

1 **A macroecological theory of microbial biodiversity**

2  
3 William R. Shoemaker<sup>1\*</sup>, Kenneth J. Locey<sup>1\*</sup>, Jay T. Lennon<sup>1</sup>

4 <sup>1</sup>Department of Biology, Indiana University, Bloomington, IN 47405 USA

5 \*Authors contributed equally to the study

6 Correspondence: K Locey, Department of Biology, Indiana University, 261 Jordan Hall,  
7 1001 East 3rd Street, Bloomington, IN 47405 USA. E-mail: [kjlocey@indiana.edu](mailto:kjlocey@indiana.edu)

8  
9 Supplementary Figures

- 10 1. The relationship between the predicted abundance of each rank and the observed  
11 abundance for different sequence similarity cutoffs.
- 12 2. The relationship between the predicted abundance of each rank and the observed  
13 abundance with singletons removed.
- 14 3. The relationship between model performance and the percent value of  $N$  used to  
15 sample SADs.
- 16 4. Kernel density estimates for the lognormal, log-series, and Zipf distribution from  
17 a single bootstrapped sample.
- 18 5. Kernel density estimates for the fitted parameters from the lognormal, log-series,  
19 and Zipf distribution.
- 20 6. A box-and-whisker plot of the percent of the time that each SAD model had the  
21 highest AICc weight.

24 Supplementary Tables

- 25 1. The results of the simple linear regression of each SAD model from Figure 3 for a  
26 single bootstrap iteration.
- 27 2. Comparison of the performance of each species abundance distribution (SAD)  
28 model for different sequence similarity cutoffs.
- 29 3. Comparison of the performance of species abundance distribution (SAD) models  
30 for microbial datasets with singletons removed.
- 31 4. A comparison of the mean and standard deviation of the log-likelihood values for  
32 the lognormal, log-series, and Zipf distribution from 10,000 bootstrapped  
33 samples.
- 34 5. The mean and standard deviation of the fitted parameters for the lognormal, log-  
35 series, and Zipf distribution.
- 36 6. The mean and standard deviation of the fitted parameters for the lognormal, log-  
37 series, and Zipf distribution.

