

# Zipf's Law: Balancing Signal Usage Cost and Communication Efficiency

Christoph Salge  
Nihat Ay  
Daniel Polani  
Mikhail Prokopenko

SFI WORKING PAPER: 1013-10-033

SFI Working Papers contain accounts of scientific work of the author(s) and do not necessarily represent the views of the Santa Fe Institute. We accept papers intended for publication in peer-reviewed journals or proceedings volumes, but not papers that have already appeared in print. Except for papers by our external faculty, papers must be based on work done at SFI, inspired by an invited visit to or collaboration at SFI, or funded by an SFI grant.

©NOTICE: This working paper is included by permission of the contributing author(s) as a means to ensure timely distribution of the scholarly and technical work on a non-commercial basis. Copyright and all rights therein are maintained by the author(s). It is understood that all persons copying this information will adhere to the terms and constraints invoked by each author's copyright. These works may be reposted only with the explicit permission of the copyright holder.

[www.santafe.edu](http://www.santafe.edu)



SANTA FE INSTITUTE

# Zipf's Law: Balancing Signal Usage Cost and Communication Efficiency

Christoph Salge <sup>\*1</sup>, Nihat Ay <sup>†2,3</sup>, Daniel Polani <sup>‡1</sup>, and Mikhail Prokopenko <sup>§4,5,6</sup>

<sup>1</sup>*School of Computer Science, University of Hertfordshire, Hatfield AL10 9AB, UK*

<sup>2</sup>*Max Planck Institute for Mathematics in the Sciences, Inselstraße 22-26, Leipzig 04103, Germany*

<sup>3</sup>*Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA*

<sup>4</sup>*CSIRO Computational Informatics, PO Box 76, Epping, NSW 1710, Australia*

<sup>5</sup>*School of Physics, University of Sydney, Physics Rd, Camperdown NSW 2050, Australia*

<sup>6</sup>*Department of Computing, Macquarie University, E6A Level 3, Eastern Rd, Macquarie Park, NSW 2113, Australia*

October 25, 2013

## Abstract

We propose a model that explains the reliable emergence of power laws (e.g., Zipf's law) during the development of different human languages. The model incorporates the principle of least effort in communications, minimizing a combination of the information-theoretic communication inefficiency and direct signal cost. We prove a general relationship, for all optimal languages, between the signal cost distribution and the resulting distribution of signals. Zipf's law then emerges for logarithmic signal cost distributions, which is the cost distribution expected for words constructed from letters or phonemes.

*Keywords:* Zipf's law, power-law distribution, communication efficiency, signal cost, least effort, information theory.

## 1 Introduction

Zipf's law [20] for natural languages states that the frequency  $p(s)$  of a given word  $s$  in a large enough corpus of a (natural) language is inversely proportional to the word's frequency rank. Zipf's law postulates a power-law distribution for languages with a specific power law exponent  $\beta$ , so if  $s_t$  is the  $t$ -th most common word, then its frequency is proportional to

$$p(s_t) \sim \frac{1}{t^\beta}, \quad (1)$$

with  $\beta \approx 1$ . Empirical data suggests that the power law holds across a variety of natural languages [2], but the exponent  $\beta$  can vary, depending on the language and the context, with a usual value of  $\beta \approx 2$  [8]. While the adherence to this "law" in different languages suggests a underlying common principle or mechanism, a generally accepted explanation for this phenomenon is still lacking [15].

Several papers [12, 13, 18] suggest that random texts already display a power law distribution sufficient to explain Zipf's law, but a detailed analysis [6] with different

---

\*c.salge@herts.ac.uk

†nay@mis.mpg.de

‡d.polani@herts.ac.uk

§mikhail.prokopenko@csiro.au

statistical tests rejects this hypothesis and argues, that there is a “meaningful” mechanism at play, which causes this distribution across different natural languages.

An alternative model by Cancho and Solé [11] follows the original idea of Zipf [20], by modelling the evolution of language based on the principle of least effort, where the assignment of words to concepts is optimized to minimize a weighted sum of speaker and listener effort. While simulations of the model produce distributions which qualitatively resemble power laws, a detailed mathematical investigation [15] reveals that the optimal solution of this model is, in fact, not following a power law; thus, the power law characteristics of the simulation results seems to be an artefact of the particular optimization model utilized.

More recent models [1, 19] prove that power laws result from an alternative cost function, but their cost function only optimizes the output distribution, and ignores any relationship between used words and intended meaning. This makes it a questionable model for human language (similar to the models with random text) as it does not account for communication efficiency, i.e., the model is not sensitive to how much information the words contain about the referenced concepts.

