

"A man is
great by
deeds, not by
birth"

-Chanakya

Welcome to IIMK



INDIAN INSTITUTE OF MANAGEMENT KOZHIKODE



Working Paper

IIMK/WPS/223/EA/2017/07

January 2017

A Survey into Evidence of Zipf's Law among Indian Socio-Economic Variables

Kausik Gangopadhyay¹

¹Associate Professor, Economics Area, Indian Institute of Management Kozhikode, IIMK Campus P.O, Kerala – 673570, India, E-mail: kausik@iimk.ac.in, Phone: +91 - 495 - 2809118

IIMK WORKING PAPER

A Survey into Evidence of Zipf's Law among Indian Socio-Economic Variables

Kausik Gangopadhyay

Associate Professor of IIM Kozhikode

Abstract

Zipf's Law is an empirical phenomenon observed in many natural systems. The distribution of a physical variable demonstrates sharp rise at the right tail under this law. The occurrence of this law is pervasive among the physical variables. Econophysics, a discipline named so by Eugene H. Stanley, studies the application of physical principles among variables related to human action mostly related to socio-economic variables. This paper surveys the studies on the existence of Zipf's law among Indian socio-economic variables. We present the evidence on economic variables, in particular income, wealth and consumption distribution. The other socio-economic variable of our choice is city size distribution. In all cases, the Zipf's law is established with different values for the Pareto exponent.

Keywords:

Econophysics; sociophysics; Zipf's Law; Power Law; Pareto exponent; income distribution; wealth distribution; consumption distribution; city size distribution.

Introduction

Scientists attempt abstraction. Abstraction from data-points leads to theory. This is strictly deductive logic. However, often the posited theory is found to have multiple applications elsewhere and then this becomes an inductive theory.

The origin of Zipf's Law is deductive. Linguist George Kingsley Zipf (1) has observed that the second most common word in the English language, "of", appears at approximately half the times compared to the most common word, "the", in the English usage. Moreover, the third most common word, "to", appears at approximately one third the rate of the most common word. The origin of this law was deductive. However, the universality of this law among the physical variables and the variables of human action cannot be emphasised properly.

Zipf (2) himself found the wide application of this law in various human languages. Moreover, the size distribution of many physical variables such as the size distribution of islands (3), forest fires (8), lunar craters (7) and solar flares (9), are characterised by the presence of the Zipf's law. Many human actions such as websurfing (4), football goals (10), and opening moves in a chess game (11) are distributed with Zipf's law being prominent in the distribution.

Many socio-political variables have been posited and found to follow the Zipf's Law across many countries. This is highly interesting as the socio-political variables are outcome of human institutions and human mechanisms. If in spite of diversity in human institutions, the outcome is quite similar, it may mean some notion of equilibrium regarding socio-political variables. Many scientists are excited of this possibility, most notably Eugene H. Stanley, Bikas K. Chakrabarti, Victor Yakovenko, e Jean-Philippe Bouchaud, Bikas K Chakrabarti, J. Doyne Farmer, Dirk Helbing, János Kertész, Francis Longstaff, Rosario N. Mantegna, Matteo Marsili, to name a few. The discipline is named (30) by Eugene H. Stanley as Econophysics at a conference held in Kolkata, India. A salient example of success of Econophysics is its prediction regarding the distribution of wealth and income in many countries (5, 6) which demonstrate salient evidence of Zipf's Law. The other socio-economic variables include city size distribution (14, 15, 16, 17, 18, 19), expenditure distribution (12) etc.

In this article, we enunciate Zipf's law and then discuss its various properties. We, then, present our survey of the prominent available work regarding socio-political variables in India.

Our focus is mainly on Income and Wealth Distribution, City Size Distribution, consumption distribution etc.

Zipf's Law and Its Estimation

Let $f(\cdot)$ be a probability density function of a size distribution. The corresponding cumulative distribution function (CDF) and the complementary cumulative distribution function (CCDF) are given by $F(\cdot)$ and $F^C(\cdot)$, respectively. The CDF, $F(x)$, is the probability that an entity has a size less than or equal to x and the CCDF, $F^C(x)$, is the probability that an entity has a size greater than x .

$$\text{By definition, } F(x) = \int_0^x f(y)dy \text{ and } F^C(x) = 1 - F(x).$$

Zipf's law posits that $f(x)$ is proportional to $x^{-\alpha}$ and $F^C(x)$ is proportional to $x^{-\alpha+1}$, where α is a constant called the exponent of power law. This family of distributions with $\alpha > 1$ are known as the Pareto distribution. Evidently, $F^C(x)$ diverges to infinity for any $\alpha > 1$ as the argument of the CDF approaches zero. Therefore, some minimum value, say x_{\min} , lies at the lower end of the Pareto distribution as the support of the distribution; in the upper end of the Pareto distribution, no such bound exists.

