## UCSF UC San Francisco Previously Published Works

## Title

The Birth-Death-Mutation Process: A New Paradigm for Fat Tailed Distributions

## Permalink

https://escholarship.org/uc/item/1536z7vn

## Journal

PLOS ONE, 6(11)

## ISSN

1932-6203

## **Authors**

Maruvka, Yosef E Kessler, David A Shnerb, Nadav M

## **Publication Date**

2011

## DOI

10.1371/journal.pone.0026480

## **Copyright Information**

This work is made available under the terms of a Creative Commons Attribution License, available at <u>https://creativecommons.org/licenses/by/4.0/</u>

Peer reviewed

# The Birth-Death-Mutation Process: A New Paradigm for Fat Tailed Distributions

#### Yosef E. Maruvka, David A. Kessler, Nadav M. Shnerb\*

Department of Physics, Bar Ilan University, Ramat-Gan, Israel

#### Abstract

Fat tailed statistics and power-laws are ubiquitous in many complex systems. Usually the appearance of of a few anomalously successful individuals (bio-species, investors, websites) is interpreted as reflecting some inherent "quality" (fitness, talent, giftedness) as in Darwin's theory of natural selection. Here we adopt the opposite, "neutral", outlook, suggesting that the main factor explaining success is merely luck. The statistics emerging from the neutral birth-death-mutation (BDM) process is shown to fit marvelously many empirical distributions. While previous neutral theories have focused on the power-law tail, our theory economically and accurately explains the entire distribution. We thus suggest the BDM distribution as a standard neutral model: effects of fitness and selection are to be identified by substantial deviations from it.

Citation: Maruvka YE, Kessler DA, Shnerb NM (2011) The Birth-Death-Mutation Process: A New Paradigm for Fat Tailed Distributions. PLoS ONE 6(11): e26480. doi:10.1371/journal.pone.0026480

Editor: Eshel Ben-Jacob, Tel Aviv University, Israel

Received September 7, 2011; Accepted September 27, 2011; Published November 1, 2011

**Copyright:** © 2011 Maruvka et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Funding:** The authors were funded by an EU grant: the complexity pathfinder of NEST, by the Israeli Ministry of science TASHTIOT program, and by the Israeli Science Foundation BIKURA grant no. 1026/11. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

\* E-mail: nadav.shnerb@gmail.com

#### Introduction

Survival of the fittest or of the luckiest? The answer depends on the subject considered. Out of ten pairs of pants bought a year ago, the survivors are perhaps those made of a better material; if wineglasses are considered, persistence is mainly a matter of luck. In the absence of prior knowledge, statistics must be used in order to identify the role of fortune: wineglass life expectancy, for example, is described by an exponential distribution. Strong deviations from this statistics indicate to what extent "death" is a result of accumulated wear, rather than from uncorrelated random events.

In many complex systems, though, it is hard to identify relative role of fortune. Large differences in success (of investors or authors) or abundance (of bio-species) do not necessarily reflect the "quality" or the "fitness" of the rich and the frequent. Huge abundance fluctuations may be a result of accumulation of stochastic events, as exemplified by the uneven statistics of surnames in society [1].

The schism between the "neutral" (stochastic) and the "fitness" (deterministic) outlooks is most strongly pronounced in the theory of evolutionary dynamics [2]. Darwin condemned those who "attribute ... (species") proportional numbers to what we call chance. But how false a view is this! [3]" and held that the main factor shaping eco-communities is natural selection. The opposite view, that random drift plays the major role in evolution — both on the molecular (Kimura's neutral evolution [4]) and the ecological (Hubbell's community drift model [5]) levels — has sparked a series of ongoing hot and emotional debates.

In economy and social sciences the deterministic approaches tend to emphasize the tremendous inequality in income and wealth, say, as reflecting underlying "quality" (from prudence to crookedness) differences. The opposing neutral approach [6] have recently found a prominent outspoken, Nassim Taleb. In his books [7,8] he maintains that the weight of unpredictable events (what he calls "black swans") is overwhelming in determining economic and social success.

Purely deterministic and purely stochastic theories are both oversimplifications. The real scientific problem is to find the relative weight of chance versus fitness. The assumption of neutral dynamics is most useful as a null hypothesis, with which empirical statistics should be compared. Nowadays this role is played by the Yule-Simon statistics [9–11], or its approximation by a simple power law [12,13]. In the following we briefly review Yule's model and point out its major shortcoming. We suggest a correction that yields different statistics and show that the new distribution fits many "canonical" empirical datasets very nicely.

