

# LANGUAGE LEARNING, POWER LAWS, AND SEXUAL SELECTION

TED BRISCOE

*Computer Laboratory  
University of Cambridge  
JJ Thomson Ave  
Cambridge CB3 0FD, UK  
ejb@cl.cam.ac.uk*

I discuss the ubiquity of power law distributions in language organisation (and elsewhere), and argue against Miller's (2000) argument that large vocabulary size is a consequence of sexual selection. Instead I argue that power law distributions are evidence that languages are best modelled as dynamical systems but raise some issues for models of iterated language learning.

## 1. Introduction

A diagnostic of a power law distribution is that a log-log plot of frequency against rank yields a (nearly) straight line. For instance, Zipf (1935) plotted word token counts in a variety of texts against the inverse rank of each distinct word type and showed that typically such plots approximate a straight line. The characteristic 'Zipf curve' of word frequency against rank deviates from this line because the relative frequency of very common word types, such as the English determiners *the* and *a*, tend to be more similar than the power law predicts, as also does the relative frequency of very rare words in the tail of the distribution. Zipf's 'law' is often expressed as:

$$c(w) \propto \frac{1}{r(w)^B} \quad (1)$$

where  $B > 1$ , the exponent, defines the slope of the plot, frequency  $c(w)$  is the token count of word type  $w$  in text, and rank  $r(w)$  is the position of word type  $w$  in the list of word types sorted in descending order of frequency,  $c(w)$ . Guiraud's (1954) related law states that the number of word types  $V$  in a text is proportional to the length of that text  $N$ :

$$V \propto N^A \quad (2)$$

Although the models of power law distributions of which I am aware have a dynamical component, they have received little attention from evolutionary linguists.

I know of only one argument, due to an evolutionary psychologist (Miller, 2000), which utilises Zipf's observation about word frequencies in attempting to explain large, redundant vocabularies in terms of sexual selection. I argue against this explanation in §4, but, before doing so, I discuss the ubiquity of power law distributions in §2, some relevant models of them in §3, return to Miller's argument in §4, and then discuss some issues power law distributions raise for evolutionary models of iterated language learning in §5.

## **2. Manifestations of Power Law Distributions**

Power law distributions are very different from normally distributed phenomena, such as height, which yield the characteristic 'Bell Curve'. The factors that influence a person's height, such as nutrition and genetic inheritance, combine in a more linear manner so that (relatively minor) variation in height is normally distributed around a mean that can be accurately estimated from a representative sample of the population. Zipf (1949) noted that Pareto's observations about the distribution of wealth in the population could also be modelled using a version of his 'law'. What makes wealth different from height intuitively is that the factors that influence the amount of money we have combine non-linearly and there are strong (positive) feedback effects (i.e. 'the rich get richer'). We now know that power law distributions are good approximations of many other non-linguistic phenomena, such as the distribution of people within cities, citations amongst scientists, accesses of web pages, species within habitats, authors amongst scientific articles, actors within films, links between web pages, activation of genes, size of earthquakes, number of sexual partners, and many more (e.g. Albert & Barabasi, 2002).

There are similar results regarding extrinsic properties of languages: for instance, the distribution of languages within language families approximates a power law (Wichmann, 2005). In terms of inherent properties of language, Zipf also showed that plotting the length or the number of meanings of word types against their frequency also yields similar distributions. With the increasing availability of annotated electronic corpora, these observations have been extended to many other areas of language organisation, such as the frequency of contiguous sequences of words (bigrams, trigrams, and more generally ngrams), of grammatical rules, of construction types, of lexical relations between word types, as well as the length of constituents, and the association of verbs with constructions (e.g. Sharman, 1989; Manning & Schutze, 1999; Korhonen 2002; Yook *et al.*, 2001).

## **3. Models of Power Law Distributions**

So far I have used the term 'distribution' ambiguously between the linguistic and probabilistic sense. The most important insight about such distributions with large numbers of rare events (e.g. Baayen, 2001) is that converting a frequency-rank plot into a probability-rank plot via maximum likelihood (i.e. relative frequency)

estimation, and treating the result as a probability distribution is unwise. Since the counts of the tail are very low, statistical estimation theory tells us that they will be unreliable. A rare word, for instance, may suddenly become fashionable (e.g. *egregious*, *serendipity*) and thus increase in relative frequency over a given time period. Since, we always see a long tail of rare events no matter how much (more) text we sample, and the number of types grows in proportion to the size of this sample (Guiraud's law), power law distributions are often described as 'scale-free'. In statistical terms, power law distributions which remain invariant over different sample sizes are a strong indication that we may be sampling from a statistically unrepresentative non-stationary (i.e. dynamical) system.