If α assumes the value of 2, then $F^C(x)$ is proportional to $1/x$. Now, $F^C(x)$ is a linear transformation of the rank. Therefore, it means that the rank of an entity is proportional to the size of the entity. This is the reason that it is known as the "Rank-Size Distribution". This principle will be illustrated in the different contexts.

Gibrat Law and Zipf's Law

Zipf's law is applicable in a static context. The dynamic context idea related to Zipf's law is called Gibrat's Law. Gibrat's law postulates that the mean and variance of the growth rate of an entity is independent of its size.

Zipf's law in a particular time is observed if Gibrat's law is observed over time. (31) Essentially Gibrat's Law is a sufficient condition for observation of Zipf's law but not a necessary condition. It is quite possible that one can observe Zipf's law at one point in time but that may be an aberration rather than being a rule. We can confirm our observation of Zipf's Law as a usual phenomenon for a variable rather than an aberration if we can observe both Zipf's law and Gibrat's law in a variable under investigation.

The other name of Zipf's law is Power Law and a distribution adhering to the Zipf's law is called a Power Law distribution. This is because of the fact that a power law is independent of the memory as is the case for Zipf's law.

Non-Validity in a Subset

Another aspect of Zipf's law is that it is valid for a complete system but usually not for a subset of the system if the subset is a random subsample of the system. For example, the income and wealth of all agents in an economic system may follow a Power Law distribution but the income and wealth of the women in that economic system may not be a Pareto distribution. This also means that Zipf's Law may also help the researcher identify a system.

Tsallis Distribution

Constantino Tsallis, proposed a generalised version of the Pareto distribution, which being named after him is called the Tsallis q-exponential distribution (21,22). For this distribution, the probability density function, the CDF and the CCDF are given as,

$$f(x) = \frac{\theta}{\sigma} \left(1 + \frac{x}{\sigma}\right)^{-\theta-1} ; F(x) = 1 - \left(1 + \frac{x}{\sigma}\right)^{-\theta} ; F^C(x) = \left(1 + \frac{x}{\sigma}\right)^{-\theta} .$$

One may note that for large values of x, the Tsallis q-Exponential becomes close to the Pareto Distribution. For the purpose of comparability between two distributions, one may assume θ as equal to $\alpha-1$ and σ as a parameter instead of x_{\min} .

Statistical Methods

For estimating the coefficient of the power law, α , different statistical methodologies are employed. The basic way to estimate the parameters of Pareto distribution is called the “Linear Fit method”. We note that

$$F^C(x) = C \times x^{-\alpha+1}$$

$$\log F^C(x) = \log C + (1-\alpha) \log x = c + (1-\alpha) \log x$$

Under this method, one can plot (29) the log of rank of an entity plotted against the log of its size. The coefficient of the regression line yields the estimate $(1-\alpha)$ from which we can derive the exponent of power law, α . Though very simple, this estimate is biased in the sense that on average the estimated coefficient is lower than actual estimate. (23) In other words, suppose we simulate a data from a known value of α and estimate the value of α using the linear fit method from the simulated data, the estimate will, on average, be lower than the actual value.

For unbiased estimation, one can use the “Hill method” or the Maximum Likelihood Estimator (MLE). For a finite sample, one can calculate the probability of observing a particular data-point in that sample by using the assumed probability functions underlying the sample, say probability density function and CDF. Using independence of the sample data points, the probability of the entire sample can be computed. This computed probability depends on the parameters present in the particular probability distribution which was assumed for calculation of the sample probability. This probability is also known as the likelihood of the sample which is to be maximised by choice of the values for the distribution parameters. The set of values for the parameters for which the likelihood is maximized, is called the MLE. (20)

The asymptotic variance matrix of the MLE is given by the Fisher's Information matrix (24). Asymptotically, the standard error of the MLE is the smallest among all estimates for a parameter. This is the reason, the MLE is called the efficient estimate in the statistical parlance. The standard error can also be estimated by the technique of Bootstrapping. In this method, we draw sub-samples from our original sample and compute the maximum likelihood estimates for those sub-samples. The standard error in the estimates obtained from different sub-samples is our estimate for the standard error of the MLE.

Distribution of Economic Variables: Income, Wealth and Consumption

Wealth and Income Distribution

Income distributions in many countries are documented but the critical problem is modeling the upper tail regarding which observations are scarce. Moreover, availability of the reliable income and wealth data is an additional concern in India. The studies from India should be seen from this general backdrop. Last but not least, the absence of a long reliable all India micro panel data, empirical verification of Gibrat law is not a feasible exercise.