Thus, to our knowledge, the question of how to achieve power laws in human language from a least effort principle is still not satisfactorily solved. Nevertheless, the idea from [11, 20] to explain power laws as the result of an evolutionary optimization process that minimizes some form of language usage cost remains attractive. In this vein, we here present an alternative model for the least effort principle in language: we minimize a cost function consisting of communication inefficiency and an inherent cost for each signal (word). To avoid past pitfalls of statistical analysis when looking for power laws [3], we offer mathematical proof that any optimal solution for our cost function necessarily realizes a power law distribution, as long as the underlying cost function for the signals increases logarithmically, if the signals are ordered according to cost rank. The result generalizes beyond this as we can state a general relationship between the cost structure of the individual signals and the resulting optimal distribution of the language signals.

The resulting insights may be of interest beyond the confines of power-law structures and offer an opportunity to study optimality conditions in other types of self-organizing coding systems, for instance in the case of the genetic code.

## 2 Model

We will use a model, similar to that used by Ferrer i Cancho and Solé [11], which considers languages as an assignment of symbols to objects, and then optimizes this assignment function in regard to some form of combined speaker and listener effort. The language emerging from our model is also based on the optimality principle of least effort in communication, but uses a different cost function.

The model has a set of  $n$  signals  $S$  and a set of  $m$  objects  $R$ . Signals are used to reference objects, and a language is defined by how the speaker assigns signals to objects, i.e. by the relation between signals and objects. The relation between  $S$  and  $R$  in this model can be expressed by a binary matrix  $A$ , where an element  $a_{i,j} = 1$  if and only if signal  $s_i$  refers to object  $r_j$ .

This model allows one to represent both *polysemy* (that is, the capacity for a signal to have multiple meanings by referring to multiple objects), and *synonymy*, where multiple signals refer to the same object. The relevant probabilities are then defined as follows:

$$p(s_i|r_j) = \frac{a_{i,j}}{\omega_j} \quad (2)$$

where  $\omega_j$  is the number of synonyms for object  $r_j$ , that is  $\omega_j = \sum_i a_{i,j}$ . Thus, the probability of using a synonym is equally distributed over all synonyms referring

to a particular object. Importantly, it is also assumed that  $p(r_j) = \frac{1}{m}$  is uniformly distributed over the objects, leading to a joint distribution:

$$p(s_i, r_j) = p(r_j) p(s_i|r_j) = \frac{a_{i,j}}{m\omega_j} . \quad (3)$$

In the previous model [11] each language has a cost based on a weighted combination of speaker and listener effort. The effort for the listener should be low if the received signal  $s_i$  leaves little ambiguity as to what object  $r_j$  is referenced, so there is little chance that the listener misunderstands what the speaker wanted to say. In the model of Ferrer i Cancho and Solé [11], the cost for listening to a specific signal  $s_i$  is expressed by the conditional entropy:

$$H_{R|s_i}(p) \equiv - \sum_{j=1}^m p(r_j|s_i) \log_m p(r_j|s_i) . \quad (4)$$

The overall effort for the listener is then dependent on the probability of each signal and the effort to decode it, that is

$$H_{R|S}(p) \equiv \sum_{i=1}^n p(s_i) H_{R|s_i} . \quad (5)$$

Ferrer i Cancho and Solé argue that the listener effort is minimal when this entropy is minimal, in which case there is a deterministic mapping between signals and objects.

The effort for the speaker is expressed by the entropy  $H_S$ , which is, as the term in (5), bound between 0 and 1, via the log with respect to  $n$ :

$$H_S(p) \equiv - \sum_{i=1}^n p(s_i) \log_n p(s_i) . \quad (6)$$

Ferrer i Cancho and Solé then combine the listener's and speaker's efforts within the cost function  $\Omega_\lambda$  as follows:

$$\Omega_\lambda = \lambda H_{R|S} + (1 - \lambda) H_S , \quad (7)$$

with  $0 \leq \lambda \leq 1$ .

It can be shown that the cost function  $\Omega_\lambda$  given by (7) is a specific case of a more general *energy* function that a communication system must minimize [9, 15]

$$\Omega_\lambda^0 = -\lambda I(S; R) + (1 - \lambda) H_S , \quad (8)$$

where the mutual information  $I(S; R) = H_R - H_{R|S}$  captures the communication efficiency, i.e. how much information the signals contain about the objects. This energy function better accounts for subtle communication efforts [10], since  $H_S$  is arguably both a source of effort for the speaker and the listener because the word frequency affects not only word production but also recognition of spoken and written words. The component  $I(S; R)$  also implicitly accounts for both  $H_{S|R}$  (a measure of the speaker's effort of coding objects) and  $H_{R|S}$  (i.e., a measure of the listener's effort of decoding signals). It is easy to see that

$$\Omega_\lambda^0 = -\lambda H_R + \lambda H_{R|S} + (1 - \lambda) H_S = -\lambda H_R + \Omega_\lambda , \quad (9)$$

and so when the entropy  $H_R$  is constant, e.g. under the uniformity condition  $p(r_j) = \frac{1}{m}$ , the more generic energy function  $\Omega_\lambda^0$  reduces to the specific  $\Omega_\lambda$ .