Baayen (1991), following in the tradition of Mandelbrot (1953) and Simon (1955), develops a stochastic Markovian model of phonotactically legal Dutch word strings and relates it to empirical data on similarities between words by phonological form and by relative frequency. He finds that to model these effects accurately, it is necessary to add a second 'dynamical' stochastic model which introduces or removes word types with probability proportional to their token frequency. This has the effect of increasing overall frequency-based and decreasing form-based similarity. For present purposes, it is indicative that the second dynamical word 'birth-death' process is required even though it says nothing directly about the relationships between word types.

Albert & Barabasi (2002) provide a recent survey of work on 'small world' networks in which most nodes of a network can be reached by any other in a small number of (node) steps, though the overall number of nodes can be arbitrarily high. They define a dynamical algorithm for generating such networks, by continuously adding new nodes and attaching them to old nodes with probability proportional to their number of existing links. They prove that such networks evolve to a scale-free organisation obeying a power law distribution in which there is a long tail of nodes with low numbers of links and a small number of 'popular' nodes with many links. They also prove that both 'growth', the dynamical component, and 'preferential attachment' are necessary for this pattern to emerge. Such networks have been applied to models like that of Baayen (1991), described above (e.g. Bornholdt and Ebel, 2001), and to lexical semantic organisation (e.g. Yook *et al.*, 2001).

#### **4. Power Laws and Sexual Selection**

Miller (2000:369f), in the context of a more general argument that human language evolved by sexual selection, argues that large vocabulary size, in comparison with those of other (artificial and natural) animal communication systems, evolved through sexual selection. Women preferred men with large active vocabularies but needed to acquire large passive vocabularies themselves to assess the trait. Miller offers, as evidence for the non-functional nature of much of this vocabulary, Zipf's observation that vocabulary distributes like a power law and

contains many near synonyms:

...any of the words we know is likely to be used on average about once in every million words we speak... Why do we bother to learn so many rare words that have practically the same meanings as common words, if language evolved to be practical? (Miller, 2000:370)

He argues that human variation in vocabulary acquisition correlates with intelligence and has a heritable component, and thus is an (indirect) fitness indicator, triggering an 'arms race' in which advertising excessive vocabulary size is a 'display' of fitness akin to the peacock's tail, precisely because it does not contribute usefully to communication.

In §2 we saw that power law distributions manifest themselves in many areas of linguistic organisation. For instance, there is a tail of rare long constituents in text samples (Sharman, 1989). However, there is no evidence that 'display' of such forms is a particular feature of courtship, nor that such forms are non-functional. As we saw in §3 models predicting such distributions need only a dynamical component and no element of natural or sexual selection whatever. Evidence of power law distributions in both idiolects and language forces us to conclude that both are best modelled as dynamical systems – rather than well-formed sets, as in generative linguistics (e.g. Sampson, 2001:165f) – but nothing more.

If vocabulary size were non-functional, we might expect there to be many truly synonymous words. What we find in the organisation of vocabulary is that partially synonymous words have different distributions in terms of specificity of reference, syntactic potential, or genre and register. There is, in fact, considerable evidence that children avoid hypothesising synonyms in language acquisition (e.g. Clark, 2003) and that language users adhere to the convention of preemption by synonymy, except where discourse or syntactic context triggers a non-synonymous reading (e.g. Briscoe *et al.*, 1995; Copestake & Briscoe, 1995). For instance, *cow*, unlike *chicken*, is not generally used to refer to the meat because of the existence of *beef*. However, in an appropriate context *cow* can be used this way and triggers an implicature of 'disgust':

There were five thousand extremely loud people on the floor eager to tear into roast cow with both hands and wash it down with bourbon whiskey. (Tom Wolfe, 1979. *The Right Stuff*, Farrar, Straus and Giroux, New York (p. 298, Picador edition, 1991))

Similarly, the word *stealer*, formed by the fairly productive derivational rule of agentive *+er* nominalisation, is blocked by *thief*, except in syntactic contexts where the specificity of reference is narrowed:

He is an inveterate \*stealer / thief / stealer of Porsche 911s

These and many similar observations suggest that partial synonymy is communicatively useful and actively exploited to convey meaning.

To understand why we have so many words and how the cognitive ability to cope with them (co-)evolved, consider the likely environment of adaptation for language. In a foraging, scavenging or hunter society, the ability to discriminate – and thus name more and more species, according to nutritional value, location, method of capture or harvesting, and so forth – would be of value for survival because it would allow efficient transmission of these skills to kin as well as survival over larger and more varied habitats. Modern hunter-gatherers are known to have large vocabularies specialised in this way (Diamond, 1997). This may not have been the sole driver for increasing vocabulary size, but it has the advantage that it predicts that vocabulary will be to a large extent organised by specificity of reference. It is useful not only to be able to talk about plants in general but also species and subgroups (e.g. by location or edible part) in order to discriminate the edible, find the source, and harvest effectively. Once we accept such a pressure to name in an increasingly complex and multifaceted environment, then the tendency for there to be smaller numbers of high frequency words of generic reference and a larger number of rarer words with highly specific denotations is just a case of the structure of vocabulary mirroring (our perception of) this environment.

