

Hilberg's Conjecture: a Challenge for Machine Learning

Łukasz Dębowski
ldebowsk@ipipan.waw.pl



Institute of Computer Science
Polish Academy of Sciences

- 1 Introduction to Hilberg's conjecture
- 2 Inefficiency of Lempel-Ziv code
- 3 Vocabulary growth
- 4 Random descriptions of a random world
- 5 Conclusion

Probabilistic model of texts?

May generation of texts in natural language be described by a probabilistic model?

- This interdisciplinary question inspired a few concepts in applied mathematics:
 - Markov chains (Markov),
 - entropy (Shannon),
 - fractals (Mandelbrot),
 - algorithmic complexity (Kolmogorov).
- Common intuition: Texts are a result of a process that is neither purely **deterministic** nor purely **random** (Zipf).
- Some empirical statistical laws of language:
Zipf-Mandelbrot's law, Herdan's law, Menzerath's law.
- Some well-defined stochastic process may describe the text generation in spite of great problems with its identification.

Practical importance of language modeling

- Statistical language modeling is highly relevant for practical applications (e.g. speech recognition, machine translation).
- Bayes theorem in speech recognition:
 - $P(\mathbf{A}|\mathbf{W})$ —probability of speech \mathbf{A} corresponding to text \mathbf{W} ,
 - $P(\mathbf{W})$ —probability of text \mathbf{W} .

$$\max_{\mathbf{W}} P(\mathbf{W}|\mathbf{A}) = \max_{\mathbf{W}} P(\mathbf{A}|\mathbf{W})P(\mathbf{W}).$$

- Quality of speech recognition system depends on both models $P(\mathbf{A}|\mathbf{W})$ and $P(\mathbf{W})$.
- State-of-the-art modeling of $P(\mathbf{W})$ consists in approximating the hypothetical stochastic process by Markov chains
— far from optimal.

Need for fundamental research?

- Can theoretical research provide some insight into the practical task of statistical language modeling?
- We suppose that **randomness** of texts is constrained by the existence or the search for **meaning**...
- ...but the **meaning** itself may manifest as a form of apparent **randomness**. (Just recall halting probability, the number Ω , which is apparently random but stores an infinite amount of mathematical knowledge)
- Let's not lose hope that we may get some **theoretical** insight into probabilistic modeling of texts. We have **both** to look at empirical data and to build idealized mathematical models.

There seems to be a fundamental property of natural language, called **Hilberg's conjecture**, which can model some observations.

Stochastic processes

- Probability space: $(\Omega, \mathcal{J}, \mathbf{P})$.
- Alphabet: \mathbb{X} .
- Random variables: $\mathbf{X}_i : \Omega \rightarrow \mathbb{X}$.
- Stochastic process: $(\mathbf{X}_i)_{i \in \mathbb{Z}}$.
- Blocks: $\mathbf{X}_k^l = (\mathbf{X}_i)_{k \leq i \leq l}$.

Process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ is stationary $\stackrel{\text{def}}{\iff} \mathbf{P}(\mathbf{X}_{i+1}^{i+n})$ doesn't depend on i .

Entropy

- Entropy of a random variable:

$$H(\mathbf{X}_k^l) := \mathbb{E} \left[-\log P(\mathbf{X}_k^l) \right].$$

- It measures uncertainty of a random variable.
 - $H(\mathbf{X}_1^n) = n \log \text{card } \mathbb{X}$ — all values are equally probable.
 - $H(\mathbf{X}_1^n) = 0$ — the random variable is almost surely constant.
- Block entropy of a stationary process:

$$H(n) := H(\mathbf{X}_{i+1}^{i+n}).$$

- Entropy rate of a stationary process:

$$h = \lim_{n \rightarrow \infty} H(n)/n.$$

Hilberg's conjecture

- Shannon (1951) estimated entropy of text in English, assuming it is drawn from a stationary process.
- Hilberg (1990) replotted these estimates in the log-log scale and observed a straightish line.
- This line corresponds to

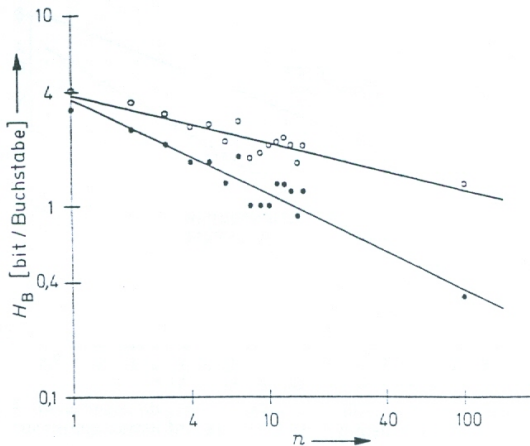
$$\mathbf{H(n)} \approx \mathbf{Bn}^\beta + \mathbf{hn}, \quad (1)$$

where $\beta \approx \mathbf{0.5}$ and $\mathbf{h} \approx \mathbf{0}$.

