

Zipf Distribution of U.S. Firm Sizes

Robert L. Axtell

Analyses of firm sizes have historically used data that included limited samples of small firms, data typically described by lognormal distributions. Using data on the entire population of tax-paying firms in the United States, I show here that the Zipf distribution characterizes firm sizes: the probability a firm is larger than size s is inversely proportional to s . These results hold for data from multiple years and for various definitions of firm size.

Firm sizes in industrial countries are highly skewed, such that small numbers of large firms coexist alongside larger numbers of smaller firms. Such skewness has been robust over time, being insensitive to changes in political and regulatory environments, immune to waves of mergers and acquisitions (1), and unaffected by surges of new firm entry and bankruptcies. It has even survived large-scale demographic transitions within work forces (e.g., women entering the labor market in the United States) and widespread technological change. The firm size distribution within an industry indicates the degree of industrial concentration, a quantity of particular interest for antitrust policy.

Beginning with Gibrat (2), firm sizes have often been described by lognormal distributions. This distribution is a consequence of the “law of proportional effect,” also known as Gibrat’s law, whereby firm growth is treated as a random process and growth rates are independent of firm size (3). Such distributions are skewed to the right, meaning that much of the probability mass lies to the right of the modal value. Thus, the modal firm size is smaller than the median size, which, in turn, is smaller than the mean.

The upper tail of the firm size distribution has often been described by the Yule (1) or Pareto (also known as power law, or scaling) distributions (4, 5). For a discrete Pareto-distributed random variable, S , the tail cumulative distribution function (CDF) is

$$\Pr[S \geq s_i] = \left(\frac{s_0}{s_i}\right)^\alpha, \quad s_i \geq s_0, \quad \alpha > 0 \quad (1)$$

where s_0 is the minimum size (6). Recent analysis of data on the largest 500 U.S. firms gives α as ~ 1.25 , whereas it is closer to 1 for many other countries (7). The special case of $\alpha = 1$ is known as the Zipf distribution and has somewhat unusual properties insofar as its moments do not exist (8). This distribution describes surprisingly diverse natural and so-

cial phenomena, including percolation processes (9), immune system response (10), frequency of word usage (4), city sizes (4, 11), and aspects of Internet traffic (12).

From an analysis using a sample of firms in Standard & Poor’s COMPUSTAT, a commercially available data set, it has been reported that U.S. firm sizes are approximately lognormally distributed (13). The COMPUSTAT data cover nearly all publicly traded firms in the United States—some 10,776 firms in 1997, almost 4300 of which had more than 500 employees. Firms covered by COMPUSTAT collectively employed over 52 million people, approximately one-half of the U.S. work force. However, these data are unrepresentative of the overall population of U.S. firms. Data from the U.S. Census Bureau put the total number of firms that had employees sometime during 1997 at about 5.5 million, including over 16,000 having more than 500 employees. Furthermore, the Census data have a qualitatively different character than the COMPUSTAT data. Census data display monotonically increasing numbers of progressively smaller firms, a shape the lognormal distribution cannot reproduce, and suggesting that a power law distribution may apply. As shown in Table 1 (14), the mean firm size in the COMPUSTAT data is 4605 employees (6349 for firms larger than 0), whereas in the Census data it is

Table 1. U.S. firm size distribution in 1997, compared across data sources. Number of firms in various size categories, with size defined as the number of employees, comparing COMPUSTAT and U.S. Census Bureau data for 1997. Note that there are monotonically decreasing numbers of progressively larger firms in the Census data, whereas this is not the case in the COMPUSTAT data (29).

Size class	COMPUSTAT	Census
0	2,576	719,978
1 to 4	123	2,638,070
5 to 9	149	1,006,897
10 to 19	251	593,696
20 to 99	1,287	487,491
100 to 499	2,123	79,707
500+	4,267	16,079
Total	10,776	5,541,918

Center on Social and Economic Dynamics, The Brookings Institution, 1775 Massachusetts Avenue, NW, Washington, DC 20036, USA.

