Q} over a finite alphabet, random Hilberg exponents $\gamma_{\mathbf{Q}}^\pm$ are almost surely constant. Moreover, we have \mathbf{Q} -almost surely

$$\delta_{\mathbf{Q}}^+ \geq \gamma_{\mathbf{Q}}^+ \geq \delta_{\mathbf{Q}}^+ - \epsilon_{\mathbf{Q}},$$

$$\delta_{\mathbf{Q}}^- \geq \gamma_{\mathbf{Q}}^- \geq \delta_{\mathbf{Q}}^- - \epsilon_{\mathbf{Q}},$$

where the left inequalities hold without restrictions, whereas the right inequalities hold for $\delta_{\mathbf{Q}}^- > \mathbf{0}$.

- 1 Introduction
- 2 Relating γ_P^\pm and δ_P^\pm
- 3 Relating δ_P^+ and δ_R^+
- 4 Evaluating γ_Q^\pm and δ_Q^\pm
- 5 Conclusion

An alternative expression for δ_P^+

Denote $\mathbf{H}^P(\mathbf{n}) := -\log \mathbf{P}(\mathbf{X}_1^n)$. Suppose that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{E}_Q \mathbf{H}^P(\mathbf{n}) = \mathbf{E}_Q h_Q$$

and suppose that $\mathbf{E}_Q I^P(\mathbf{n}) \geq -\mathbf{D}$ for a certain $\mathbf{D} > 0$. Then

$$\delta_P^+ = \limsup_{n \rightarrow \infty} \frac{\log^+ \mathbf{E}_Q I^P(\mathbf{n})}{\log n} = \limsup_{n \rightarrow \infty} \frac{\log^+ \mathbf{E}_Q [\mathbf{H}^P(\mathbf{n}) - h_Q n]}{\log n}$$

follows from the telescope sum

$$\mathbf{E}_Q \mathbf{H}^P(\mathbf{n}) - h_Q n = \sum_{k=0}^{\infty} \frac{\mathbf{E}_Q I^P(2^k \mathbf{n})}{2^{k+1}}.$$

A hierarchy of approximations of block entropy

- 1 $\mathbf{R}(\mathbf{X}_1^n) = 2^{-|\mathbf{C}(\mathbf{X}_1^n)|}$ — a computable universal code.
- 2 $\mathbf{P}(\mathbf{X}_1^n) = 2^{-\mathbf{K}(\mathbf{X}_1^n)}$ — the unconditional Kolmogorov code.
- 3 $\mathbf{Q}(\mathbf{X}_1^n)$ — the underlying measure.
- 4 $\mathbf{E}(\mathbf{X}_1^n) = \mathbf{Q}(\mathbf{X}_1^n | \mathcal{I})$ — the random ergodic measure.
- 5 $\mathbf{H}^T(\mathbf{n}) = \mathbf{H}^T(\mathbf{n}; \mathbf{X}_1^{n(|\mathbb{X}|+\epsilon)^n})$ — the plugin entropy estimator.

We have

$$\mathbf{E}_Q \mathbf{H}^R(\mathbf{n}) \geq \mathbf{E}_Q \mathbf{H}^P(\mathbf{n}) \geq \mathbf{E}_Q \mathbf{H}^Q(\mathbf{n}) \geq \mathbf{E}_Q \mathbf{H}^E(\mathbf{n}) \geq \mathbf{E}_Q \mathbf{H}^T(\mathbf{n}),$$

whereas the common rate of these is $\mathbf{E}_Q \mathbf{h}_Q$. Hence

$$\delta_R^+ \geq \delta_P^+ \geq \delta_Q^+ \geq \delta_E^+ \geq \delta_T^+.$$

The difference $\delta_P^+ - \delta_E^+$ can be arbitrarily close to $\mathbf{1}$.

- 1 Introduction
- 2 Relating $\gamma_{\mathbf{P}}^{\pm}$ and $\delta_{\mathbf{P}}^{\pm}$
- 3 Relating $\delta_{\mathbf{P}}^{+}$ and $\delta_{\mathbf{R}}^{+}$
- 4 Evaluating $\gamma_{\mathbf{Q}}^{\pm}$ and $\delta_{\mathbf{Q}}^{\pm}$
- 5 Conclusion

Memoryless sources and hidden Markov processes

- For IID processes, $\delta_Q^\pm = \mathbf{0}$ and hence $\gamma_Q^\pm = \mathbf{0}$ since there is no dependence in the process.
- For Markov processes over a finite alphabet and hidden Markov processes with a finite number of hidden states, we also have $\delta_Q^\pm = \mathbf{0}$ and hence $\gamma_Q^\pm = \mathbf{0}$, since the expected mutual information is bounded for measures of those processes by the data-processing inequality.