Sitabhra Sinha (33) analysed the upper tail of the Indian income distribution. He gathered data for the 125 wealthiest individuals and households in India were obtained from a special report by the Indian business magazine, Business Standard, on Indian Billionaires. The wealth was reported on two dates, Dec 31, 2002 and Aug 31, 2003. The same data set also reported the gross salary of the 67 highest-paid executives in India including the foreign nationals. Another list of 40 richest Indians published by the international business magazine Forbes in Dec 10, 2004 was also used in the study. The latter list was based on Indian nationality rather than the residency status.

The two figures 1 and 2, describe the incidence of power law among the extreme upper end of the wealth and income distribution of India, for both cases when the term was defined geographically as well as politically to cover the residents. The rank-size plot indicates linearity at the right tail with a value of α ranging between 1.66 and 2.23.

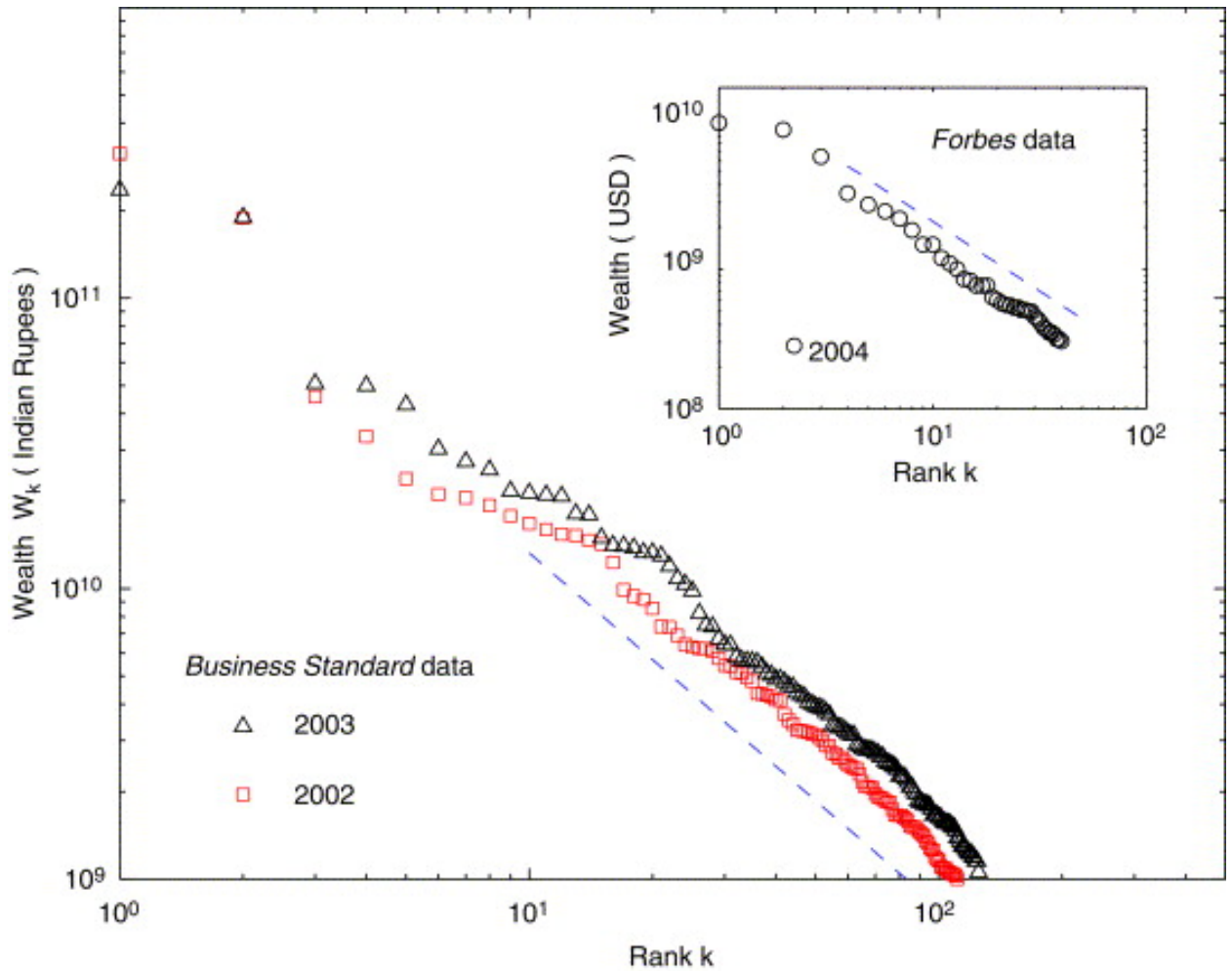


Figure 1: Rank ordered plots of the wealth of the richest Indians during the period 2002–2004 on a double-logarithmic scale. The main figure shows the wealth of the k th ranked richest person (or household) against the rank k (with rank 1 corresponding to the wealthiest person) as per two surveys conducted by Business Standard in Dec 31, 2002 (squares) and Aug 31, 2003 (triangles). The broken line having a slope of -1.23 is shown for visual reference. The inset shows the rank ordered plot of wealth based on data published by Forbes in Dec 10, 2004, with the broken line having a slope of -1.08. Source: Sinha (33)