Correspondence should be addressed to raxtell@brookings.edu.

14. Elastic upper crust and viscoelastic lower crust thicknesses are 16 km and 14 km, respectively. The assigned 16-km thickness is consistent with the cutoff depth of seismicity between about 15- and 20-km depth [K. B. Richards-Dinger, P. M. Shearer, *J. Geophys. Res.* **105**, 10939 (2000)], thought to coincide with the brittle-to-ductile transition. Depth-dependent elastic parameters are constrained by seismic information [J. Qu, T. L. Teng, J. Wang, *Bull. Seismol. Soc. Am.* **84**, 596 (1994)], and viscosities η_c and η_m are variable.
15. The 16-km upper depth is chosen to coincide with the base of the coseismic slip zone, and the 36-km depth is chosen to permit postseismic slip that is deep enough to reproduce the long-wavelength pattern of surface displacement seen on the interferograms.
16. F. F. Pollitz, G. Peltzer, R. Bürgmann, *J. Geophys. Res.* **105**, 8035 (2000).
17. We tested possible stratification within the lower crust by considering an additional model in which the lower crust consists of two uniform layers of identical thickness, with the upper layer three times as viscous as the lower layer. In this case, we found that optimal mantle viscosity is in the same range as the present model (3 to 8×10^{17} Pa s), whereas the lower crustal viscosities are 1.7×10^{19} Pa s and 0.6×10^{19} Pa s. The weaker crustal layer is still much more viscous than the mantle.
18. B. G. Bills, D. R. Currey, G. A. Marshall, *J. Geophys. Res.* **99**, 22059 (1994).
19. T. S. James, J. J. Clague, K. Wang, I. Hutchinson, *Quat. Sci. Rev.* **19**, 1527 (2000).
20. G. Kaufmann, F. Amelung, *J. Geophys. Res.* **105**, 16341 (2000).
21. G. L. Farmer et al., *J. Geophys. Res.* **100**, 8399 (1995).
22. B. L. Beard, C. M. Johnson, *J. Geophys. Res.* **102**, 20149 (1997).
23. A. F. Glazner et al., *J. Geophys. Res.* **96**, 13673 (1991).
24. S. H. Kirby, A. K. Kronenberg, *Rev. Geophys.* **25**, 1219 (1987).
25. S. M. Nakiboglu, K. Lambeck, *J. Geophys. Res.* **88**, 10439 (1983).
26. L. Block, L. H. Royden, *Tectonics* **9**, 557 (1990).
27. J. McCarthy et al., *J. Geophys. Res.* **96**, 12259 (1991).
28. L. H. Royden et al., *Science* **276**, 788 (1997).
29. M. K. Clark, L. H. Royden, *Geology* **28**, 703 (2000).
30. P. B. Gans, W. A. Bohrsen, *Science* **279**, 66 (1998).
31. D. McKenzie et al., *J. Geophys. Res.* **105**, 11029 (2000).
32. Sixteen continuous GPS time series are provided by Southern California Integrated GPS Network (SCIGN) (<http://pasadena.wr.usgs.gov/scign/cgi-bin/datafile.cgi>), and 13 campaign GPS times series are provided by U.S. Geological Survey (USGS) (<http://quake.usgs.gov/research/deformation/gps/auto/HectorMine>) covering the time period 20 October 1999 to 21 June 2000 or a portion thereof. The SCIGN measurements are referenced to ITRF97 [C. Boucher, Z. Altamimi, P. Sillard, in *IERS Technical Note 27* (Observatoire de Paris, Paris, France, 1999)]. The USGS measurements are referenced to the SCIGN measurements by aligning the velocity fields at three common sites, achieving a consistency of about 1 mm/year between the two velocity fields.
33. D. Massonnet, W. Thatcher, H. Vadon, *Nature* **382**, 612 (1996).
34. We thank D. Sandwell, W. Prescott, J. Svare, K. Hudnut, N. King, and S. Kirby for helpful discussions; J. Savage and R. Stein for constructive reviews; D. Dreger for providing the Hector Mine coseismic model; and G. Bawden for assistance with graphics. We acknowledge the Southern California Integrated GPS Network and its sponsors, the W. M. Keck Foundation, NASA, NSF, the U.S. Geological Survey, and the Southern California Earthquake Center, for providing GPS data used in this study. This paper benefited from constructive criticisms by two anonymous reviewers.