- Shannon provided estimates of $\mathbf{H(n)}$ for $\mathbf{n} \leq \mathbf{100}$ characters.
- Hilberg supposed that relationship (1) can be extrapolated and $\mathbf{h} = \mathbf{0}$ holds asymptotically.

Shannon's data in log-log scale (Hilberg 1990)

Conditional entropy $H(n+1) - H(n)$ vs. context length n :



The original vs. the relaxed Hilberg conjecture

Process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ is asymptotically deterministic $\stackrel{\text{def}}{\iff} \mathbf{h} = \mathbf{0}$.
 (Each random variable is a function of infinite past.)

Consider mutual information

$$I(\mathbf{X}; \mathbf{Y}) = H(\mathbf{X}) + H(\mathbf{Y}) - H(\mathbf{X}, \mathbf{Y}).$$

The **original Hilberg conjecture** is

$$H(n) \propto n^\beta, \tag{2}$$

whereas the **relaxed Hilberg conjecture** is

$$I(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) = 2H(n) - H(2n) \propto n^\beta. \tag{3}$$

Relationship (3) follows from (2) but it does not imply $\mathbf{h} = \mathbf{0}$.

Why is Hilberg's conjecture (HC) important?

- HC corroborates Zipf's insight that texts produced by humans diverge from both pure **randomness** and pure **determinism**. (In a sense, they would be both random and deterministic.)
- Relaxed HC also distinguishes natural language from **k-parameter sources**. (Some basic model of statistics.)
- HC, in its original form, implies that texts are in a sense deterministic and **infinitely compressible**. (We have to explain why modern text compressors cannot achieve that.)
- HC can be linked with **Zipf's law** and **Herdan's law**. (These are celebrated laws of quantitative linguistics.)
- Stochastic processes that satisfy HC are mathematical **terra incognita**. (Understanding their construction and properties can lead to a progress both in mathematics and applications like computational linguistics and **machine learning**.)

- 1 Introduction to Hilberg's conjecture
- 2 Inefficiency of Lempel-Ziv code
- 3 Vocabulary growth
- 4 Random descriptions of a random world
- 5 Conclusion

A skeptic's remark

- Hilberg's conjecture in its original form implies that a typical text of one **million** letters could be theoretically compressed into a string of roughly one **thousand** letters. This is far beyond the power of any known text compressor!
- How is it possible? What blocks the optimal compression? How does the optimal compression look like?
- Some idea: Modern text compressors work mostly by detecting repeated strings and replacing them with shorter identifiers. They **cannot** compress texts beyond the **maximal repetition**.
- Another idea: Giving the **ISBN number** is sufficient to identify a printed literary text that remains in cultural circulation. Thus given enough memory, "**hypercompression**" is achievable.
- Is something **similar** possible in the world of stationary stochastic processes? We suppose that it is.

Maximal repetition

Definition

The maximal repetition in text \mathbf{w} is defined as

$$L(\mathbf{w}) := \max \{ |s| : \mathbf{w} = \mathbf{x}_1 \mathbf{s} \mathbf{y}_1 = \mathbf{x}_2 \mathbf{s} \mathbf{y}_2 \text{ and } \mathbf{x}_1 \neq \mathbf{x}_2 \},$$

where \mathbf{s} , \mathbf{x}_i , and \mathbf{y}_i are substrings of text \mathbf{w} .

Maximal repetition and Hilberg's conjecture

Definition

For a random variable \mathbf{X} , topological entropy is

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

Theorem

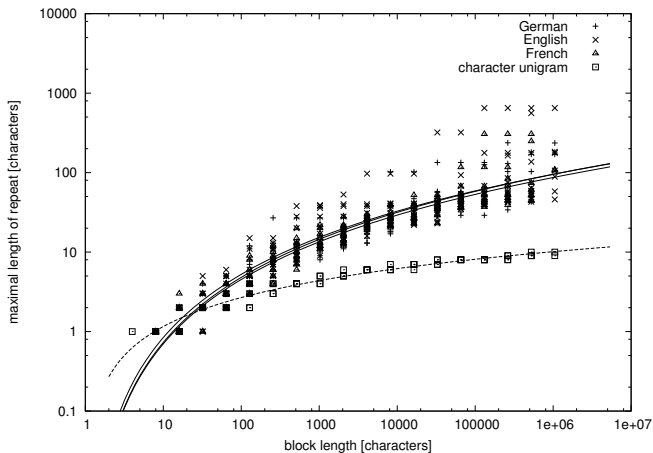
If a stationary stochastic process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ satisfies