We propose instead another cost function that not only produces optimal languages exhibiting power laws, but also retains the clear intuition of generic energy functions which typically reflect the global quality of a solution. Firstly, we represent the communication inefficiency by the information distance, the Rokhlin metric,  $H_{S|R} + H_{R|S}$  [17, 4]. This distance is often more sensitive than  $-I(S; R)$  in measuring the “disagreements” between variables, especially in the case when one information source is contained within another [16].

Secondly, we define the signal usage effort by introducing an explicit cost function  $c(s_i)$ , which assigns each signal a specific cost. The signal usage cost for a language is then the weighted average of this signal specific cost:

$$\sum_{i=1}^n p(s_i) c(s_i). \quad (10)$$

This is motivated by the basic idea that words have an intrinsic cost associated with using (speaking, writing, hearing, reading) them. To illustrate, a version of English where each use of the word “P” is replaced with “Antidisestablishmentarianism” and vice versa should not have the same signal usage cost as normal English. The optimal solution considering the signal usage cost alone would be to reference every object with the cheapest signal.

The overall cost function for a language  $\Omega_\lambda^c$  is the energy function trading off the communicative inefficiency with the signal usage cost, with  $0 \leq \lambda \leq 1$  trading off the efforts as follows:

$$\Omega_\lambda^c(p) = \lambda(H_{S|R}(p) + H_{R|S}(p)) + (1 - \lambda) \sum_{i=1}^n p(s_i) c(s_i). \quad (11)$$

A language can be optimized for different values of  $\lambda$ , weighting the respective costs. In the extreme cases, only the signal usage cost ( $\lambda = 0$ ) or the communication inefficiency ( $\lambda = 1$ ) is considered. Following the principle of least effort, we aim to determine the properties of those languages that have minimal cost according to  $\Omega_\lambda^c$ .

### 3 Results

First of all, we establish that all local minimizers, and hence all global minimizers, of the cost function (11) are solutions without synonyms. Formally, we obtain the following result.

**Theorem 1.** *Each local minimizer of the function*

$$\mathcal{C} \rightarrow \mathbb{R}, \quad p \mapsto \Omega_\lambda^c(p),$$

where

$$\mathcal{C} := \{p \in \mathcal{P}(S \times R) : p(r_j) = \sum_i p(s_i, r_j) = \frac{1}{m} \text{ for all } j\},$$

and  $\Omega_\lambda^c(p)$  is specified by the equation (11),  $0 < \lambda \leq 1$ , can be represented as a function  $f : R \rightarrow S$  such that

$$p(s_i, r_j) = \begin{cases} 1/m & \text{if } s_i = f(r_j); \\ 0 & \text{otherwise.} \end{cases} \quad (12)$$

The proof is given in Appendix 1. Note that each solution, i.e. each distribution  $p$  in expression (3), corresponds to a matrix  $A$  (henceforth called *minimizer matrix*) which is given in terms of function  $f$  as follows:

$$a_{i,j} = \begin{cases} 1 & \text{if } s_i = f(r_j); \\ 0 & \text{otherwise.} \end{cases} \quad (13)$$

The main outcome of this observation is that the analytical minimization of the suggested cost function results in solutions without synonyms — since any function  $f$  precludes multiple signals  $s$  referring to the same object  $r$ . That is, each column in the minimizer matrix has precisely one non-zero element. Polysemy is allowed within the solutions.

We need the following lemma as an intermediate step towards deriving the analytical relationship between the specific word cost  $c(s)$  and the resulting distribution  $p(s)$ .

**Lemma 2.** *For each solution  $p$  minimizing the function  $\Omega_\lambda^c$ ,*

$$H_{R|S} + \frac{1}{\log_n m} H_S = 1 . \quad (14)$$

The proof follows from the joint entropy representations

$$H_{S,R} = \frac{H_{R|S}}{1 + \log_m n} + \frac{H_S}{1 + \log_n m} = \frac{H_{S|R}}{1 + \log_n m} + \frac{H_R}{1 + \log_m n} , \quad (15)$$

noting that for each minimal solution  $H_{S|R} = 0$ , while  $H_R = 1$  under the uniformity constraint  $p(r_j) = \frac{1}{m}$ .

