

*"This book is an introduction to the statistical analysis of word frequency distributions, intended for linguists, psycholinguists, and researchers working in the field of quantitative stylistics and anyone interested in quantitative aspects of lexical structure.*

The following online article has been derived mechanically from an MS produced on the way towards conventional print publication. Many details are likely to deviate from the print version; figures and footnotes may even be missing altogether, and where negotiation with journal editors has led to improvements in the published wording, these will not be reflected in this online version. Shortage of time makes it impossible for me to offer a more careful rendering. I hope that placing this imperfect version online may be useful to some readers, but they should note that the print version is definitive. I shall not let myself be held to the precise wording of an online version, where this differs from the print version. Published in Computational Linguistics Word Frequency Distributions R. Harald Baayen University of Nijmegen Dordrecht: Reviewed by Geoffrey Sampson University of Sussex This book must surely in future become the standard point of departure for statistical studies of vocabulary. Baayen begins with a puzzle that has troubled many investigators who have studied vocabulary richness, for instance people hoping to find stylistic constants characteristic of individual authors for use in literary or forensic authorship disputes. It is not, because it is not independent of sample size. In most domains, sample means fluctuate randomly around population means while getting closer to them as sample sizes increase. In natural-language vocabulary studies, mean word frequencies systematically increase with sample size even when samples of tens of millions of words are examined. One might hypothesize that greater experience would lead a writer to use a richer vocabulary in a later book, but mean word frequency is actually higher. However, *Through the Looking-Glass* is a somewhat longer book. If just the first words are used this is the length of the earlier book, the direction of the difference in mean word frequencies is reversed: Normally, more data give a more accurate picture of anything; but here the direction of change in frequency, from 9. Can we conclude that Carroll was using a richer vocabulary in the later book, because of the figures for equal-sized samples; or that he was using a less rich vocabulary, because of the figures for total available samples; or can we make no inference either way? Gustav Herdan argued in a series of works which were influential in the 1950s that the ratio of the logarithms of number of types and number of tokens was such a constant. Conversely, Baayen quotes Naranan and Balasubrahmanyam. Eventually, Baayen is able to show that this negative position is also unjustified; but between that conclusion, and the statement of the puzzle, lie some two hundred pages of fairly dense mathematics. This is certainly not a book for the mathematically faint-hearted. Baayen does a great deal, though, to help the reader follow him through the thickets. Not only does each chapter end with a summary of its findings, but – unusually for a work that is not a student textbook – Baayen also gives lists of test questions which the diligent reader can work through to consolidate his understanding of the material. What lies behind the unusual relationship between type frequencies and sample sizes in the case of vocabulary? Baayen clarifies the situation by an analogy with dice-throwing. Baayen plots a graph showing how the expected frequency spectrum that is, the number of vocabulary elements observed once, the number of vocabulary elements observed twice, etc. changes as the sequence is extended. For hapax legomena elements of the vocabulary observed once each, the expected figure rises to a maximum of about 2. For successive elements of the spectrum, the waves are successively lower and later, but the pattern is similar: Meanwhile, a plot on the same graph of expected sample vocabulary size rises rapidly and is close to the population vocabulary size. In most domains to which statistical techniques are applied, sample sizes are large enough to involve areas far out to the right of this kind of graph a serious examination of possible bias in a dice would surely involve hundreds of throws, so the special features of its left-hand end are irrelevant. With natural-language vocabulary studies, on the other hand, even the largest practical samples leave us in an area analogous to the extreme left-hand end of the dice-throwing graph, with numbers of hapax legomena and consequently also dis legomena, tris legomena, etc. The intuitive meaning of this is fairly clear, and it is made exact via alternative formal definitions. Using