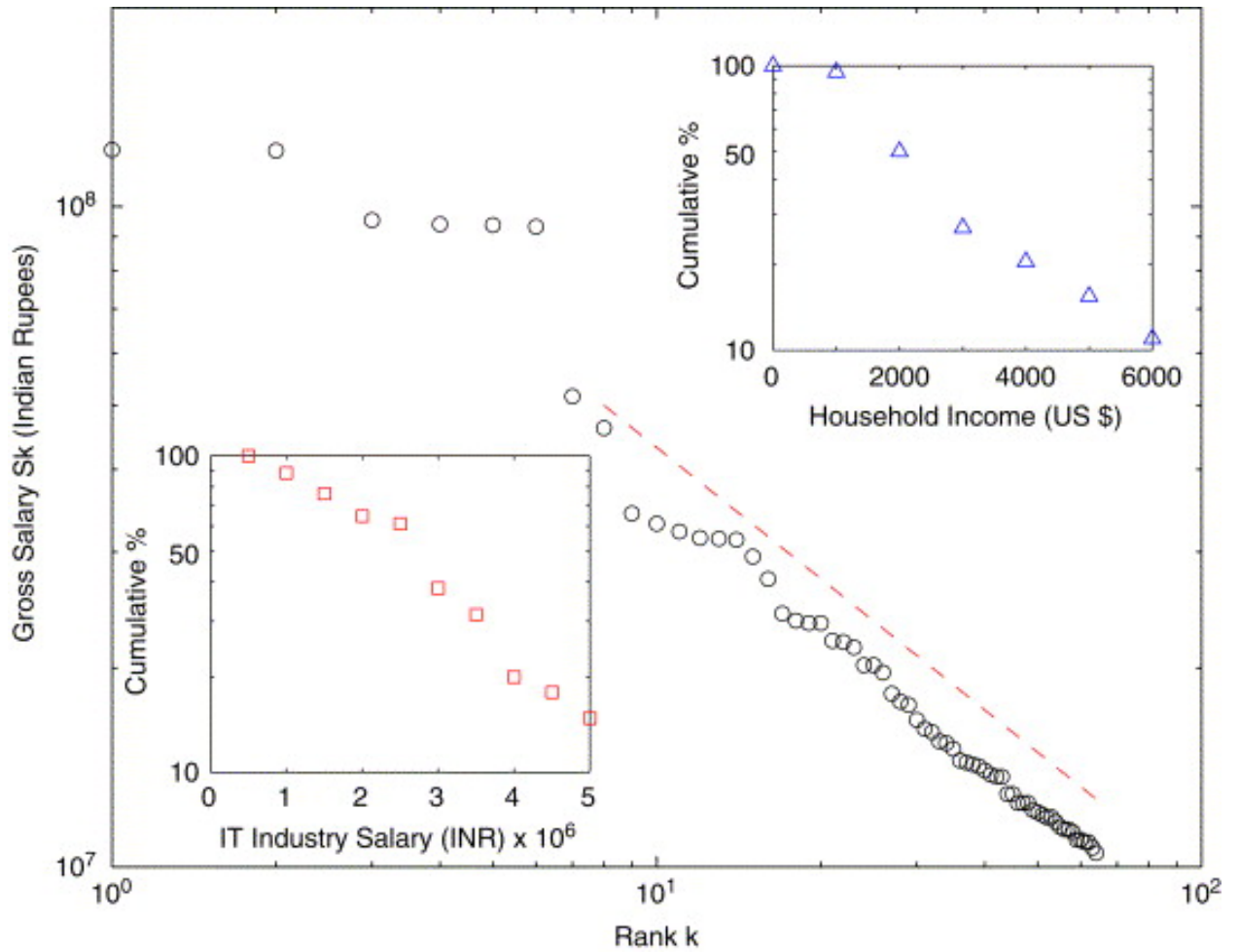


Figure 2: The rank ordered plot of the gross salary (in Indian Rupees) of the k th ranked highest paid executive against the rank k on a double-logarithmic scale. The broken line of slope -0.66 is shown for visual reference. The upper inset shows, on a semi-logarithmic scale, the cumulative percentage of Indian households at income level I (i.e., the percentage with household income greater than I) plotted against I (in US Dollars; 1 US Dollar ≈ 37 Indian Rupees during this period), for the lower-end of the income distribution. The lower inset shows, on a semi-logarithmic scale, the percentage of individuals in the Information Technology industry with 10 years or more experience, having a salary S or more (in Indian Rupees). Source: Sinha (33)