38

39

40

41

42

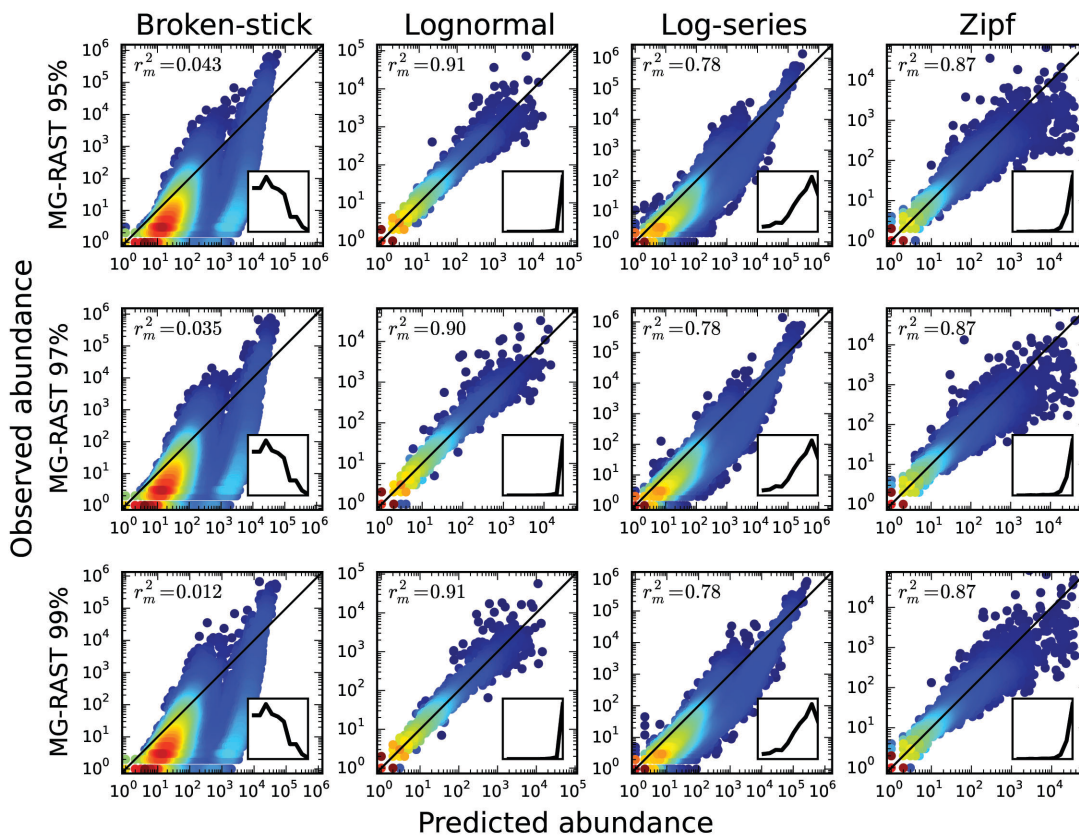
43

44

45

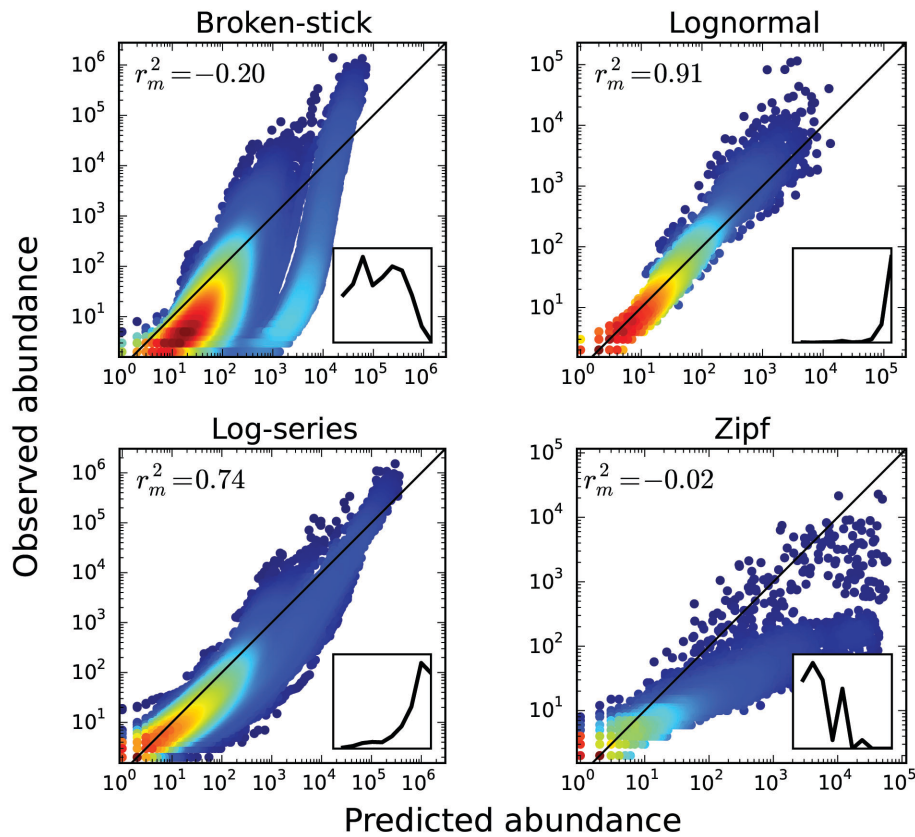
46

47 **Supplementary Figure 1.** The relationship between the predicted abundance of each  
 48 rank and the observed abundance for each SAD model for sequence similarity cutoffs of  
 49 95, 97, and 99%. Sequence similarity had no measurable effect on model performance.  
 50 Points are color-coded by the density of adjacent points. Hotter colors indicate a higher  
 51 density of points.  
 52