5 April 2001; accepted 30 July 2001

REPORTS

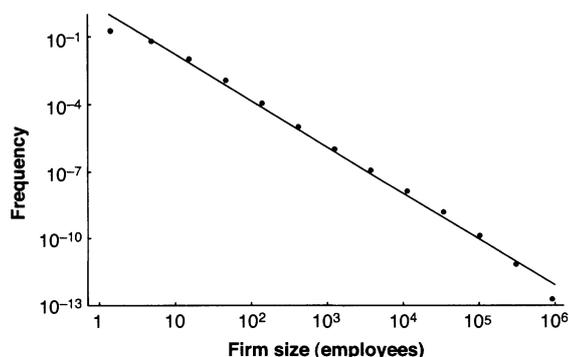


Fig. 1. Histogram of U.S. firm sizes, by employees. Data are for 1997 from the U.S. Census Bureau, tabulated in bins having width increasing in powers of three (30). The solid line is the OLS regression line through the data, and it has a slope of 2.059 (SE = 0.054; adjusted $R^2 = 0.992$), meaning that $\alpha = 1.059$; maximum likelihood and nonparametric methods yield similar results. The data are slightly concave to the origin in log-log coordinates, reflecting finite size cutoffs at the limits of very small and very large firms.

19.0 (21.8 for firms larger than 0). Clearly, the COMPUSTAT data are heavily censored with respect to small firms. Such firms play important roles in the economy (15, 16).

For further analysis, I used a tabulation from Census in which successive bins are of increasing size in powers of three. The modal firm size is 1, whereas the median is 3 (4 if size 0 firms are not counted) These data are approximately Zipf-distributed ($\alpha = 1.059$), as determined by ordinary least squares (OLS) regression in log-log coordinates (Fig. 1). There are too few very small and very large firms with respect to the Zipf fit, presumably due to finite size effects, yet the power law distribution well describes the data over nearly six decades of firm size (from 10^0 to 10^6 employees). This result suggests both that a common mechanism of firm growth operates on firms of all sizes, and that the fundamental unit of analysis is the individual employee.

But firms having a single employee are not the smallest economic entities in the U.S. economy. Although there were some 5.5 million firms that had at least one employee at some time during 1997, there were another 15.4 million business entities in that year with no employees. These are predominantly self-employed individuals and partnerships, and are called “nonemployer” firms by Census. These smallest of firms account for nearly \$600 billion in receipts in 1997. Yet, if these firms are included in the overall firm size distribution, the Zipf distribution still fits the data well. To see this, Eq. 1 must be modified to accommodate firms having no employees

$$\Pr[S \geq s_i] = \left(\frac{s_0}{s_i + 1} \right)^\alpha, \quad s_i \geq 0, \alpha > 0 \quad (2)$$

Table 2. Power law exponent for U.S. firms in 1992, firms with employees and all firms. Results using OLS regression on Census data, with standard errors in parentheses.

Type	Estimated α	Adjusted R^2
Firms with employees	0.994 (0.043)	0.995
All businesses	0.995 (0.031)	0.994

Here, OLS yields an estimate of $\alpha = 1.098$ (SE = 0.064), and the adjusted $R^2 = 0.977$. Including self-employment drives the average firm size down to 5.0 employees/firm, and makes the median number of employees 0.

An interesting property of firm size distributions noted in previous studies of large firms is that the qualitative character of such distributions is independent of how size is defined (1). Although the position of individual firms in a size distribution does depend on the definition of size, the shape of the distribution does not. This also holds for the Census data. Basing firm size on receipts, a Zipf distribution describes the data ($\alpha = 0.994$) (Fig. 2). Here, modal and median firm revenues are each less than \$100,000, and the average is \$173,000/firm.