Yule-Simon theory [9] arose from a study of the the highly skewed distribution of biological species within genera. One of the graphs studied by Yule — for the family of long-horn beatles Cerambycinea — is plotted in the left inset of Fig. 1. This is a Pareto plot showing  $n_m$ , the fraction of genera with m species, vs. m on a log-log scale. One observes a few "wealthy" genera to which many species belong, and many "poor" genera with apparent linear dependence that suggests a power-law distribution.

Yule's neutral model posited that the rate of speciation is fixed for all species. Upon speciation, the new species stays in the same genus with probability  $1-\mu$ .  $\mu$ , the chance that the offspring species branches out to form a new genus, is also fixed, ensuring perfect neutrality (no fitness). This simple process generates a steady state distribution that converges rapidly to a power law for the relative species abundance  $n_m$ ,

$$n_m = CB(m, 2+\mu) \sim Cm^{-(2+\mu)}.$$
 (1)



**Figure 1. Species within genera statistics for Cerambycinea beatles.** The original species within genera statistics used by Yule (blue squares), based on 1024 genera known at 1925 for the Cerambycinea family (down left). On a log-log scale this graph looks very much straight, suggesting a power-law statistics (black line). In the main figure, the black circles show the contemporary statistics as obtained for 4411 genera (27221 species of Cerambycinea [24]), where a pronounced "shoulder" appears. The red line is the best fit of BDM function (2),  $\gamma$  is the diversification rate and  $\mu$  reflects the chance of a new species to initiate a new genus. The blue line shows the prediction of our theory for a sample of  $R_0 = 5719$  species chosen at random out of the 22271 known today with the same  $\mu$  and  $\gamma$ , as obtained from Eq. 3. This is now a prediction *without any fitting parameters*, to be compared with the original Yule statistics.

where C is a normalization factor. Note that this fat-tailed distribution has nothing to do with the "quality" differences among species, instead it is a result of the multiplicative character of the noise.

As pointed out by Herbert Simon [10], Yule's argument goes far beyond its original context. Simon considered power-laws for the number of occurrences of words in a text, scientific publications and wealth distribution. Subsequently, the appearance of powerlaws has been recognized as a fundamental feature of eco-, econo-, bio- and socio-systems, with countless of examples from protein family statistics [14], surname abundance ratio [1,15], internet connections [16], firm sizes [17], casualties in terror attacks [18] and so on. In addition the common scenario considered in the new popular theory of scale free networks - the preferential attachment dynamics - is indeed mathematically equivalent to Yule's process [see methods (A)] where small families are generated by a source, not by mutations [11].

#### **Results and Discussion**

As a starting point for the presentation of our new neutral model, let us stick for the moment to the original context of Yule theory, the species within genera statistics. The main panel of Figure 1 reveals a major failure of the Yule-Simon model. The original distribution observed by Yule for Cerambycinea beatles, based on the 1024 genera (5719 species) is compared with the current data with 27221 species and 4411 genera. Clearly, something bad has happened to the simple power-law: it characterizes now only the tail of the distribution, and a very pronounced "shoulder" appears for the small genera.

This shoulder appears in almost any fat-tailed distribution [11]. Accordingly, a "power law fit" indeed involves *two parameters*: a threshold  $x_{min}$  marking the end of the shoulder and the tail's slope.

Unfortunately, the large argument tail tends to be of poor quality, noisy, brutish and short. Very rarely one finds a reliable dataset that allows for a good quality fit. Indeed, a recent metaanalysis by Clauset, Shalizi and Newman [19] reveals that, among 20 canonical datasets considered, only in one case a power law fit is really convincing and in most cases other distributions are doing better.

We suggest that these obstacles reflect an essential shortcoming of the Yule-Simon theory: the neglect of "death" events. In reality species go extinct, individuals die and links break down, yet in the Yule-Simon theory this never happens. A death process cannot be taken into account by simply introducing a net birth rate; it also accounts for the stochastic extinction of existing families (genera). Yule theory thus overestimates the fraction of small families, which explains the typical "shoulder" that appears at small *m*'s.