53  
 54  
 55  
 56  
 57

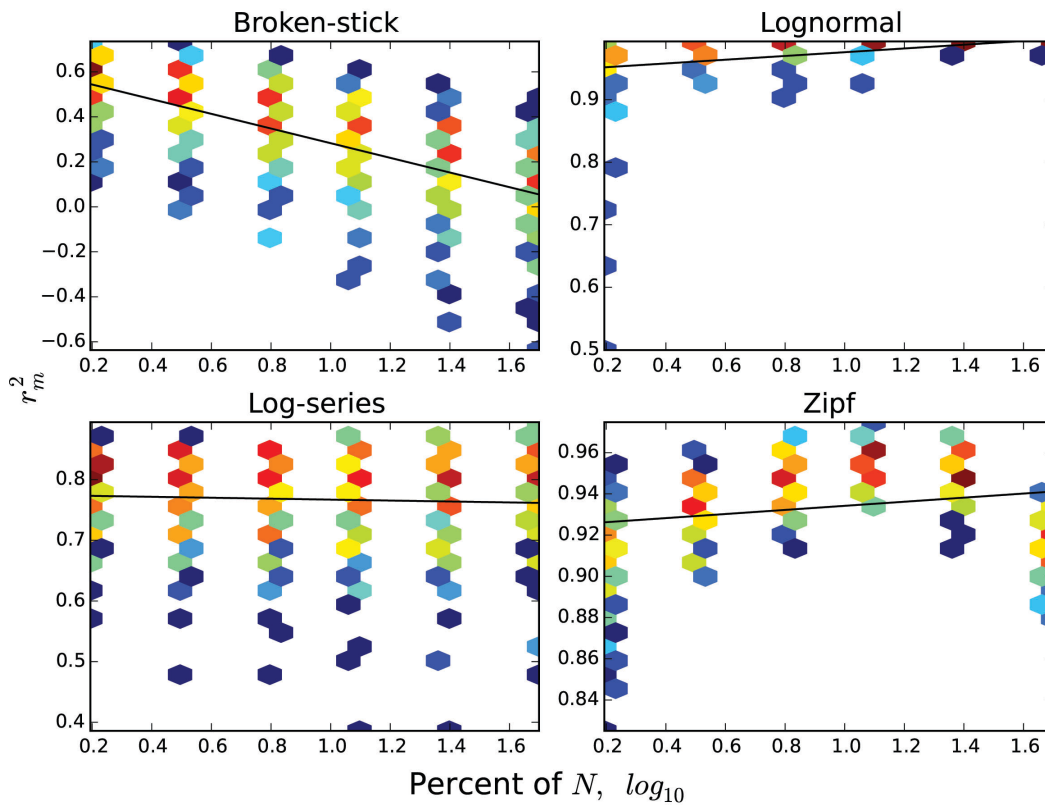
58 **Supplementary Figure 2.** The relationship between the predicted abundance of each  
 59 rank and the observed abundance for each SAD model with singletons (i.e. species that  
 60 have an abundance of one) removed. Removing singletons had little effect on the  
 61 performance of the lognormal distribution. The performance of the log-series and  
 62 Broken-stick increased while the performance of the Zipf distribution decreased, likely as  
 63 a result of the models tendency to under and over predict, respectively, the number of  
 64 singletons in a site. The value at the top-left of each sub-plot is the mean  $r_m^2$  value for  
 65 10,000 bootstrapped samples (see Methods). Points are color-coded by the density of  
 66 adjacent points. Hotter colors indicate a higher density of points.  
 67



68

69

70 **Supplementary Figure 3.** The relationship between model performance (i.e.  $r_m^2$ ) and the  
71 percent value of  $N$  used to sample SADs. The performance of the Broken-stick and the  
72 log-series decreases as  $N$  increases while the performance of the log-series increases as  $N$   
73 increases. While the performance of the Zipf displays a humped-shaped relationship,  
74 performance increases up until a cut-off of 25%. Each sub-plot contains the results of 100  
75 bootstrap iterations (see Methods).  
76

