

Finite-Size Scaling Law for Word-Frequency Distributions

Alvaro Corral^{1,2,3,4} and Francesc Font-Clos⁵

¹ Centre de Recerca Matemàtica, Barcelona, Spain

² Departament de Matemàtiques, Universitat Autònoma de Barcelona, Spain

³ Barcelona Graduate School of Mathematics, Barcelona, Spain

⁴ Complexity Science Hub Vienna, Austria

⁵ ISI Foundation, Torino, Italy

The dependence on text length of the statistical properties of word occurrences has long been considered a severe limitation for the usefulness of quantitative linguistics. We propose a simple scaling law for the distribution of word frequencies that brings to light the robustness of this distribution as text grows. In this way, the shape of the distribution is always the same, and it is only a scale parameter that increases (linearly) with text length. We give evidence for the validity of such scaling law, both using analytical arguments based on the generalized central-limit theorem and careful statistical tests.