## 5. The Real Challenge – Iterated Learning

One achievement of recent evolutionary models of language is the demonstration that treating languages as complex adaptive systems responding to conflicting selection pressures (e.g. Briscoe, 2000) leads to insightful accounts of typological and other linguistic universals without the need to invoke innateness. These accounts rely heavily on the iterated learning model (ILM, e.g. Kirby, 2001) in which linguistic traits must undergo repeated relearning by successive generations of language learners acquiring their language from that of the previous generation. For instance, Kirby (2001) demonstrates that languages in the ILM evolve to have compositional structure in which only high frequency irregular form-meaning mappings are stable, given the following assumptions:

1. an *invention strategy* for form-meaning pairs,
2. a *production bias* to express meanings using short forms,
3. an *inductive bias* to learn small grammars and lexicons,
4. a *learning period* in which not all form-meaning pairs appear
5. and *environmental structure* which favours some meanings

In the simulation, initial (proto)languages are holistic and non-compositional but chance regularities which emerge in form-meaning mappings are acquired by

learners, who then reliably exemplify them for the next generation of learners, because regularities are, by definition, more frequent in data. Thus, over time the language evolves to be mostly compositional and regular. However, (short) irregular mappings can survive provided they are associated with meanings which are expressed frequently and, therefore, also occur reliably during the learning period.

This instantiation of the ILM neatly explains the observation that irregularity correlates with high frequency in attested languages: children would continue to say *goed* into adulthood if *went* were not a high frequency form. The corollary, however, is that rare unpredictable properties of language which do not follow from some regularity manifest during the learning period should be unstable and, therefore, rarely observed.

Rare word-meaning associations are unpredictable and may also influence lexico-grammatical behaviour. For example, the verb *obsess* is a stable lexeme of English, but does not appear in any of the 40 or so case studies of child-directed speech in CHILDES<sup>a</sup>. It is transitive but usually appears in the passive in adult speech accompanied by a PP headed by *by*, *with* or *over*. However, vocabulary acquisition continues through adulthood, so the ILM (and other models) simply predict that such vocabulary will be acquired later (and less universally).

Marked but predictable constructions, such as multiple centre-embeddings, which Sampson (2001:21) estimates occur once in every 250K words on average, are also not counter-examples if one believes that they are a consequence of learners acquiring, in the basis of more frequent constructions, grammatical rules which correctly predict the appropriate form-meaning mapping for these constructions.

A more challenging case for the ILM is diathesis alternation, in which verbs of certain semantic classes semi-predictably occur in alternant constructions often with predictable meaning changes. For instance, *eat* can appear in intransitive and transitive constructions but when it occurs intransitively the theme of the action is 'understood'. However, verbs with similar senses, such as *devour* or *consume* do not undergo this alternation. There is evidence that children learn at least some of these alternation rules by 3 years because they produce errors, such as *Don't fall my dolly down* – the causative-inchoative alternation. However, the rate at which such errors occur also suggest that alternation rules are learnt conservatively and only rarely overapplied. There are on the order of 100 such alternation rule types in English, when productive meaning change is taken into account.

Figure 1 shows log-log plots of the unconditional probability of over 150 verb-headed constructions against their inverse rank on the left and of the conditional probability of these constructions when headed by any form of the verb *believe* on the right, calculated from 30M words of automatically parsed text along with the closest fit straight line derived using (1) above with  $B$  set appropriately. Both distributions loosely approximate power laws with long tails of rare events, but

---

<sup>a</sup><http://childes.psy.cmu.edu/data/>

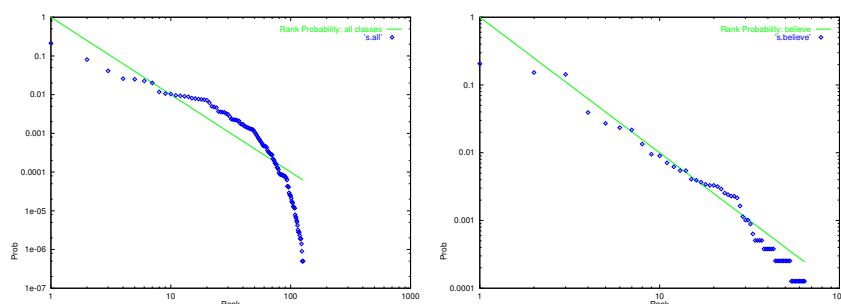


Figure 1. AllVerb/Believe Constructions

critically the correlation between the unconditional distribution and conditional ones for individual verbs is low (0.47 Spearman-Rank Coefficient for 14 verbs). This means that it is not possible to predict the individual association of verbs and constructions on the basis of the unconditional distribution. For instance, sentential complements are rare overall but the most common construction with *believe*, accounting for over 90% of occurrences. This lack of correlation, taken together with the fact that analysis of CHILDES shows that child-directed speech only exemplifies common verb-construction associations (e.g. Buttery & Korhonen, 2005), suggests that children do not have reliable evidence for the existence of most alternation *rules* – assuming that evidence would be several exemplars of the same alternation involving several different verbs.