**Corollary 3.** *If  $n = m$ ,  $H_{R|S} + H_S = 1$ .*

Using this lemma, and noting that each such solution represented as a function  $f : R \rightarrow S$  has the property  $H_{S|R} = 0$ , we reduce the equation (11) to

$$\Omega_\lambda^c(p) = \lambda \left( 1 - \frac{1}{\log_n m} H_S(p) \right) + (1 - \lambda) \sum p(s_i) c(s_i) \quad (16)$$

$$= \lambda + \frac{\lambda}{\log_n m} \sum p(s_i) \log_n p(s_i) + (1 - \lambda) \sum p(s_i) c(s_i). \quad (17)$$

Varying with respect to  $p(s_i)$ , under the constraint  $\sum p(s_i) = 1$ , yields the extremality condition

$$\frac{\lambda}{\log_n m} (\log_n p(s_i) + 1) + (1 - \lambda) c(s_i) - \kappa' = 0 \quad (18)$$

for some Lagrange multiplier  $\kappa'$ . The minimum is achieved when

$$p(s_i) = \kappa e^{-\beta c(s_i)} , \quad (19)$$

where

$$\beta = \frac{1 - \lambda}{\lambda} \log_n m , \quad (20)$$

$$\kappa = \frac{1}{\sum e^{-\beta c(s_j)}} . \quad (21)$$

In addition, we require

$$c(s_i) = \ln m_i \quad (22)$$

for some integer  $m_i$  such that  $\sum m_i = m$ . The last condition ensures that the minimal solutions  $p(s_i)$  correspond to functions  $p(s_i, r_j)$  (i.e., minimizer matrices without synonyms). In other words, the marginal probability (19) without the condition (22) may not concur with the probability  $p(s_i, r_j)$  that represents a minimizer matrix under the uniformity constraint  $p(r_j) = \frac{1}{m}$ .

Under the condition (22), we have  $p(s_i) = \kappa m_i^{-\beta}$ , while  $\kappa = 1/\sum m_i^{-\beta}$ . In general, one may relax the condition (22), specifying instead an upper-bounded error of approximating the minimal solution by any  $p(s_i) = \kappa e^{-\beta c(s_i)}$  which would then allow for arbitrary cost functions  $c(s)$ .

Interestingly, the optimal marginal probability distribution (19) is the Gibbs measure with the energy  $c(s_i)$ , while the parameter  $\beta$  is, thermodynamically, the inverse temperature. It is well-known that the Gibbs measure is the unique measure maximizing the entropy for a given expected energy, and appears in many solutions outside of thermodynamics, leading to representations of the respective energy function as a product of factors.

Let us now consider some special cases. For the case of equal effort, i.e.  $\lambda = 0.5$ , and  $n = m$ , the solution simplifies to  $\beta = 1$  and  $p(s_i) = \kappa m_i^{-1}$ , where  $\kappa = 1/\sum m_i^{-1}$ .

Another important special case is given by the cost function  $c(s_i) = \ln \rho_i/N$ , where  $\rho_i$  is the rank of symbol  $s_i$ , and  $N$  is a normalization constant equal to  $\frac{n(n+1)}{2m}$  (so that  $\sum \rho_i/N = m$ ). In this case, the optimal solution is attained when

$$p(s_i) = \frac{\kappa N^\beta}{\rho_i^\beta} \quad (23)$$

with

$$\kappa = \frac{1}{N^\beta \sum \rho_j^{-\beta}}. \quad (24)$$

This means that a power law with the exponent  $\beta$ , specified by equation (20), is the optimal solution in regard to our cost function (11) if the signal usage cost increases logarithmically. In this case, the exponent  $\beta$  depends on the system's size ( $n$  and  $m$ ) and the efforts' trade-off  $\lambda$ . Importantly, this derivation shows a connection between scaling in languages and thermodynamics: if the signal usage cost increases logarithmically, then the scaling exponent of the resulting power law is given by the corresponding inverse temperature.

Zipf's law (a power law with exponent  $\beta = 1$ ) is then nothing but a special case for systems that satisfy  $\log_n m = \frac{\lambda}{1-\lambda}$ . For instance, for square matrices, Zipf's law results from the optimal languages which satisfy equal efforts, i.e.,  $\lambda = 0.5$ . The importance of equal cost was emphasized in earlier works [15, 5]. The exponent defined by (20) changes with the system size ( $n$  or  $m$ ), and so the resulting power law "adapts" to linguistic dynamics and language evolution in general.

In summary, the derived relationship expresses the optimal probability  $p(s)$  in terms of the usage cost  $c(s)$ , yielding Zipf's law when this cost is logarithmically distributed over the symbols.

## 4 Discussion

To explain the emergence of power laws for signal selection, we need to explain why the cost function of the signals would increase logarithmically, if the signals are ordered by their cost rank. This can be motivated, across a number of languages, by assuming that signals are in fact words, which are made up of letters from a finite alphabet; or in regard to spoken language, are made of from a finite set of phonemes. Compare [14], in which Nowak and Krakauer demonstrate how the error

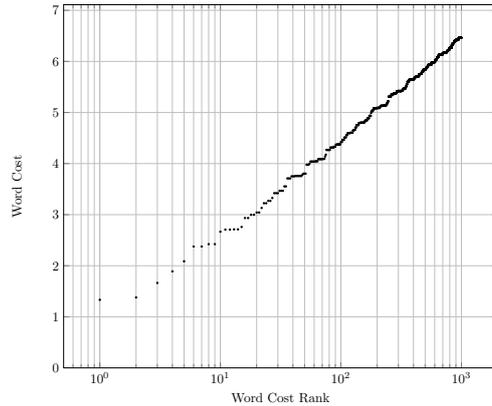


Figure 1: A log-plot of the 1000 cheapest words created from a 10 letter alphabet, ordered by their cost rank. Word cost is a sum of individual letter cost, and letter cost is between 1.0 and 2.0 units.

limits of communication with a finite list of phonemes can be overcome by combining phonemes into words.