Jayadev (34) analysed All India Debt and Investment Survey (in 1991–1992 and 2002–2003) data collected by the National Sample Survey Organization (NSSO). The focus on a standard measure of wealth: total household assets which is defined as the total household assets comprising of “physical assets like land, buildings, livestock, agricultural machinery and implements, non-farm business equipment, all transport equipment, durable household goods and financial assets like dues receivable on loans advanced in cash or in kind, shares in companies and cooperative societies, banks, etc., national saving certificates and the like, deposits in companies, banks, post offices and with individuals”. The values of the Pareto exponent α varies between 1.85 and 2.45.

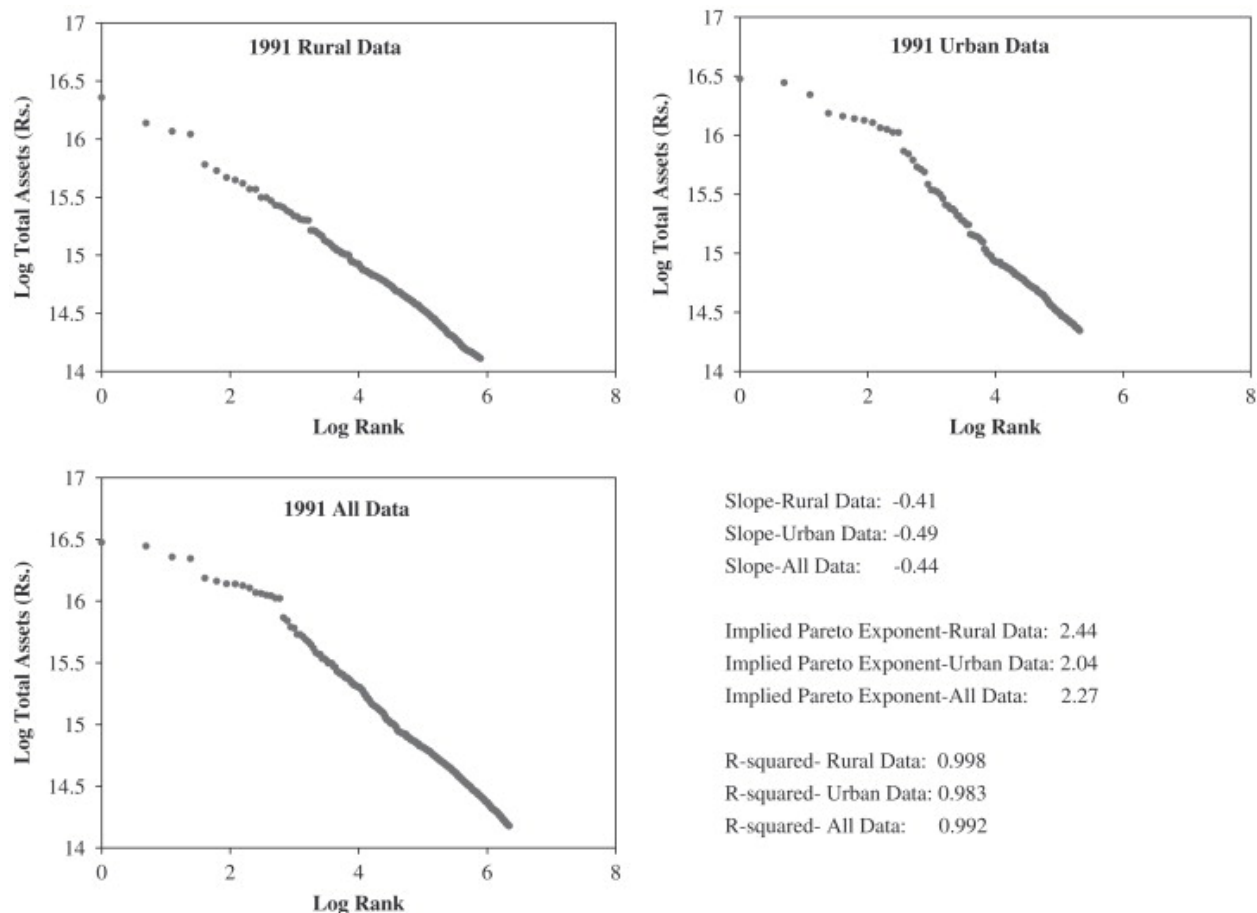


Figure 3: Rank ordered plots of the wealth of the top 1% of surveyed households from the All India Debt and Investment Survey in 1991 on a double natural logarithmic scale. The top left panel is for rural households only, the top right for urban households only, and the bottom left for both rural and urban households together. Source: Jayadev (34)