Mixture Bernoulli process

Some simple example of a process with unbounded mutual information is the mixture of Bernoulli processes over the alphabet $\mathbb{X} = \{0, 1\}$, which we will call the mixture Bernoulli process:

$$\mathbf{Q}(\mathbf{x}_1^n) = \int_0^1 \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i} d\theta = \frac{1}{n + 1} \binom{n}{\sum_{i=1}^n x_i}^{-1}.$$

Mixture Bernoulli process (continued)

$$\mathbf{T}_n = \sum_{i=-n+1}^0 \mathbf{X}_i, \quad \mathbf{S}_n = \sum_{i=1}^n \mathbf{X}_i.$$

Theorem

For the mixture Bernoulli process, $\delta_Q^\pm = \gamma_Q^\pm = 0$.

Proof: \mathbf{X}_{-n+1}^0 and \mathbf{X}_1^n are independent given \mathbf{T}_n and \mathbf{S}_n . Hence

$$I^Q(n) = -\log \frac{Q(\mathbf{T}_n)Q(\mathbf{S}_n)}{Q(\mathbf{T}_n, \mathbf{S}_n)}$$

so $\mathbf{E}_Q I^Q(n) = I_Q(\mathbf{T}_n; \mathbf{S}_n)$. Variable \mathbf{S}_n assumes under \mathbf{Q} each value in $\{0, 1, \dots, n\}$ with equal probability $(n+1)^{-1}$. Hence $0 \leq I_Q(\mathbf{T}_n; \mathbf{S}_n) \leq H_Q(\mathbf{S}_n) = \log(n+1)$, which implies the claim.

Santa Fe process

The Santa Fe process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ is a sequence of variables

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

where processes $(\mathbf{K}_i)_{i \in \mathbb{Z}}$ and $(\mathbf{Z}_k)_{k \in \mathbb{N}}$ with

$$\begin{aligned} \mathbf{Q}(\mathbf{Z}_k = 0) = \mathbf{Q}(\mathbf{Z}_k = 1) = 1/2, & \quad (\mathbf{Z}_k)_{k \in \mathbb{N}} \sim \text{IID}, \\ \mathbf{Q}(\mathbf{K}_i = k) = k^{-1/\beta} / \zeta(\beta^{-1}), & \quad (\mathbf{K}_i)_{i \in \mathbb{Z}} \sim \text{IID}, \end{aligned}$$

where $\beta \in (0, 1)$ is a parameter and $\zeta(x) = \sum_{k=1}^{\infty} k^{-x}$.

Variable $\mathbf{Y} = \sum_{k=1}^{\infty} 2^{-k} \mathbf{Z}_k$ could be considered a random real parameter of the process but the distribution of the process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ is not a differentiable function of this parameter.

So, the Santa Fe process is not a $\mathbf{1}$ -parameter source.

Santa Fe process (continued)

Theorem

For the Santa Fe process, $\delta_Q^\pm = \gamma_Q^\pm = \beta$.

A process with $\delta_Q^+ \neq \delta_Q^-$

Consider a sequence of numbers $(\mathbf{a}_k)_{k \in \mathbb{N}}$ where $\mathbf{a}_k \in \{\mathbf{0}, \mathbf{1}\}$. Let

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

where $\mathbf{Y}_k = \mathbf{a}_k \mathbf{Z}_k$, whereas processes $(\mathbf{K}_i)_{i \in \mathbb{Z}}$ and $(\mathbf{Z}_k)_{k \in \mathbb{N}}$ are independent and distributed as for the original Santa Fe process.

Theorem

There exists such a sequence $(\mathbf{a}_k)_{k \in \mathbb{N}}$ that for the modified Santa Fe process, we have $\delta_Q^+ = \beta$ and $\delta_Q^- = \mathbf{0}$.

- 1 Introduction
- 2 Relating γ_P^\pm and δ_P^\pm
- 3 Relating δ_P^+ and δ_R^+
- 4 Evaluating γ_Q^\pm and δ_Q^\pm
- 5 Conclusion

Conclusion

- We have defined Hilberg exponents — the bounding rates of the power-law growth of mutual information in a process.
- There are surprisingly many meaningful Hilberg exponents, for different measures and different codes.
- We have begun sorting out order in this menagerie but surely there are some interesting hard open problems.

www.ipipan.waw.pl/~ldebowsk