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

for certain constants $0 < \beta < 1$ and $B > 0$ then there exists $A > 0$ such that for $\alpha = 1/\beta$ almost surely we have

$$L(\mathbf{X}_1^m) \geq A(\log m)^\alpha.$$

35 texts in 3 languages + unigram text



$$A \approx 0.093$$

$$\alpha \approx 2.64$$

Regular Hilberg processes

Definition

A stationary process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ is called a regular Hilberg process if

$$\begin{aligned} \mathbf{H}(n) &= \Theta(n^\beta), \\ \mathbb{E} \mathbf{L}(\mathbf{X}_1^n) &= \Theta((\log n)^\alpha) \end{aligned}$$

for a certain $\beta \in (0, 1)$ and $\alpha \geq 1/\beta$.

Bound for the Lempel-Ziv code

- The Lempel-Ziv (LZ) code is the oldest known universal code.

Theorem

The length of the LZ code satisfies

$$|\mathbf{C}(\mathbf{w})| \geq \frac{|\mathbf{w}|}{L(\mathbf{w}) + 1} \log \frac{|\mathbf{w}|}{L(\mathbf{w}) + 1}.$$

- Similar bounds hold for a few other known universal codes.
- Hence, for regular Hilberg processes, the length of the LZ code is orders of magnitude larger than the block entropy,

$$\mathbf{H}(\mathbf{n}) = \Theta(\mathbf{n}^\beta), \quad \mathbb{E} |\mathbf{C}(\mathbf{X}_1^n)| = \Omega \left(\frac{\mathbf{n}}{(\log \mathbf{n})^\alpha} \right).$$

Hilberg's conjecture constitutes a **challenge** for machine learning.

- 1 Introduction to Hilberg's conjecture
- 2 Inefficiency of Lempel-Ziv code
- 3 Vocabulary growth
- 4 Random descriptions of a random world
- 5 Conclusion

Herdan's law (an integrated version of Zipf's law)

- Consider texts in a natural language (such as English):
 - \mathbf{V} — the number of different words in the text,
 - \mathbf{n} — the length of the text.
- We observe Herdan's law, i.e., the relationship

$$\mathbf{V} \propto \mathbf{n}^\gamma,$$

where γ is between **0.5** a **1** depending on a text.

- We will show that Hilberg's conjecture implies Herdan's law.

A context-free grammar that generates one text

$$\left. \begin{array}{l} \mathbf{A_1} \rightarrow \mathbf{A_2A_2A_4A_5dear_childrenA_5A_3all.} \\ \mathbf{A_2} \rightarrow \mathbf{A_3youA_5} \\ \mathbf{A_3} \rightarrow \mathbf{A_4_to_} \\ \mathbf{A_4} \rightarrow \mathbf{Good_morning} \\ \mathbf{A_5} \rightarrow \mathbf{,-} \end{array} \right\} .$$

*Good morning to you,
 Good morning to you,
 Good morning, dear children,
 Good morning to all.*

The grammar-based codes

- A function Γ such that $\Gamma(\mathbf{w})$ is an admissible grammar that generates text \mathbf{w} is called a grammar transform.
- Certain grammar transforms can be turned into universal codes if we apply a certain encoding of an arbitrary grammar into a string.
- We may suppose that the number of distinct words in text \mathbf{X}_1^n can be approximated by the number of distinct nonterminals $\mathbf{V}(\mathbf{X}_1^n)$ in an admissibly minimal grammar-based code $\mathbf{C}(\mathbf{X}_1^n)$ for text \mathbf{X}_1^n .

The first result (non-zero entropy rate)

- Admissibly minimal grammar-based codes satisfy:

$$|\mathbf{C}(\mathbf{u})| + |\mathbf{C}(\mathbf{v})| - |\mathbf{C}(\mathbf{uv})| \leq \mathbf{BV}(\mathbf{uv})(1 + \mathbf{L}(\mathbf{uv})).$$

- The left hand side is an estimate of mutual information

$$\mathbf{I}(\mathbf{X}; \mathbf{Y}) = \mathbf{H}(\mathbf{X}) + \mathbf{H}(\mathbf{Y}) - \mathbf{H}(\mathbf{X}, \mathbf{Y}).$$

Theorem

Let $\mathbf{V}(\mathbf{X}_1^n)$ be the number of distinct nonterminals in an admissibly minimal grammar-based code $\mathbf{C}(\mathbf{X}_1^n)$. If for a stationary process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ over a finite alphabet with a strictly positive entropy rate we have $\mathbf{I}(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) = \Omega(n^\beta)$ for some $\beta \in (0, 1)$ then

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E} [\mathbf{V}(\mathbf{X}_1^n)(1 + \mathbf{L}(\mathbf{X}_1^n))]}{n^\beta} > 0.$$

The second result (zero entropy rate)

- We suppose that admissibly minimal grammar-based codes are universal, i.e., they satisfy $\mathbb{E} |\mathbf{C}(\mathbf{X}_1^n)| - H(n) = o(n)$. This would guarantee $\mathbb{E} |\mathbf{C}(\mathbf{X}_1^n)| = o(n)$ if $H(n) = o(n)$.