Recently Manrubia and Zannete [1] studied the distribution of surnames in a population, using a model which is a specific example of the birth-death-mutation (BDM) process (see also [20]). We [15] then extended these results, showing that the resulting distribution is independent of the particular details of the process. In the spirit of Simon's realization that the Yule model results are applicable in a much broader context, we here propose, and demonstrate by numerous examples, that the BDM process and its resulting statistics should be applicable to a very wide range of empirical datasets.

#### BDM statistics: results and applications

1

Here is a list of the main results for the statistics of the BDM process, where the total population is growing/decaying at rate  $\gamma$ . In the supplementary material we resent a detailed description of the BDM dynamics and establish the equivalence between this process and preferential attachment [16] with the possibility of link removal.

- 1. The probability distribution function (the chance  $n_m$  to pick at random a family of size m) is described by the Kummer function U(a,b,c) [21].
  - (a) If the growth rate γ is larger than the mutation rate μ, an asymptotic power-law tail appears:

$$n(m) = \frac{vR_c\Gamma(2+v)}{m} U\left(1+v,0,\frac{R_cm}{N_0}\right)$$
$$\stackrel{m \to \infty}{\sim} m^{-(1+\frac{\gamma}{2}-\mu)}.$$

where  $v \equiv \mu/(\gamma - \mu)$  and  $R_c \equiv 2N_0|\gamma - \mu|/\sigma^2$ ,  $N_0$  is the current population size.

(b) For  $\mu > \gamma$ , the BDM dynamics supports a truncated power-law distribution [here  $v \equiv \gamma/(\mu - \gamma)$ ],

$$n(m) = \frac{R_c \Gamma(1+\nu)}{m} U\left(\nu, 0, \frac{R_c m}{N_0}\right) e^{-\frac{R_c m}{N_0}}$$
(2)  
$$\overset{m \to \infty}{\sim} m^{-1-\nu} e^{-\frac{2}{\sigma^2}(\mu-\gamma)m}.$$

2. When  $R_0$  individuals are sampled the effective strength of the sampling is  $s = R_0/R_c$ . In the strong sampling limit,  $s \gg 1$ , the new distribution is just a rescaled Kummer [15]. On the other hand if  $s \ll 1$ ,

$$n^{R}(m) \approx B(m-1-v,2+v)vR_{o}s^{v}.$$
 (3)



Figure 2. Tour de force of BDM statistics: Pareto plots are presented for empirical datasets obtained from independent studies across many disciplines. The best fit values of  $\gamma$  and  $\mu$  are given for each item. (a) Distribution of number of chromosome abberations in cancer tumors [26].  $\gamma = 0.28 \ \mu = 0.37$  (b) Surname statistics from the 1790 US census. The growth rate ( $\gamma = 0.034$ ) was inferred [15] from historical censuses in England, and the fit retrieves the "mutation" (surname changes) rate to be  $\mu = 0.011$ . (c) WWW: number of sites with certain degree of links as a function of the degree. The set of 200 million web pages with 1,500 million hyperlinks first considered by Broder et. al. [31] has been analyzed.  $\gamma = 0.27 \ \mu = 0.065$ . (d) Internet (physical structure) - number of nodes with *m* links vs. m. Data obtained from DIMES web site (www.netdimes.org).  $\gamma = 0.72 \ \mu = 0.51$ . (e) Clusters of trees in the tropical forest. Shown here is the number  $n_s$  of clusters of size *s* for *Hybanthus proinfolius*, the most frequent species in the Barro-Colorado Island plot [32]. (f) Species abundance ratio in the tropical forest [32]. Here  $\gamma = 5.4 \cdot 10^{-5} \ \mu = 1.5 \cdot 10^{-4}$ . (g) Human insurgency: number of Norwegian firms with *m* employees, as obtained from statistics Norway website, www.ssb.no. (Data for 2010).  $\gamma = 0.1\mu = 0.051$ . (h) Number of Norwegian firms with *m* employees, as obtained from statistics Norway website, www.ssb.no. (Data for 2010).  $\gamma = 0.11 \ \mu = 0.04$ . (i) Species within genera statistics for the Plantae kingdom [24]  $\gamma = 0.055\mu = 0.017$ .