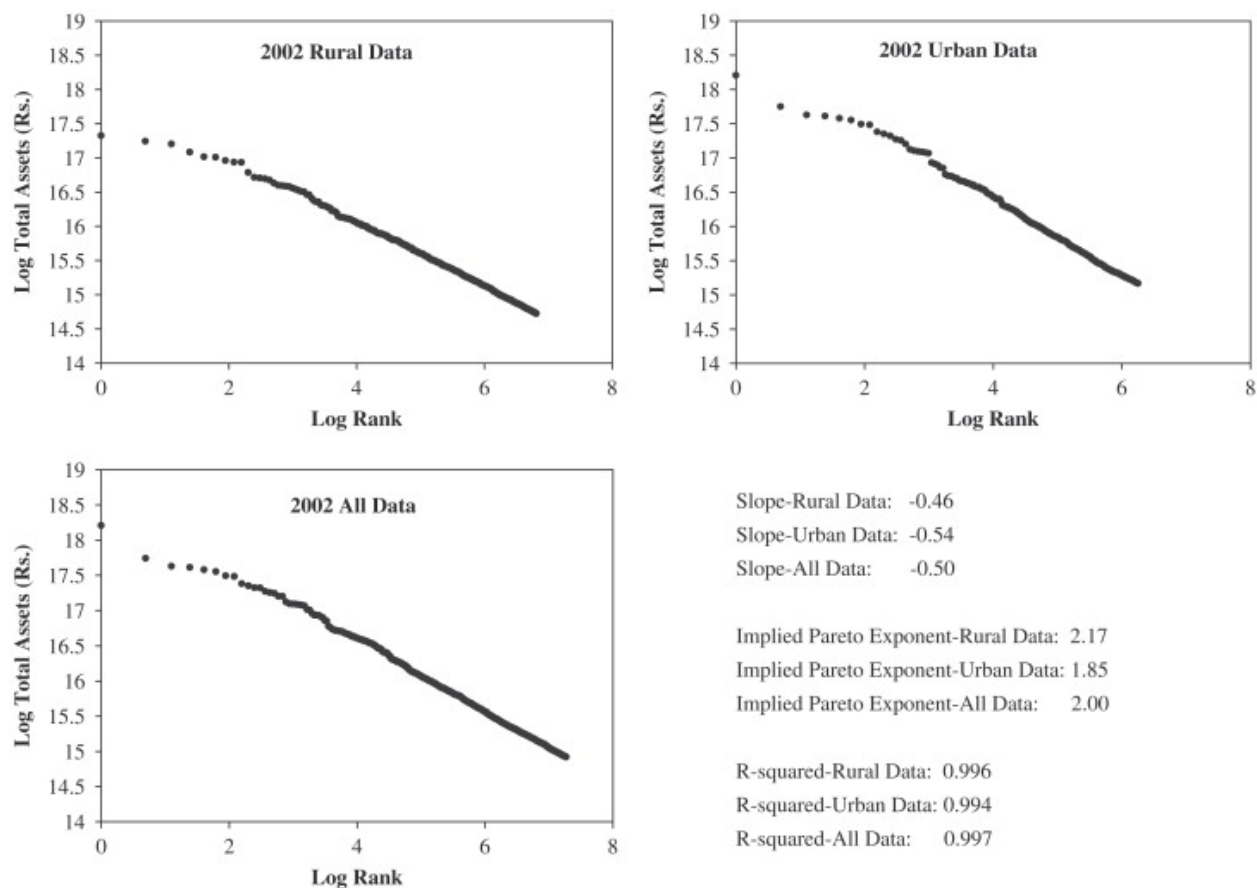


Figure 4: Rank ordered plots of the wealth of the top 1% of surveyed households from the All India Debt and Investment Survey in 2002 on a double natural logarithmic scale. The top left panel is for rural households only, the top right for urban households only, and the bottom left for both rural and urban households together. Source: Jayadev (34)

Consumption Distribution

Consumption data in every year is collected by NSSO and usually in every fifth year, the sample size is much more, called quinquennial rounds. We (13) gather our data consisting of expenditure for individuals and families grouped in different “MPCE” classes along with the average expenditure in each class from the NSSO rounds during 1983 to 2007. We modelled it as a mixture of lognormal distribution and a Pareto upper tail.