As a further test on the robustness of these results, I repeated these analyses for Census data from 1992. Average firm size was slightly smaller then, at 20.9 employees/firm (excluding size 0 firms). But overall, the Zipf distribution is as strong (Table 2).

Virtually all U.S. firms experienced significant changes in revenue and work force from 1992 to 1997. Thus, individual firms migrated up and down the Zipf distribution, but economic forces seem to have rendered any systematic deviations from it short-lived. Even the substantial merger and acquisition activity of this period seemed to have little

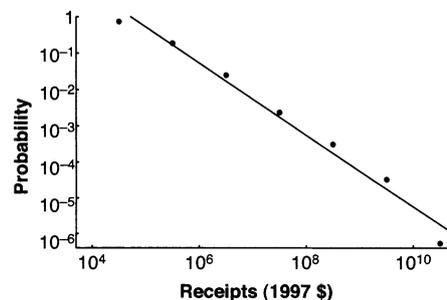


Fig. 2. Tail cumulative distribution function of U.S. firm sizes, by receipts in dollars. Data are for 1997 from the U.S. Census Bureau, tabulated in bins whose width increases in powers of 10. The solid line is the OLS regression line through the data and has slope of 0.994 (SE = 0.064; adjusted $R^2 = 0.976$).

effect on the overall firm size distribution.

There are a variety of stochastic growth processes that converge to Pareto and Zipf distributions (1, 5, 17, 18). Empirically, there is support for Gibrat-like processes in which average growth rates are independent of size (19, 20) and growth rate variance declines with size (21, 22). Consider a variation of the Gibrat process known as the Kesten process (23-25), in which sizes are bounded from below; i.e.,

$$s_i(t + 1) = \max[s_0, \gamma(t)s_i(t)] \quad (3)$$

where γ is a random growth rate. For nearly any growth rate distribution, this process yields Pareto distributions that have the exponent α defined implicitly by (26)

$$N = \frac{\alpha - 1}{\alpha} \left[\frac{\left(\frac{s_0}{A} \right)^\alpha - 1}{\left(\frac{s_0}{A} \right)^\alpha - \left(\frac{s_0}{A} \right)} \right] \quad (4)$$

where N is the total number of firms and A is the number of employees. For $N = 5.5 \times 10^6$ and $A = 105 \times 10^6$, as in 1997 (excluding self-employment), $s_0 = 1$ implies $\alpha \approx 0.997$, a value close to my empirical finding. Similar results are obtained for each year back through 1988 (Table 3).

Table 3. Theoretical power law exponents for U.S. firms over a 10-year period. Note that even though the number of firms and total employees each increased over this period, as did the average firm size, the value of α was approximately unchanged.

Year	Firms	Employees	Mean firm size	α , from (4)
1997	5,541,918	105,299,123	19.00	0.9966
1996	5,478,047	102,187,297	18.65	0.9986
1995	5,369,068	100,314,946	18.68	0.9983
1994	5,276,964	96,721,594	18.33	1.0004
1993	5,193,642	94,773,913	18.25	1.0008
1992	5,095,356	92,825,797	18.22	1.0009
1991	5,051,025	92,307,559	18.28	1.0004
1990	5,073,795	93,469,275	18.42	0.9995
1989	5,021,315	91,626,094	18.25	1.0006
1988	4,954,645	87,844,303	17.73	1.0039

REPORTS

The Zipf distribution is an unambiguous target that any empirically accurate theory of the firm must hit. This result, taken together with those in (21) and (27), place important limits on models of firm dynamics. That is, (i) firm growth rates follow a Laplace distribution, (ii) the standard deviation in growth rates falls with initial firm size according to a power law, and (iii) large firms pay higher wages for the same job according to yet another power law (the so-called wage-size effect). Because the Zipf distribution obtains all the way down to the smallest sizes, it should be possible to derive Kesten-type processes and, hence, the Zipf distribution from a microeconomic model in which individual agents interact to form productive teams. Although today no analytically tractable models of this type exist, agent-based computational results have achieved significant success according to these criteria (28).