Lets assume that each letter (or phoneme) has an inherent cost which is approximate to a unit letter cost. Furthermore, assume that the cost of a word roughly equals the sum of its letter costs. A language with an alphabet of size  $a$  then has  $a$  unique one letter words which the approximate cost of one,  $a^2$  two letter words with an approximate cost of two,  $a^3$  three letter words with a cost of three, etcetera. If we rank these words by their cost, then their cost will increase approximately logarithmically with their cost rank. To illustrate, Figure 1 is a plot of the 1000 cheapest unique words formed with a ten letter alphabet (with no word length restriction), where each letter has a random cost between 1.0 and 2.0. The first few words deviate from the logarithmic cost function, as their cost only depends on the letter cost itself, but the latter words closely follow a logarithmic function.

This signal usage cost can be interpreted in different ways. In spoken language it might simply be the time needed to utter a word, which makes it a cost both for the listener and the speaker. In written language it might be the effort to write a word, or the bandwidth needed to transmit it, in which case it is a speaker cost. On the other hand, if one is reading a written text, then the length of the words might translate into “listener” cost again. In general, the average signal usage cost corresponds to the effort of using a specific language to communicate for all involved parties. This differs from the original least effort idea, which balances listener and speaker effort [20]. In our model we balance the general effort of using the language with the communication efficiency, which creates a similar tension, as described in [11], between using a language that only uses one signal, and a language that references every object with its own signal. If only communication efficiency was relevant, then each object would have its own signal. Conversely, if only cost mattered, then all objects would be referenced by the same cheapest signal. Balancing these two components with a weighting factor  $\lambda$  yields power laws, where  $\beta$  varies with changes in the weighting factor. This is in contrast to the model in [11], where power laws were only found in a phase transition along the weighting factor. Also, in [8] Cancho discusses how some variants of language (military, children) have  $\beta$  values that deviate from the  $\beta$  value of their base language, which could indicate that the effort of language production or communication efficiency is weighted differently in these cases, resulting in different optimal solutions, which are power laws with other values for  $\beta$ .

We noted earlier that there are other options to produce power laws, which are insensitive to the relationship between objects and signals. Baek et al. [1] obtain a power law by minimizing the cost function  $I_{cost} = -H_S + \langle \log s \rangle + \log N$ , where  $\langle \log s \rangle = \sum p(s_i) \log(s_i)$ , and  $\log(s_i)$  is interpreted as the logarithm of the index of  $s_i$  (specifically, its rank). Their argument that this cost function follows from a more general cost function  $H_{R|S} = -I(S; R) + H_R$ , where  $H_R$  is constant, is undermined by their unconventional definition of conditional probability (cf. Appendix A [1]). Specifically, this probability is defined as  $p(r|s) = \frac{\delta_{s^t(r),s}}{sN(s)}$ , where  $N(s)$  is the number of objects to which signal  $s$  refers. This definition not only requires some additional assumptions in order to make  $p(r|s)$  a conditional probability, but also implicitly embeds the “cost” of symbol  $s$  within the conditional probability  $p(r|s)$ , by dividing it by  $s$ . Thus, we are left with the cost function  $I_{cost}$  *per se*, not rigorously derived from a generic principle, and this cost function ignores joint probabilities and the communication efficiency in particular.

A very similar cost function was offered by Visser [19], who suggested to maximize  $H_S$  subject to a constraint  $\langle \log s \rangle = \chi$ , for some constant  $\chi$ . Again, this maximization produces a power law, and again we may note that the cost function and the constraint used in the derivation do not capture communication efficiency or trade-offs between speaker and listener, omitting joint probabilities as well.

Finally, we would like to point out that the cost function  $-H_S + \langle \log s \rangle$  is equivalent to the cost function  $H_{R|S} - H_{S|R} + \langle \log s \rangle$ , under constant  $H_R$ . This expression reveals another important drawback of minimizing  $-H_S + \langle \log s \rangle$  directly: while minimizing  $H_{R|S}$  reduces the ambiguity of polysemy, minimizing  $-H_{S|R}$  explicitly “rewards” the ambiguity of synonyms. In other words, languages obtained by minimizing such a cost directly do exhibit a power law, but mostly at the expense of potentially unnecessary synonyms. It may be argued that abundance of synonyms is a feature in some early and proto-languages. This issue has been a somewhat controversial topic: for example, Herder’s conjecture that the wealth of synonyms is a trait in “original” languages: “alongside of real poverty it has the most unnecessary superfluity” [7], dates back to 18th century and is based on an uninformed time-perspective. Nevertheless, it is now more or less accepted that proto-languages did offer a multitude of synonyms, lessening the effort of speakers, and shifting the communication burden to listeners. However, as languages mature, the communicative efficiency and the balance between speaker’s and listener’s efforts become a more significant driver, and so the simplistic cost function  $-H_S + \langle \log s \rangle$  can no longer be justified.