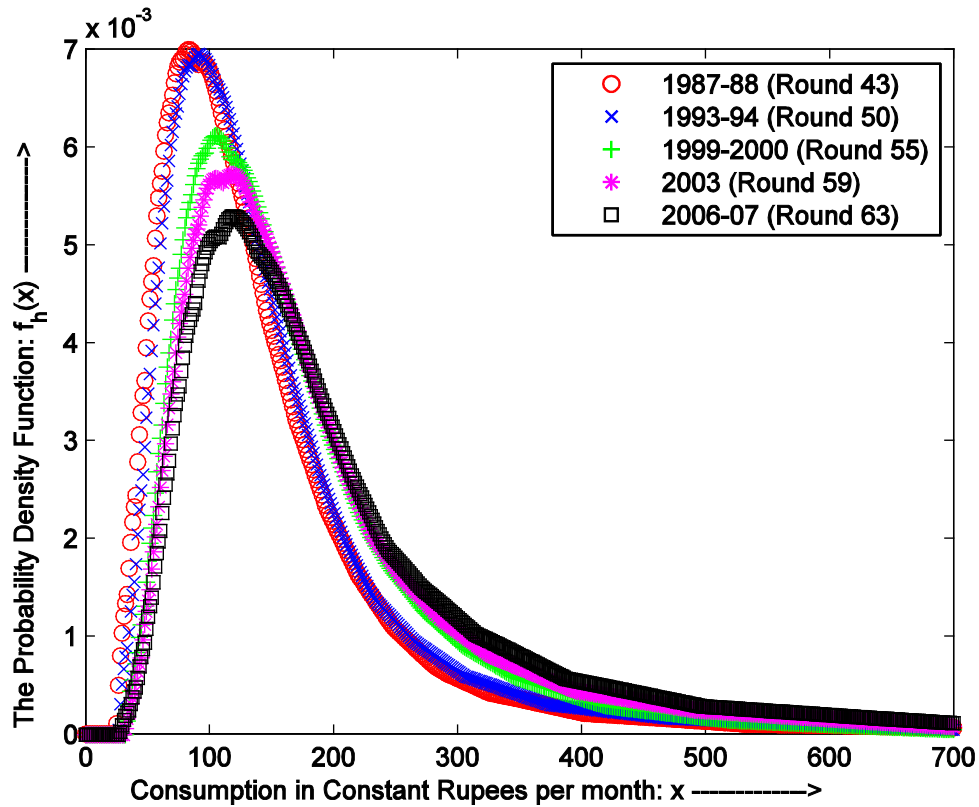


Figure 5: Kernel Density Estimate for the Expenditure Distribution in India plotted in linear scale: 1987-2007 (for Households) using the NSS data for Rural India.

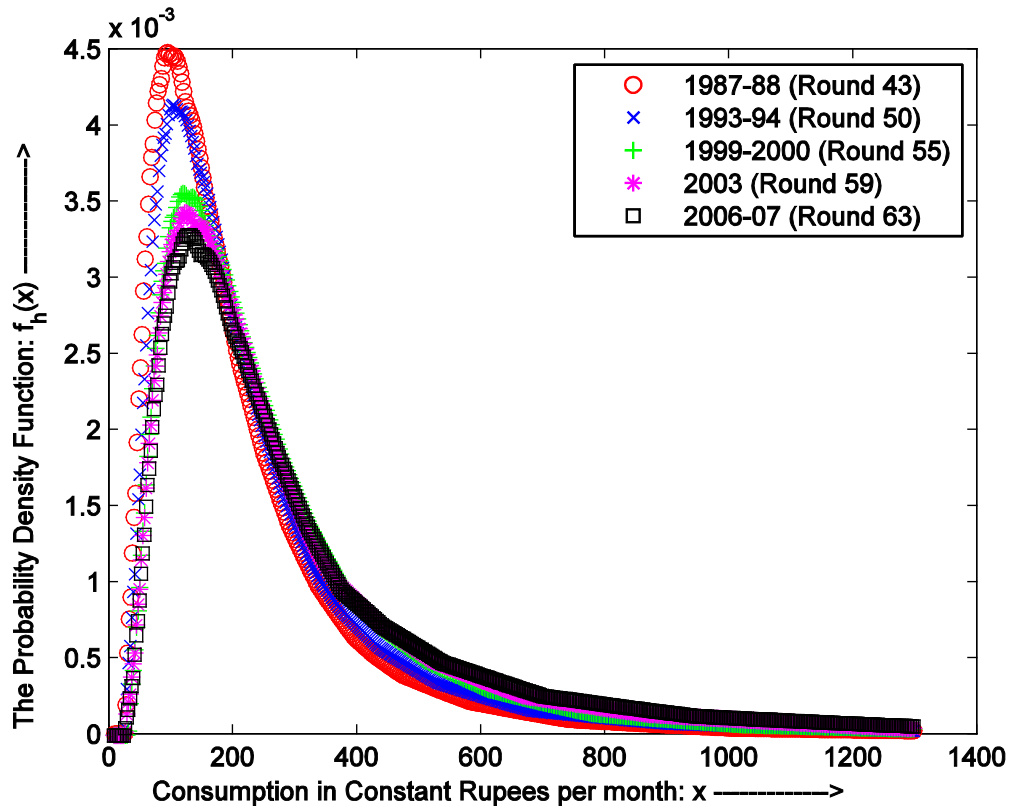


Figure 6: Kernel Density Estimate for the Expenditure Distribution in India plotted in linear scale: 1987-2007 (for Households) using the NSS data for Urban India