these techniques, it turns out that the growth in vocabulary richness between *Alice in Wonderland*, and *Through the Looking-Glass* after truncation to make its length the same, is marginally significant. Not all of the exposition is original with Baayen. One of the many virtues of his book lies in drawing together in one convenient location a clear statement of relevant analyses by others over several decades, often published relatively obscurely. Carroll, H. Sichel, and J. A. A point which emerges from the book and which readers of this review may have begun to infer from names cited is the extent to which, in the late 20th century, this mathematical approach to natural language was a scholarly speciality of the former Soviet Union; in consequence it was largely unknown in the West. There are other channels through which this work has become accessible to the English-speaking world in recent years, notably the *Journal of Quantitative Linguistics* but that German-based journal, though published in English, has to date attracted limited attention in Britain and North America. The book under review may well be the most significant route by which important ex-Soviet research in our area will become known to English-speaking scholars. It would be beyond the scope of this review to survey all issues relating to LNRE distributions which Baayen investigates. For linguists, one particularly interesting area concerns departures from the randomness assumption made by the simpler LNRE models. These pretend, for the sake of mathematical convenience, that texts are constructed by drawing successive words blindly out of an urn containing different numbers of tokens of all possible words in the vocabulary, so that the difficulties to be addressed relate only to the vast size of the urn. Real life is not like that, of course: If we are primarily interested in overall vocabulary size, one problem that is repeatedly produced by the urn model is that inferences from vocabulary size in observed samples to vocabulary sizes for other, so-far unobserved sample sizes turn out to be overestimates when samples of the relevant size are examined. Many linguists, particularly after the above discussion of the urn model, will be professionally inclined to assume that this problem stems from ignoring syntactic constraints within sentences, as the urn model does. Baayen demonstrates that this is not the source of the problem. If the sentences of *Alice in Wonderland* are permuted into a random order while preserving the sequence of words within each individual sentence, the overestimation bias disappears. Different passages of a document deal with different topics, so topic-sensitive words are not distributed evenly through the text. No doubt what Baayen gives us is not always the last word to be said on some of the questions he takes up, but as already suggested it is hard to think that future analyses will not treat Baayen as the standard jumping-off point for further exploration. But, in the first place, this seems disconnected from the body of the book, because the relevant distributions are not LNRE. Furthermore, the only newspaper identified by name is the *Frankfurter Allgemeine Zeitung*, and although we are told that other papers show the same pattern, we are not told which papers these are. Finding that Germans perceive a unique historical discontinuity in the *s* might be a very different thing from finding that Europeans, or Westerners, in general do so. Nevertheless, this last chapter does also contain important findings that relate more closely to the central concerns of the book. But there are exceptions: The last chapter also contains a number of misprints, which are not self-correcting and may be worth listing here. In a discussion of word-length distribution, there are repeated confusions between length 4, length 5, and length 6, on p. The volume is accompanied by a CD-ROM containing numerous relevant software programs; these and various data sets are detailed in a series of four Appendices to the book. Acknowledgment I am grateful to fellow members of the Sussex University probabilistic parsing reading group for insights gained during the weeks when the book under review was our target volume. Responsibility for errors and inadequacies in this review is entirely mine. On sampling from a lognormal model of word frequency distribution. Brown University Press, Providence, R. Some simple measures of richness of vocabulary. *Association of Literary and Linguistic Computing Bulletin*, pp. Centre for Mathematics and Computer Science. Models for power law relations in linguistics and information science. *Journal of Quantitative Linguistics*, 5: On a distribution law for word frequencies. *Journal of the American Statistical Association*, *Bulletin of the Academy of Sciences*, Georgia,

## Chapter 2 : Bag of Words and Frequency Distributions in C# | Nick Grattan's Data Science Blog

*This book is an introduction to the statistical analysis of word frequency distributions, intended for linguists, psycholinguistics, and researchers working in the field of quantitative stylistics and anyone interested in quantitative aspects of lexical structure.*

The aim of this book is to make these techniques accessible to nonspecialists. Harald Baayen prepares readers to conduct practical analyses of word frequencies in text samples. Good instruction with just the right amount of literary whimsy. The central graphical and statistical representation of word frequency analysis is the "frequency spectrum" defined--adequately for even beginners--on page 8. Frequency spectra are lists of all words in a text sample along with how many times each word appears in the sample. These lists are sorted from most to least frequent words. High frequency words are usually function words such as "the" and "a. Word frequency spectra also contain numerous very low frequency words, many appearing only once. The tail of low frequency words stretches off to seeming infinity. From a sampling perspective, this long tail is made theoretically longer by the virtual existence of "zero" frequency words which would achieve a frequency of "one" if the text sample size were sufficiently increased. The first chapter further defines frequency spectra and presents other concepts basic to word frequency analysis. The long tail property underlies the first problem in frequency spectra analysis. The number of unique words in a text sample systematically increases with the word length of the sample. This introduces bias into statistical procedures which compare different-sized text samples. This bias distorts analysis at the full range of document lengths, from brief open-ended survey responses to book-length narratives. Baayen demonstrates these biasing effects in commonly used frequency comparison statistics. He then introduces corrective procedures using such information as the "vocabulary growth rate" observed as a text sample size is incrementally increased in length. The well-organized narrative of describing, demonstrating and solving this problem plays out across chapters two, three and four. A second problem stems from the statistical assumption that words are randomly distributed within text samples. Obviously this is never true in naturally-occurring text. Less obvious is the type of bias this assumption introduces and how to correct for it. Using clever experimentation, Baayen traces the primary biasing effects to our mid-frequency key content words, which tend to cluster in all text, more so in longer documents. He concludes chapter five with recommendations for adjusted procedures that reduce the biasing effects of these "underdispersed" words. The reader is now prepared to conduct unbiased analysis of word frequency distributions. The software on the included CD implements the improved frequency analyses devised by the author. It requires Linux, but his web site at the University of Alberta contains a downloadable Windows version. More recent research on the properties of word frequency distributions is available from this web site and is worth reading. I personally found this book quite valuable, even though I do not conduct any of the specific analyses demonstrated. Much of my word frequency work is semi-automated content analysis of open-ended survey responses primarily using QDA Miner and WordStat from Provalis Research. He achieves his stated goal of making his corrected statistical procedures understandable and accessible to nonlinguists.

## Chapter 3 : How to Make a Frequency Table on Microsoft Excel | [www.nxgvision.com](http://www.nxgvision.com)

*Zipf's law in fact refers more generally to frequency distributions of "rank data," in which the relative frequency of the  $n$ th-ranked item is given by the Zeta distribution,  $1/(n^s \zeta(s))$ , where the parameter  $s > 1$  indexes the members of this family of probability distributions.*

## Chapter 4 : chi squared - Comparing two word distributions - Cross Validated

*This book is an introduction to the statistical analysis of word frequency distributions, intended for linguists, psycholinguistics, and researchers working in the field of quantitative stylistics and anyone interested in quantitative*