Eq. (3) implies that the BDM statistics crosses over to the Yule-Simon result when the sampling is weak [see Eq. (1) and the discussion in methods (B)]. Since weak sampling yields mainly members of large families for which the chance of extinction is small, Yule's theory with a net birth rate becomes adequate. Indeed, in the main part of Fig. 1 we show how the BDM Kummer statistics fits the contemporary data for Cerambycinea and how one can reconcile the Yule result by taking into account the effect of sampling. Note that our theory [15] is based on a Fokker-Planck equation that fails when the size of the family is of order unity [22], thus here and in the following figures the curve fails to fit the number of singletons.

Fig. 2 demonstrates the power of our technique using many paradigmatic fat-tailed distributions from the social sciences (surnames, insurgency, WWW), engineering (internet), ecology (species within genera, species abundance ratio, clusters of trees), biology (cancer abberations statistics) and economy (firms size distribution). In all cases presented here a two parameter fit is shown, thus we are not using more fitting parameters than a standard power-law fit. In some cases the relevance of the BDM dynamics to the underlying process is clear; in other cases (terror attacks) the underlying process is not well understood, and more studies are needed in order to prove, or disprove, the relevance of BDM, perhaps along the lines suggested by [23]. The agreement of theory and data is impressing with respect to other fits on loglog scale; some examples of other fitting functions and distributions are given in the methods section (D).

Clearly the BDM theory is much stronger than a simple powerlaw fit, yielding sharper predictions and fitting almost perfectly many paradigmatic empirical datasets. Its amazing success, even



**Figure 3. A Pareto plot for the species within genera statistics for the Animalia kingdom.** The fit of the BDM theory to the data is surprisingly good, given the existence of different taxonomical classifications for genera. The fit suggests a diversification (speciation minus extinction) rate of about 0.063; this value falls within the confidence intervals obtained by Ricklefs [25] for North and South American clades of passerine birds. doi:10.1371/journal.pone.0026480.g003

where the BDM process is certainly a crude approximation for the real dynamics, suggests that this distribution behaves like a central limit for many multiplicative neutral processes.

For any of the topics of Fig. 2 a comprehensive discussion is needed in order to put our new results for  $\gamma$  and  $\mu$  in the context of the specific field. This is beyond the scope of this Letter, and short specific comments are presented in methods, subsection (C).

Let us conclude by demonstrating the quality of our results using one example. Figure 3 shows the species within genera statistics for all the Animalia kingdom [24]. The Kummer function fits almost exactly the empirical data, much better than other distributions conjectured (see SM). The rate of diversification (speciation minus extinction),  $\gamma = 0.063 \pm 0.02$ , is consistent with the range of values estimated from lineage through time plots [25], and our confidence intervals are much tighter.

#### **Materials and Methods**

#### A. The birth-death-mutation process

The birth-death-mutation (BDM) process, in its simplest form, governs the dynamics of S families of agents. Each family is characterized by m, the number of agents in it. For the sake of concreteness let us consider a population of species (agents), each of which belongs to a genus (family).

At every time step a species is chosen at random among all species, independent of its genus. This agent is removed with probability 1-p and reproduces (speciates) with probability p. The offspring belongs to the same genus as its parent species with probability  $1-\mu$ , and "mutates" to form a new genus with probability  $\mu$ . Note that we use the word "mutation" to indicate an offspring that forms a new family (genus, surname), rather than belonging to the same clan as its parent. The parameter  $\gamma = 2p - 1$  defines the growth rate (if positive) or the decay rate (if negative) of the population. This is the overlapping generations (Moran) version of the process.

Many other processes support the same steady state distribution of family sizes [15]. Of particular importance is the nonoverlapping generations (Wright-Fisher) version of this dynamics. In this case all agents produces offspring at once and then are removed. An agent produce *n* offspring with probability  $P_n$ . The average number of offspring per individual is thus given by  $\bar{n} = \sum nP_n$ ,



Figure 4. Animalia kingdoms statistics: Modified pareto (Zipf-Mandelbrot, dashed line) best fit vs. Kummer best fit. doi:10.1371/journal.pone.0026480.g004

and the growth/decay rate is  $\gamma = \bar{n} - 1$ . Again  $\mu$  is the mutation rate as described above.