Theorem

Let $\mathbf{V}(\mathbf{X}_1^n)$ be the number of distinct nonterminals in an admissibly minimal grammar-based code $\mathbf{C}(\mathbf{X}_1^n)$. If for a stationary process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ over a finite alphabet the code satisfies $\mathbb{E} |\mathbf{C}(\mathbf{X}_1^n)| = o(n)$ then

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E} [\mathbf{V}(\mathbf{X}_1^n)(1 + L(\mathbf{X}_1^n))]}{n/\mathbb{E} L(\mathbf{X}_1^n)} > 0.$$

The second result for regular Hilberg processes

For a regular Hilberg process we have:

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E} [V(X_1^n)(1 + L(X_1^n))]}{n / (\log n)^\alpha} > 0.$$

- 1 Introduction to Hilberg's conjecture
- 2 Inefficiency of Lempel-Ziv code
- 3 Vocabulary growth
- 4 Random descriptions of a random world
- 5 Conclusion

Processes satisfying Hilberg's conjecture?

- Hilberg's conjecture:

$$H(n) \approx Bn^\beta + hn.$$

- There are a few processes that satisfy HC with $h > 0$.
- We have some idea how to construct a process that satisfies HC with $h = 0$ but have not completed the construction yet.

The Santa Fe process

A linguistic interpretation

Process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ is a sequence of random statements **consistently** describing the state of an “earlier drawn” random object $(\mathbf{Z}_k)_{k \in \mathbb{N}}$. $\mathbf{X}_i = (\mathbf{k}, \mathbf{z})$ asserts that the \mathbf{k} -th bit of $(\mathbf{Z}_k)_{k \in \mathbb{N}}$ has value $\mathbf{Z}_k = \mathbf{z}$.

Let a process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ have the form

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

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

$$\begin{aligned} P(\mathbf{K}_i = \mathbf{k}) &= \mathbf{k}^{-1/\beta} / \zeta(\beta^{-1}), & \beta &\in (0, 1), \\ P(\mathbf{Z}_k = \mathbf{z}) &= \frac{1}{2}, & \mathbf{z} &\in \{0, 1\}. \end{aligned}$$

We have $\lim_{n \rightarrow \infty} I(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) / n^\beta > 0$.

A mixing Santa Fe process

A linguistic interpretation

Object $(\mathbf{Z}_{ik})_{k \in \mathbb{N}}$ described by text $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ is a function of time i .

Let a process $(\mathbf{X}_i)_{i \in \mathbb{Z}}$ have the form

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

where $(\mathbf{K}_i)_{i \in \mathbb{Z}}$ and $(\mathbf{Z}_{ik})_{i \in \mathbb{Z}}$, $\mathbf{k} \in \mathbb{N}$, are independent,

$$\mathbf{P}(\mathbf{K}_i = \mathbf{k}) = \mathbf{k}^{-1/\beta} / \zeta(\beta^{-1}), \quad (\mathbf{K}_i)_{i \in \mathbb{Z}} \sim \text{IID},$$

whereas $(\mathbf{Z}_{ik})_{i \in \mathbb{Z}}$ are Markov chains with

$$\begin{aligned} \mathbf{P}(\mathbf{Z}_{ik} = \mathbf{z}) &= \frac{1}{2}, \\ \mathbf{P}(\mathbf{Z}_{ik} = \mathbf{z} | \mathbf{Z}_{i-1, \mathbf{k}} = \mathbf{z}) &= 1 - \mathbf{p}_{\mathbf{k}}. \end{aligned}$$

We have $\lim_{n \rightarrow \infty} \mathbf{I}(\mathbf{X}_1^n; \mathbf{X}_{n+1}^{2n}) / n^\beta > 0$ for $\mathbf{p}_{\mathbf{k}} \leq \mathbf{P}(\mathbf{K}_i = \mathbf{k})$.

- 1 Introduction to Hilberg's conjecture
- 2 Inefficiency of Lempel-Ziv code
- 3 Vocabulary growth
- 4 Random descriptions of a random world
- 5 Conclusion

Conclusions

- Hilberg's conjecture is a hypothesis about a power law growth of block entropy for texts in natural language.
- It may have profound implications for text compression, statistical natural language modeling, and machine learning.
- Further fundamental mathematical research is needed (models of processes, entropy estimation methods).

www.ipipan.waw.pl/~ldebowsk