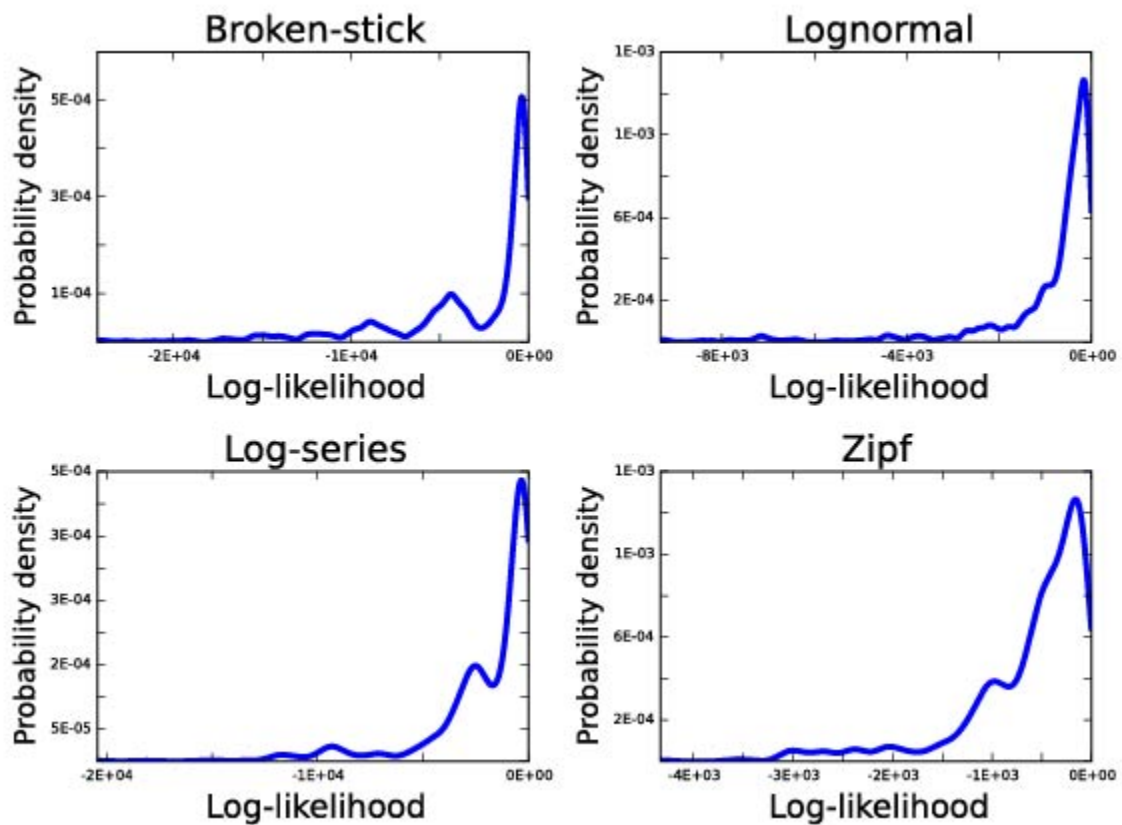


77  
78  
79  
80  
81

82 **Supplementary Figure 4.** Kernel density estimates for the lognormal, log-series, and  
83 Zipf distribution from a single bootstrapped sample. The lognormal and Zipf distribution  
84 had similar levels performance, while the log-likelihood of the log-series was larger by  
85 several orders of magnitude (Supplementary Table 4). In addition, the distribution of  
86 parameter values for all biodiversity models with a fitted parameter were highly peaked  
87 and unimodal. Kernel density estimates were selected based on cross-validation.