In previous work [15] we have shown that all these processes yield the same steady-state distribution of family sizes, which is independent of the "microscopic" details. The final distribution depends only on the growth rate  $\gamma$ , the mutation rate  $\mu$ , and the variance  $\sigma^2 = Var(n)$ . For the Moran case  $\sigma^2 = 2$ . It turns out that n(m) satisfies the Kummer differential equation

$$\frac{\partial n(m)}{\partial t} = \frac{\sigma^2}{2} \frac{\partial^2}{\partial m^2} [mn(m)] + (\mu - \gamma) \frac{\partial}{\partial m} [mn(m)]. \tag{4}$$

Note that this equation resembles a diffusion-convection process for mn(m).

The same statistics emerges if agents are removed with probability 1-p, reproduce into the parent set with probability  $p(1-\mu)$ , and new agents, each deposited into an empty set (family), are added with probability  $p\mu$  (we refer to this as the birth-death-source process, BDS). This is the case, e.g., if nodes, each carrying a certain number of links, are added to an already existing network and the chance of a link to be attached to an already existing node is proportional to the degree of the node. If links are removed at a different rate, the process yields the same statistics as the BDM (up to slight modifications since new families appear, in realistic networks, with size which is greater than one).

The BDM process is a generalization of the famous Yule process which has no death in it; i.e., agents are only born and mutate. In the same sense, the BDS version generalizes the preferential attachment process [16] in which links are only added to the network but are never removed.

#### B. Yule-Simon statistics as a weak sampling limit of BDM

In the process defined by Yule there is no death, and the mutation rate  $\mu_{Yule}$  is simply the ratio between the average number of new surnames (or genera) that appear during a period of time and the number of new individuals added, during the same period, to already existing families (see the detailed discussion in [11]).

In the BDM process the rate in which new families are generated is  $\mu bN$  (*N* is the total population at certain time, *b* is the



**Figure 5. Out-degree statistics: The best fit to Kummer fails systematically at small** *ms.* doi:10.1371/journal.pone.0026480.g005

birth rate) and the rate in which the total population in the already existing families grows is  $[b(1-\mu)-d]N$ . Without loss of generality we can choose d=1 such that  $b=1+\gamma$ , since the growth rate  $\gamma \equiv b-d$ . The ratio between the new families generation rate and the old families growth rate is, (to the first order in the small parameters  $\gamma$  and  $\mu$ ),  $v \equiv \mu/(\gamma - \mu)$ . This implies that for small growth and mutation rates, which is the regime of validity of the Kummer theory, Yule theory is equivalent to BDM iff stochastic extinction is neglected and  $\mu$  Yule is replaced by v. For that reason, Eq. (5) of the main text is equivalent to Yule statistics (Eq. 1) with v instead of  $\mu$ .

#### C. 3 Remarks for Figure 2 of the main text

The remarks below refer to the panels of Fig. 2:

General: The binning of the data was done using a half logarithmic scale, which means that for small families ( $m \le 10$ ) we had a bin for every number, while for large families we used logarithmic binning with a bin size  $2^k$  (k is the bin number). We have found this to be optimal in terms of presentation clarity, but the Kummer fit has been checked using other binning schemes and the differences are negligible. For two datasets (surname panel (b), and firms panel (h)) the data was available only in a binned form, so the existing binning scheme has been retained.

(a) Cancer: The data we present here is the distribution of the number of chromosome abberations in cancer tumors [26], includes all different types of cancer. See [27] for analysis of different types of cancer.

- (b) Surname: The size of a family was defined as the number of households having the same surname. Data refer to the US census of 1790, when the US population shared the same genealogic and demographic histories with the British population. The English demography is roughly documented since the Domesday Book census carried out by William the Conquerer. For more details see [15].
- (c) WWW links statistics. There is some ambiguity about the kind of sampling involved in the collection of the data. In principle one should make a distinction between building a surname statistics by sampling *individuals* and asking for their surname, in which case Eq. (5) of the method section is applicable, and sampling surnames and asking for the number of individuals having this specific surname. In the internet case the sampling is done by crawlers moving from node to node along the links; here a link is an individual and a node is a "surname". In any case, the success of our fit to a full census theory means that the effect of sampling, if any, is weak (i.e., that we are in the strong sampling regime).
- (d) We present here the nodes in-degree distribution (i.e. the size of a node is determined by the number of links pointing to it). The nodes out-degree distribution does not follow Kummer. This difference needs further analysis.
- (c) The data presented here is for the most frequent species in the plot, *Hybanthus pronifolius*. There are about 40000 individual trees of this species in a rectangular area of  $1000 \times 500$  meters. We have covered the area by a

 $2 \times 2$  meters grid and consider any square that contains at least a single *Hybanthus* tree as black, other squares are white. We then identified and tracked black cluster using the standard Hoshen-Kopelman algorithm. We have checked that the results are not sensitive to slight modifications of the lattice constant (grids with 1–3 meters squares were checked) and have gotten fits of similar quality for the other frequent species in the plot.