The Zipf distribution may describe firm sizes in other countries as well, a conjecture that can only be tested once individual governments make available—and in some cases gather for the first time—data that purport to be comprehensive.

References and Notes

1. Y. Ijiri, H. A. Simon, *Skew Distributions and the Sizes of Business Firms* (North-Holland, New York, 1977).
2. R. Gibrat, *Les Inégalités Économiques; Applications: aux inégalités des richesses, à la concentration des entreprises, aux populations des villes, aux statistiques des familles, etc., d'une loi nouvelle, la loi de l'effet proportionnel* (Librairie du Recueil Sirey, Paris, 1931).
3. J. Sutton, *J. Econ. Lit.* **XXXV**, 40 (1997).
4. G. K. Zipf, *Human Behavior and the Principle of Least Effort* (Addison-Wesley, Reading, MA, 1949).
5. J. Steindl, *Random Processes and the Growth of Firms* (Hafner, New York, 1965).
6. N. L. Johnson, S. Kotz, N. Balakrishnan, *Continuous Univariate Distributions* (Wiley, New York, ed. 2, 1994).
7. J. J. Ramsden, G. Kiss-Haypál, *Physica A* **277**, 220 (2000), who use a functional form slightly different from (1).
8. Although any finite sample will have moments, by definition, the nonexistence of moments in the context of real data implies that the moments give no indication of convergence as the number of data increase.
9. M. S. Watanabe, *Phys. Rev. E* **53**, 4187 (1996).
10. J. D. Burgos, P. Moreno-Tovar, *Biosystems* **39**, 227 (1996).
11. M. Gell-Mann, *The Quark and the Jaguar* (Freeman, New York, 1994), pp. 92–97.
12. L. Breslau et al., *Proceedings of INFOCOM '99*, 21 to 25 March 1999, New York (IEEE Press, Piscataway, NJ, 1999), vol. 1, pp. 126–134.
13. M. H. R. Stanley et al., *Econ. Lett.* **49**, 453 (1995).
14. Census data is based on Small Business Administration (SBA) tabulations; available at www.sba.gov/advo/stats/data.html.
15. Z. Acs, D. Audretsch, *Innovation and Small Firms* (MIT Press, Cambridge, MA, 1990).
16. Z. Acs, Ed., *Are Small Firms Important? Their Role and Impact* (Kluwer Academic, Boston, 1999).
17. M. Marsili, Y.-C. Zhang, *Phys. Rev. Lett.* **80**, 2741 (1998).
18. H. Takayasu, K. Okuyama, *Fractals* **6**, 67 (1998).
19. P. E. Hart, S. J. Prais, *J. R. Stat. Soc. Ser. A* **119**, 150 (1956).
20. P. E. Hart, N. Oulton, *Econ. J.* **106**, 1242 (1996).
21. M. H. R. Stanley et al., *Nature* **379**, 804 (1996).
22. L. A. N. Amaral et al., *J. Phys. I France* **7**, 621 (1997).
23. H. Kesten, *Acta Mathematica* **131**, 207 (1973).
24. O. Biham, O. Malcai, M. Levy, S. Solomon, *Phys. Rev. E* **58**, 1352 (1998).
25. X. Gabaix, *Q. J. Econ.* **CXIV**, 739 (1999).
26. O. Malcai, O. Biham, S. Solomon, *Phys. Rev. E* **60**, 1299 (1999).
27. C. Brown, J. Medoff, *J. Pol. Econ.* **97**, 1027 (1989).
28. R. L. Axtell, in preparation; available at www.brookings.edu/dynamics/papers/firms.
29. The Census data were gathered in March of 1997. Firms that had receipts during 1997 but no employees as of March are shown in the size 0 category. Such firms should be in one of the other size classes. One might assume that it is possible to adjust the data by including these firms in the overall distribution by having them follow the Zipf distribution, for instance. However, any such procedure leads to the unrealistic conclusion that some of these temporarily size 0 firms actually have thousands or tens of thousands of employees. The firms in the size 0 category in COMPUSTAT are ostensibly holding companies.
30. These data were created by the U.S. Census Bureau under contract to the Brookings Institution. Given that the bins are not equally sized, construction of the probability mass function shown in Fig. 1 proceeds by taking the number of firms in each bin and dividing by the width of the bin. The resulting adjusted frequency is then located at the geometric mean of the bin endpoints. The tail CDF shown in Fig. 2 was constructed by cumulating the raw population data.
31. The late H. A. Simon initiated my interest in this subject. Lectures on Zipf's law by M. Gell-Mann at the Santa Fe Institute and conversations with B. Mandelbrot, B. Morel, and P. Bak were formative in my thinking. I thank Z. Acs, T. Åstebro, W. Dickens, C. Graham, J. Lanjouw, F. Pryor, J. Roth, and H. P. Young for suggestions, and T. Cole, R. Constantino, R. Hammond, and K. Landis for assistance. Support from the National Science Foundation, the Alex C. Walker Foundation, and the John D. and Catherine T. MacArthur Foundation is gratefully acknowledged.