88

89



90

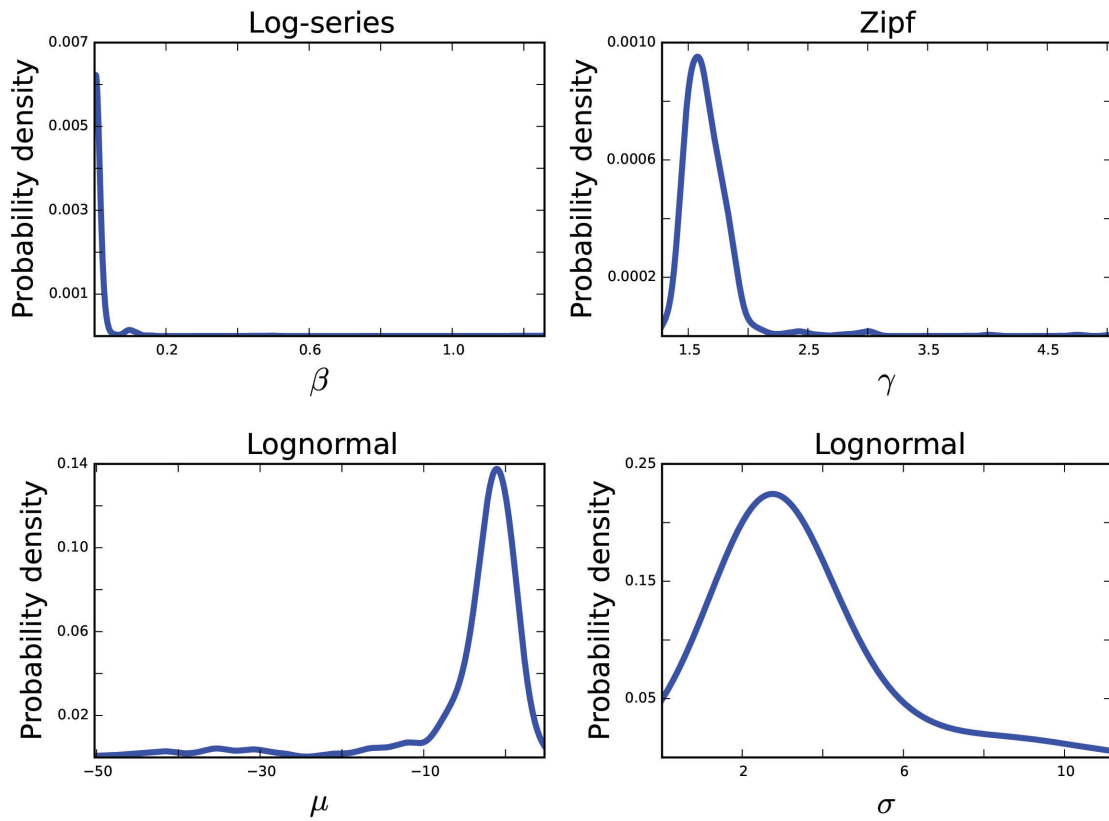
91

92

93

94 **Supplementary Figure 5.** Kernel density estimates for the fitted parameters from the  
95 lognormal, log-series, and Zipf distribution. The mean and standard deviation for each  
96 parameter can be found in Supplementary Table 5.

97



98

99

100

101

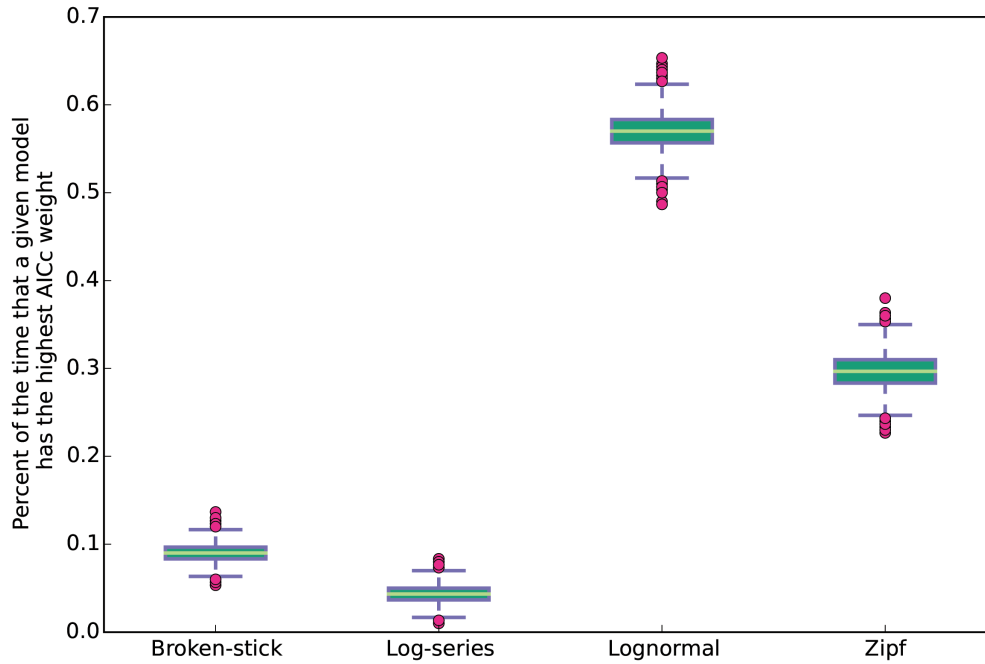
102

103

104

105

106 **Supplementary Figure 6.** A box-and-whisker plot of the percent of the time that each  
107 SAD model had the highest AICc weight. The mean and standard deviation of this data  
108 can be found in Supplementary Table 6.  
109



110  
111  
112  
113  
114  
115  
116  
117



118 **Supplementary Table 1.** The results of the simple linear regression of each SAD model  
119 from Figure 3 for a single iteration. Positive slopes indicate that model performance ( $r^2_m$ )  
120 increased as abundance ( $N$ ) increased. The lognormal and Zipf provide better  
121 explanations of microbial SADs as  $N$  increases than the Broken-stick or the log-series.