(f) The data was averaged over six different censuses. Time between consecutive censuses is five years, to be compared with the lifetime of a tree which is typically about 100 years. Our best fit yields  $\gamma = 4.310^{-5}$  and  $\mu = 2.910^{-4}$ . This suggests that the total population of the meta-community isn't really fixed but rather grows extremely slowly. Although the model is neutral, the overall effect of adaptation may very slowly increase the carrying capacity of the forest.

While we are not trying to claim that our fit is actually conclusive, this result opens an interesting possibility for refutation of the critics of the "point mutation" version of Hubbell's theory, who base themselves on turnover rates. As pointed out by Ricklefs [28] and by Nee [29] the time to origination of a species with N individuals is about 2N. This leads to ridiculously large timescales when applied to realistic species abundance. One implication of our work is that the introduction of a very weak growth rate does not kill the statistics, yet it clearly shortens the time to origination significantly. For example for 10 million trees with generation time of a 100 years, the time to origination if the total population is fixed will be of order of a billion years, while for the  $\gamma$  above it will be 40 million years.

- (g) The datasets had also some non-integers values (the meaning of which is unclear to us) that we rounded up to the closest integer number.
- (h) The dataset includes the number of establishments with m employees, starting from m=0. In order to avoid this zero we have shifted  $m \rightarrow m+1$ , counting the owner also as an employee.
- (i) The statistics of the Plantae kingdom. This dataset is similar to the Animalia displayed and analyzed in Fig. 3; we have preferred to present a more detailed analysis of Animalia since this is the largest kingdom.

#### References

- Manrubia SC, Zanette DH (2002) At the boundary between biological and cultural evolution: The origin of surname distributions. J Theor Biol 216: 461–477.
- Raup DM (1992) Extinction: bad genes or bad luck. New York: WW Norton & Company.
- Darwin C (1859) The origin of species by means of natural selections. London: John Murrary.
- Kimura M (1985) The neutral theory of molecular evolution. Cambridge: Cambridge Univ. Press.
- Hubbell SP (2001) The unified neutral theory of biodiversity and biogeography. Princeton, NJ: Princeton Univ. Press.
- 6. Gibrat R (1931) Les inégalités économiques. Librairie du Recueil Sirey.
- 7. Taleb N (2005) Fooled by randomness: The hidden role of chance in the markets and life. New York: Random House.
- 8. Taleb N (2007) The Black Swan: The Impact of the Highly Improbable. New York: Random House.
- Yule GU (1925) A mathematical theory of evolution, based on the conclusions of dr. jc willis, frs. Phil Trans R Soc Lond B 213: 21–87.
- 10. Simon H (1955) On a class of skew distribution functions. Biometrika 42: 425–440.
- Newman MEJ (2005) Power laws, Pareto distributions and Zipf's law. Contemporary Phys 46: 323–351.

#### D. The adequacy of Kummer

When dealing with fat-tailed distributions that are extended over many orders of magnitude, a log-log plot must be used. However, these plots are notoriously known to smear out some fine details of the distribution, and sometimes this feature blurs the actual mismatch between the theory and the empirical data. The level of exactness is thus a crucial factor in determining the adequacy of a fit. Here we describe two examples.

First, in Fig. 4 the Kummer best fit is compared with the best fit obtained for the modified Pareto (Zipf - Mandelbrot) distribution, which is a two parameter law with the same concave shape,

$$n_m = N_0 \frac{(a+1)^{b-1}(b-1)}{(a+m)^b} \tag{5}$$

where  $N_0$  is the population size. The best fit for the parameters a and b is shown together with the best Kummer fit. One can see that, although the mismatch is never large in a loglog plot once the function captures the general trend, there are systematic deviations in the modified Pareto case but not from the Kummer function (note again that the singletons are not covered by our theory so the mismatch at m = 1 is irrelevant).