30 April 2001; accepted 9 August 2001

Segregation of Human Neural Stem Cells in the Developing Primate Forebrain

Václav Ourednik,^{1*†} Jitka Ourednik,^{1*} Jonathan D. Flax,¹ W. Michael Zawada,² Cynthia Hutt,² Chunhua Yang,¹ Kook I. Park,^{1,3} Seung U. Kim,⁴ Richard L. Sidman,⁵ Curt R. Freed,^{2†‡} Evan Y. Snyder^{1†‡}

Many central nervous system regions at all stages of life contain neural stem cells (NSCs). We explored how these disparate NSC pools might emerge. A traceable clone of human NSCs was implanted intraventricularly to allow its integration into cerebral germinal zones of Old World monkey fetuses. The NSCs distributed into two subpopulations: One contributed to corticogenesis by migrating along radial glia to temporally appropriate layers of the cortical plate and differentiating into lamina-appropriate neurons or glia; the other remained undifferentiated and contributed to a secondary germinal zone (the subventricular zone) with occasional members interspersed throughout brain parenchyma. An early neurogenetic program allocates the progeny of NSCs either immediately for organogenesis or to undifferentiated pools for later use in the "postdevelopmental" brain.

As cells with stemlike qualities have come to be identified within a widening range of organs [e.g., (1, 2)], new questions have arisen about their relevance to normal development. The central nervous system (CNS) may serve as a bellwether for insights in this field. NSCs have been identified in the mammalian CNS, including humans (3–9), at stages from fetus to adult in a surprisingly wide range of regions (10–13). NSCs, defined as self-renewing, propagatable primordial cells each with the capacity to give rise to differentiated progeny within all neural lineages in all regions of the neuraxis, are posited to exist in the embryonic and fetal ventricular germinal zone (VZ) where they participate in CNS organogenesis (5, 14, 15). Cells equally "stemlike" in their potential have been identified at later stages (including old age) from

a variety of regions: subventricular (SVZ) (13–17) and ependymal (18) zones of the forebrain, subgranular zone of the hippocampus (6–10, 19), retina (20) and optic nerve (10, 11), cerebellum (12), spinal cord (21), and even cortical parenchyma (10, 15, 22). How might these observations be reconciled? Are such stemlike pools, particularly those isolated from various parenchymal regions at "postdevelopmental" periods, of physiological relevance or artifacts of experimental manipulation (10, 11)? Do these populations represent the same lineage or unique pools (17)? Of what relevance are these cells to normal human CNS development and repair?

We hypothesized that multiple stem cell pools, descendants of a common NSC, emerge during early cerebrogenesis as cells are used in organogenesis and concurrently