122

123

124

Model	Slope	$r^2$	$p$ - value
Lognormal	0.022	0.058	< 0.0001
Zipf	0.064	0.18	< 0.0001
Log-series	-0.010	0.00057	0.52
Broken-stick	-0.51	0.20	< 0.0001

125

126

127

128

129

130

131

132

133

134

135

136

137 **Supplementary Table 2.** Comparison of the performance of each species abundance  
 138 distribution (SAD) model for sequence similarity cutoffs of 95, 97, and 99%. Results are  
 139 reported as the reported as the mean and standard deviation of  $r_m^2$ . Sequence similarity  
 140 has little measurable effect on model performance. Additional information pertaining to  
 141 the table can be found in the description for Table 1.

142

143

Model	Sequence similarity	$\overline{r_m^2}$	$\sigma_{\overline{r_m^2}}$
Lognormal	95 %	0.91	0.11
	97 %	0.90	0.14
	99 %	0.91	0.12
Zipf	95 %	0.87	0.080
	97 %	0.87	0.11
	99 %	0.87	0.076
Log-series	95 %	0.78	0.18
	97 %	0.78	0.19
	99 %	0.78	0.18
Broken-stick	95 %	0.043	0.66
	97 %	0.034	0.67
	99 %	0.012	0.67

144

145

146

147

148

149

150 **Supplementary Table 3.** Comparison of the performance of species abundance  
151 distribution (SAD) models for microbial datasets with singletons removed. A sample of  
152 100 SADs was randomly sampled from each dataset 10,000 times and reported as the  
153 mean and standard deviation of  $r^2_m$ . Removing singletons had no measurable effect on  
154 model performance. Additional information pertaining to the results summarized in this  
155 table can be found in the description of Table 1.

156

157

158

Model	$\overline{r^2_m}$	$\sigma_{r^2_m}$
Lognormal	0.91	0.087
Zipf	-0.021	0.73
Log-series	0.74	0.28
Broken-stick	-0.20	0.82

159

160

161

162

163

164

165

166

167

168

169

170 **Supplementary Table 4.** A comparison of the mean and standard deviation of the log-  
 171 likelihood values for the lognormal, log-series, and Zipf distribution from 10,000  
 172 bootstrapped samples. Similar to the results using the  $r^2_m$ , the Zipf and lognormal  
 173 perform similarly well while the log-series performs poorly. These samples are drawn  
 174 from the set of sites that a given model was able to arrive at a prediction. The log-  
 175 likelihood values reported here are for all sites where numerical estimation arrived a  
 176 solution for a given model.

177

<b>Model</b>	$\mu_{Log-likelihood}$	$\sigma_{Log-likelihood}$
Lognormal	-970	47
Zipf	-630	19
Log-series	-2000	89
Broken-stick	-3400	130

182

183

184

185

186

187

188

189

190

191

192

193 **Supplementary Table 5.** The mean and standard deviation of the fitted parameters for  
194 the lognormal, log-series, and Zipf distribution. The Broken-stick has no fitted  
195 parameters and is not included in this table.

196

197

198

199

Model	Parameter	Mean	Standard deviation
Lognormal	$\mu$	1.8	0.044
	$\sigma$	1.9	0.029
Zipf	$\gamma$	1.4	0.0038
Log-series	$\beta$	0.0055	0.00069

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216 **Supplementary Table 6.** The mean and standard deviation for the percent of the time  
217 that a given SAD model has the highest corrected Akaike Information Criterion (AICc)  
218 weight from 10,000 bootstrap samples. Because the AICc needs to be calculated on a  
219 site-by-site basis, the values in this table were calculated using the set of sites where all  
220 models arrived at a successful prediction. AICc values were calculated using the set of  
221 sites where numerical estimation arrived at a successful prediction for all SAD models.

222

223

224

225

226

<b>Model</b>	$\mu\%$ <i>winning</i>	$\sigma\%$ <i>winning</i>
Lognormal	57.0	2.1
Zipf	30.0	2.0
Log-series	4.3	1.0
Broken-stick	9.0	1.1

227

228

229

230

231

232

233

234

235

236

237

238