It may be that such semi-productive alternations are also acquired later in life, despite the occasional errors in children’s speech. This is a general strategy that proponents of ILM-style explanations can take. But on the other hand, there must also be some learning ‘bottleneck’, caused by limited exposure to data during the learning period, for ILM accounts of linguistic evolution to work. Cases like this pose interesting challenges for the approach because they suggest that linguistic data is distributed in such a fashion that there may still be a ‘poverty of stimulus’ issue during the sensitive period for acquisition. More empirical work on language acquisition is needed to determine whether the ILM’s predictions hold up for such specific cases.

## References

- Albert, R. & A. Barabasi (2002) ‘Statistical mechanics of complex networks’, *Reviews of Modern Physics*, vol.74, 47–97.
- Baayen, H. (1991) ‘A stochastic process for word frequency distributions’, *Proceedings of the Assoc. for Computational Linguistics*, Morgan Kaufmann, Menlo Park, CA, pp. 271–278.
- Baayen, H. (2001) *Word Frequency Distributions*, Kluwer, Dordrecht.

- Bornholdt, S. & Ebel, H. (2001) 'World Wide Web scaling exponent from Simon's 1955 model', *Physical Review*, vol.64, 035104.
- Briscoe, E.J. (2000) 'Evolutionary perspectives on diachronic syntax' in (eds) Pintzuk, S., Tsoulas, G. and Warner, A. (eds.), *Diachronic Syntax: Models and Mechanisms*, Oxford University Press, Oxford, pp. 75–108.
- Copestake, A.A. and E.J. Briscoe (1995) 'Regular polysemy and semi-productive sense extension', *Journal of Semantics*, vol.12, 15–67.
- Briscoe, E.J., A.A. Copestake and A. Lascarides (1995) 'Blocking' in St. Dizier, P. and Viegas, E. (eds.), *Computational Lexical Semantics*, Cambridge University Press, Cambridge, pp. 273–302.
- Buttery, P. & A. Korhonen (2005) 'Large-scale analysis of verb subcategorization differences between child directed speech and adult speech', *Proceedings of the Interdisciplinary Workshop on the Identification and Representation of Verb Features and Verb Classes*, Saarland University.
- Clark, E. (2003) *First Language Acquisition*, Cambridge University Press, Cambridge.
- Diamond, J. (1997) *Guns, Germs and Steel: The Fate of Human Societies*, Random House, New York.
- Guiraud, H. (1954) *Les Caractères Statistiques du Vocabulaire*, Press Universitaires de France, Paris.
- Kirby, S. (2001) 'Spontaneous evolution of linguistic structure: an iterated learning model of the emergence of regularity and irregularity', *IEEE Transactions on Evolutionary Computation*, vol.5(2), 102–110.
- Korhonen, A. (Computer Laboratory, University of Cambridge) *Subcategorization Acquisition*, Technical Report UCAM-CL-TR-530. 2002
- Mandelbrot, B. (1953) 'An informational theory of the statistical structure of language' in W. Jackson (eds.), *Communication Theory*, Butterworths, London.
- Manning, C. & H. Schütze (1999) *Foundations of Statistical Natural Language Processing*, MIT Press, Cambridge MA.
- Miller, G. (2000) *The Mating Mind: How Sexual Choice Shaped the Evolution of Human Nature*, William Heinemann, London.
- Sampson, G. (2001) *Empirical Linguistics*, Continuum, London.
- Sharman, R. (1989) *Observational Evidence for a Statistical Model of Language*, IBM UKSC Report 205.
- Simon, H. (1955) 'On a class of skew distribution functions', *Biometrika*, vol.42, 435–440.
- Wichmann, S. (2005) 'On the power law distribution of language family sizes', *Journal of Linguistics*, vol.41, 117–131.
- Yook S., Jeong, H., Barabasi, A-L & Tu, Y. (2001) 'Weighted evolving networks', *Physical Review Letters*, vol.86, 5835–5838.
- Zipf, G. (1935) *The Psycho-Biology of Language: An Introduction to Dynamic Philology*, Houghton-Mifflin, New York.
- Zipf, G. (1949) *Human Behavior and the Principle of Least Effort*, Addison-Wesley, Cambridge, MA.