In contrast, the cost function proposed in this paper  $H_{R|S} + H_{S|R} + \langle \log s \rangle$  reduces to  $-H_S + \langle \log s \rangle$  only *after* minimizing over the joint probabilities  $p(s, r)$ . Importantly, it captures communication (in)efficiency and average signal usage explicitly, balancing out different aspects of the communication trade-offs and representing the concept of least effort in a principled way. The resulting solutions do not contain synonyms, which disappear at the step of minimizing over  $p(s, r)$ , and so correspond to “perfect”, maximally efficient and balanced, languages. The fact that even these languages exhibit power (Zipf’s) laws is a manifestation of the continuity of scale-freedom in language evolution, along the refinement of cost functions representing the least effort principle: as long as the language evolves closely to the optima of the prevailing cost function, power laws will be adaptively maintained.

In conclusion, our paper addresses the long-held conjecture that the principle of least effort provides a plausible mechanism for generating power laws. For this, we interpret the effort in suitable information-theoretic terms and prove that its global minimum produces Zipf’s law. Our formalization enables a derivation of languages which are optimal with respect to both the communication inefficiency and direct signal cost. The proposed combination of these two factors within a generic cost

function is an intuitive and powerful method to capture the trade-offs intrinsic to least-effort communication.

## Acknowledgments

CS and DP were supported by the European Commission as part of the CORBYS (Cognitive Control Framework for Robotic Systems) project under contract FP7 ICT-270219. The views expressed in this paper are those of the authors, and not necessarily those of the consortium.

# Appendices

## A Appendix 1

**Theorem 1.** *Each local minimizer of the function*

$$\mathcal{C} \rightarrow \mathbb{R}, \quad p \mapsto \Omega_\lambda^c(p),$$

where

$$\mathcal{C} := \left\{ p \in \mathcal{P}(S \times R) : p(r_j) = \sum_i p(s_i, r_j) = \frac{1}{m} \text{ for all } j \right\},$$

and  $\Omega_\lambda^c(p)$  is specified by the equation (11),  $0 < \lambda \leq 1$ , can be represented as a function  $f : R \rightarrow S$  such that

$$p(s_i, r_j) = \begin{cases} 1/m & \text{if } s_i = f(r_j); \\ 0 & \text{otherwise.} \end{cases} \quad (25)$$

In order to prove this theorem, we establish a few preliminary propositions.<sup>1</sup>

### A.1 Extreme points

The extreme points of  $\mathcal{C}$  are specified by the following proposition.

**Proposition 2.** *The set  $\mathcal{C}$  has the extreme points*

$$\text{Ext}(\mathcal{C}) = \left\{ p \in \mathcal{P}(S \times R) : p(s_i, r_j) = \frac{1}{m} \delta_{f(r_j)}(s_i) \right\},$$

where  $f$  is a function  $R \rightarrow S$ .

**Proof.** Consider the convex set

$$\mathcal{T} = \left\{ A = (a_{i|j})_{i,j} \in \mathbb{R}^{m \cdot n} : a_{i|j} \geq 0 \text{ for all } i, j, \right.$$

$$\left. \text{and } \sum_i a_{i|j} = 1 \text{ for all } j \right\}$$

of transition matrices. The extreme points of  $\mathcal{T}$  are given by functions  $f : j \mapsto i$ . More precisely, each extreme point has the structure

$$a_{i|j} = \delta_{f(j)}(i).$$

---

<sup>1</sup>These results are obtained by Nihat Ay.

Now consider the map  $\varphi : \mathcal{T} \rightarrow \mathcal{C}$  that maps each matrix  $A = (a_{i|j})_{i,j}$  to the probability vector

$$p(s_i, r_j) := \frac{1}{m} a_{i|j}, \quad \text{for all } i, j.$$

This map is bijective and satisfies  $\varphi((1-t)A + tB) = (1-t)\varphi(A) + t\varphi(B)$ . Therefore, the extreme points of  $\mathcal{C}$  can be identified with the extreme points of  $\mathcal{T}$ . ■