As another example let us present a case where systematic deviations from Kummer show up. In Fig. 5 the out-degree distribution of nodes in the internet (the in-degree that satisfies Kummer is shown in Fig. 2d) is shown together with the best fit to Kummer, and indeed one can see systematic deviations that makes the Kummer fit very suspicious, if not fully disqualified.

In general the Kummer function may be considered in any case where the distribution is monotonically decreasing (so it is inappropriate as an explanation to, say, scientific citation statistics where a hump appears at intermediate values of m). For a reasonable fit the slope at small m-s should be close to one, not too shallow (as in the Tsallis distribution [30]) or too steep.

#### Acknowledgments

We thank Robert Ricklefs and David Aldous for helpful comments and discussions.

#### **Author Contributions**

Conceived and designed the experiments: NS DK YM. Performed the experiments: NS DK YM. Analyzed the data: NS DK YM. Contributed reagents/materials/analysis tools: NS DK YM. Wrote the paper: NS DK YM.

- Blank A, Solomon S (2000) Power laws in cities population, financial markets and internet sites (scaling in systems with a variable number of components). Physica A: Statistical Mechanics and its Applications 287: 279–288.
- Solomon S, Richmond P (2002) Stable power laws in variable economies; lotkavolterra implies pareto-zipf. The European Physical Journal B - Condensed Matter and Complex Systems 27: 257–261.
- Zeldovich KB, Chen P, Shakhnovich BE, Shakhnovich EI (2007) A firstprinciples model of early evolution: Emergence of gene families, species, and preferred protein folds. PLoS Comp Biol 3: e139: 1224–1238.
- Maruvka YE, Shnerb NM, Kessler DA (2009) Universal features of surname distribution in a subsample of a growing population. J Theor Biol 262: 245–256.
- Barabási AL, Albert R (1999) Emergence of scaling in random networks. Science 286: 509–512.
- Simon HA, Bonini CP (1958) The size distribution of business firms. The American Economic Review 48: 607–617.
- Clauset A, Young M, Gleditsch KS (2007) On the Frequency of Severe Terrorist Events. Journal of Conflict Resolution 51: 58–87.
- Clauset A, Shalizi CR, Newman MEJ (2009) Power-law distributions in empirical data. SIAM Review 51: 661–703.
- Volkov I, Banavar JR, Hubbell SP, Maritan A (2003) Neutral theory and relative species abundance in ecology. Nature 424: 1035–1037.

- Abramowitz M, Stegun IA, eds. (1964) Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables. Number 55 in National Bureau of Standards Applied Mathematics Series. Washington, DC: Gov't. Printing Office, 10th edition.
- Kessler DA, Shnerb NM (2007) Extinction rates for fluctuation-induced metastabilities: A real-space wkb approach. Journal of Statistical Physics 127: 861–886.
- Bohorquez JC, Gourley S, Dixon AR, Spagat M, Johnson NF (2009) Common ecology quantifies human insurgency. Nature 462: 911–914.
  Bisby FA, Roskov YR, Orrell TM, Nicolson D, Paglinawan LE, et al. (2009)
- Bisby FA, Roskov YR, Orrell TM, Nicolson D, Paglinawan LE, et al. (2009) Species 2000 & ITIS Catalogue of Life: 2009 Annual Checklist. CD-ROM; Species 2000. Reading, U.K.: WW Norton & Company.
- Ricklefs RE (2007) Estimating diversification rates from phylogenetic information. Trends Ecol Evol (Amst) 22: 601–610.

- (2010) Mitelman database of chromosome aberrations and gene fusions in cancer. http://cgap.nci.nih.gov/Chromosomes/Mitelman|.
- Frigyesi A, Gisselsson D, Mitelman F, Hoglund M (2003) Power law distribution of chromosome aberrations in cancer. Cancer Research 63: 7094–7097.
- Ricklefs RE (2006) The unified neutral theory of biodiversity: Do the numbers add up? Ecology 87: 1424–1431.
- Sean N (2005) The neutral theory of biodiversity: Do the numbers add up? Functional Ecology 19: 173–176.
- Tsallis C (1988) Possible generalization of Boltzmann-Gibbs statistics. Journal of Statistical Physics 52: 479–487.
  Broder A, Kumar R, Maghoul F, Raghavan P, Rajagopalan S, et al. (2000)
- Graph structure in the web. Computer Networks 33: 309–320.
- 32. Hubbell SP, Condit R, Foster RB (2005) Barro colorado forest census plot data.