## A.2 Concavity

Consider the set  $S = \{s_1, \dots, s_n\}$  of signals with  $n$  elements and the set  $R = \{r_1, \dots, r_m\}$  of  $m$  objects, and denote with  $\mathcal{P}(S \times R)$  the set of all probability vectors  $p(s_i, r_j)$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq m$ . We define the following functions on  $\mathcal{P}(S \times R)$ :

$$H_{S|R}(p) := - \sum_j p(r_j) \sum_i p(s_i|r_j) \log_n p(s_i|r_j),$$

and

$$H_{R|S}(p) := - \sum_i p(s_i) \sum_j p(r_j|s_i) \log_m p(r_j|s_i).$$

**Proposition 3.** *All three functions  $H_{R|S}$ ,  $H_{S|R}$ , and*

$$\langle c \rangle : p \mapsto \sum_i p(s_i) c(s_i)$$

*that are involved in the definition of  $\Omega_\lambda^c$  are concave in  $p$ . Furthermore, the restriction of  $H_{S|R}$  to the set  $\mathcal{C}$  is strictly concave.*

**Proof.** The statements follow from well-known convexity properties of the entropy and the relative entropy.

(1) *Concavity of  $H_{R|S}$ :* We rewrite the function  $H_{R|S}$  as

$$\begin{aligned} H_{R|S}(p) &= - \sum_i p(s_i) \sum_j p(r_j|s_i) \log_m p(r_j|s_i) \\ &= - \sum_{i,j} p(s_i, r_j) \log_m \frac{p(s_i, r_j)}{\sum_j p(s_i, r_j)} \\ &= - \sum_{i,j} p(s_i, r_j) \log_m \frac{p(s_i, r_j)}{m \frac{1}{m} \sum_j p(s_i, r_j)} \\ &= - \sum_{i,j} p(s_i, r_j) \log_m \frac{p(s_i, r_j)}{\frac{1}{m} \sum_j p(s_i, r_j)} + 1. \end{aligned}$$

The concavity of  $H_{R|S}$  now follows from the joint convexity of the relative entropy  $(p, q) \mapsto D(p||q) = \sum_{i,j} p(s_i, r_j) \log_m \frac{p(s_i, r_j)}{q(s_i, r_j)}$ .

(2) *Concavity of  $H_{S|R}$ :* The concavity of  $H_{S|R}$  follows by the same arguments as

in (1). We now prove the strict concavity of its restriction to  $\mathcal{C}$ .

$$\begin{aligned}
H_{S|R}(p) &= - \sum_j p(r_j) \sum_i p(s_i|r_j) \log_n p(s_i|r_j) \\
&= - \sum_{i,j} p(s_i, r_j) \log_n \frac{p(s_i, r_j)}{p(r_j)} \\
&= - \sum_{i,j} p(s_i, r_j) \log_n \frac{p(s_i, r_j)}{\frac{1}{m}} \\
&= - \sum_{i,j} p(s_i, r_j) \log_n p(s_i, r_j) - \log_n m.
\end{aligned}$$

The strict concavity of  $H_{R|S}$  now follows from the strict concavity of the Shannon entropy.

(2) *Concavity of  $\langle c \rangle$* : This simply follows from the fact that  $\langle c \rangle$  is an affine function and therefore concave and convex at the same time.  $\blacksquare$

With a number  $0 < \lambda \leq 1$ , we now consider the function

$$\Omega_\lambda^c(p) = \lambda (H_{R|S}(p) + H_{S|R}(p)) + (1 - \lambda) \sum_i p(s_i) c(s_i).$$

From Proposition 3, it immediately follows that  $\Omega_\lambda^c$  also has corresponding concavity properties.

**Corollary 4.** *For  $0 \leq \lambda \leq 1$ , the function  $\Omega_\lambda^c$  is concave in  $p$ , and, if  $\lambda > 0$ , its restriction to the convex set  $\mathcal{C}$  is strictly concave.*

### A.3 minimizers

We have the following direct implication of Corollary 4.

**Corollary 5.** *Let  $0 < \lambda \leq 1$  and let  $p$  be a local minimizer of the map*

$$\mathcal{C} \rightarrow \mathbb{R}, \quad p \mapsto \Omega_\lambda^c(p).$$

*Then  $p$  is an extreme point of  $\mathcal{C}$ .*

**Proof.** This directly follows from the strict concavity of this function.  $\blacksquare$

Together with Proposition 2, this implies Theorem 1, our main result on minimizers of the restriction of  $\Omega_\lambda^c$  to the convex set  $\mathcal{C}$ .

We finish this analysis by addressing the problem of minimizing  $\Omega_\lambda^c$  on a discrete set. In order to do so, consider the set of 0/1-matrices that have at least one “1”-entry in each column:

$$\mathcal{S} := \left\{ (a_{i,j}) \in \{0, 1\}^{n \cdot m} : \sum_i a_{i,j} \geq 1 \text{ for all } j \right\}.$$

This set can naturally be embedded into the set  $\mathcal{T}$ , which we have considered in the proof of Proposition 2:

$$\iota : \mathcal{S} \hookrightarrow \mathcal{T}, \quad (a_{i,j})_{i,j} \mapsto a_{i|j} := \frac{a_{i,j}}{\sum_i a_{i,j}}.$$

Together with the map  $\varphi : \mathcal{T} \rightarrow \mathcal{C}$  we have the injective composition  $\varphi \circ \iota$ . From Proposition 2 it follows that the extreme points of  $\mathcal{C}$  are in the image of  $\varphi \circ \iota$ . Furthermore, Corollary 5 implies that all local, and therefore also all global, minimizers of  $\Omega_\lambda^c$  are in the image of  $\varphi \circ \iota$ . The previous work of Ferrer i Cancho and Solé [11] refers to the minimization of a function on the discrete set  $\mathcal{S}$ :

$$\tilde{\Omega}_\lambda^c := \Omega_\lambda^c \circ \varphi \circ \iota : \mathcal{S} \rightarrow \mathbb{R}.$$

It is not obvious how to relate local minimizers of this function, with an appropriate notion of locality in  $\mathcal{S}$ , to local minimizers of  $\Omega_\lambda^c$ . However, we have the following obvious relation between global minimizers.

**Corollary 6.** *A point  $p \in \mathcal{C}$  is a global minimizer of  $\Omega_\lambda^c$  if and only if it is in the image of  $\varphi \circ \iota$  and  $(\varphi \circ \iota)^{-1}(p)$  globally minimizes  $\tilde{\Omega}_\lambda^c$ .*

## References

- [1] Seung Ki Baek, Sebastian Bernhardsson, and Petter Minnhagen. Zipf’s law unzipped. *New Journal of Physics*, 13(4):043004, 2011.
- [2] VK Balasubrahmanyam and S Naranan. Quantitative linguistics and complex system studies\*. *Journal of Quantitative Linguistics*, 3(3):177–228, 1996.
- [3] Aaron Clauset, Cosma Rohilla Shalizi, and Mark EJ Newman. Power-law distributions in empirical data. *SIAM review*, 51(4):661–703, 2009.
- [4] J. P. Crutchfield. Information and its Metric. In L. Lam and H. C. Morris, editors, *Nonlinear Structures in Physical Systems – Pattern Formation, Chaos and Waves*, pages 119–130. Springer Verlag, 1990.
- [5] Ronald Dickman, Nicholas R. Moloney, and Eduardo G. Altmann. Analysis of an information-theoretic model for communication. *Journal of Statistical Mechanics: Theory and Experiment*, 2012(12):P12022, 2012.
- [6] Ramon Ferrer-i Cancho and Brita Elvevåg. Random texts do not exhibit the real Zipf’s law-like rank distribution. *PLoS One*, 5(3):e9411, 2010.
- [7] Johann Gottfried Herder. Abhandlung über den Ursprung der Sprache. *Johann Gottfried Herder Frühe Schriften*, pages 695–810, 1772.
- [8] Ramon Ferrer i Cancho. The variation of Zipf’s law in human language. *The European Physical Journal B-Condensed Matter and Complex Systems*, 44(2):249–257, 2005.
- [9] Ramon Ferrer i Cancho. Zipf’s law from a communicative phase transition. *The European Physical Journal B-Condensed Matter and Complex Systems*, 47(3):449–457, 2005.
- [10] Ramon Ferrer i Cancho and Albert Díaz-Guilera. The global minima of the communicative energy of natural communication systems. *Journal of Statistical Mechanics: Theory and Experiment*, 2007(06):P06009, 2007.
- [11] Ramon Ferrer i Cancho and Ricard V. Solé. Least effort and the origins of scaling in human language. *Proceedings of the National Academy of Sciences*, 100(3):788–791, 2003.

- [12] Wentian Li. Random texts exhibit Zipf’s-law-like word frequency distribution. *Information Theory, IEEE Transactions on*, 38(6):1842–1845, 1992.
- [13] George A Miller. Some effects of intermittent silence. *The American Journal of Psychology*, 70(2):311–314, 1957.
- [14] Martin A. Nowak and David C. Krakauer. The evolution of language. *Proceedings of the National Academy of Sciences*, 96(14):8028–8033, 1999.
- [15] Mikhail Prokopenko, Nihat Ay, Oliver Obst, and Daniel Polani. Phase transitions in least-effort communications. *Journal of Statistical Mechanics: Theory and Experiment*, 2010(11):P11025, 2010.
- [16] Mikhail Prokopenko, Daniel Polani, and Matthew Chadwick. Stigmergic gene transfer and emergence of universal coding. *HFSP Journal*, 3:317–327, 2009.
- [17] Vladimir A. Rokhlin. Lectures on the entropy theory of measure-preserving transformations. *Russian Mathematical Surveys*, 22(5):1–52, 1967.
- [18] Ryuji Suzuki, John R Buck, and Peter L Tyack. The use of Zipf’s law in animal communication analysis. *Animal Behaviour*, 69(1):F9–F17, 2005.
- [19] Matt Visser. Zipf’s law, power laws and maximum entropy. *New Journal of Physics*, 15(4):043021, 2013.
- [20] George Kingsley Zipf. Human behavior and the principle of least effort. 1949.