We estimated the extent of the Pareto tail both at the household level and at the person level. The tail varies widely between rural to urban samples and also across different years (See Table 1). Summarily speaking, the Pareto tail can be as low as 2.8% and as high as 35.4%. Some stylized facts are also found to be true such as urban inequality is more than rural inequality, marked by the presence of a higher Pareto tail. Another study (28) also confirmed presence of Pareto tail in consumption distribution.

The main limitation of our study is arising out of working with a grouped data. As a consequence, the variations are somewhat higher. The other problems include non-sampling errors in the NSS data as highlighted by many studies (25, 27).

City Size Distribution

Let us take a cut-off for the city population, say fourteen million. Three Indian metropolises existed in 2011, the year of the latest census, over this cut-off population limit: Greater Mumbai (18,414,288), Delhi (16,314,838) and Kolkata (14,112,536), the figures inside the parenthesis denote the population of the respective city. The CCDF of the empirical distribution of the Indian city sizes at the size of fourteen million is $3/(n+1)$, in which n denotes the number of cities in the sample. Since $F^C(x)$ is proportional to $1/x$, on halving the size, the CCDF should be doubled to $6/(n+1)$. Therefore, there should be six cities over the size of seven million. Interestingly, we find that *exactly* six Indian cities have a population over seven million. Apart from Mumbai, Delhi and Kolkata, the other three cities are Chennai (8,696,010), Bangalore (8,499,399) and Hyderabad (7,749,334), again the figures inside the parenthesis denote the population for each city.

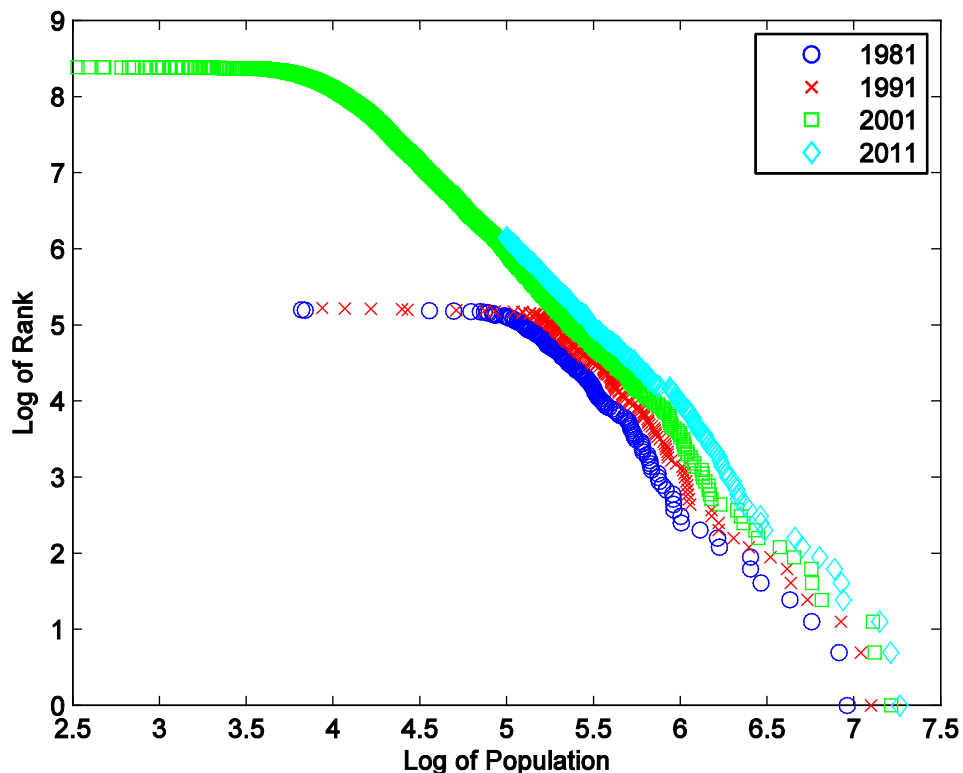


Figure 7: City sizes have been plotted against their ranks, in the logarithmic scale, 1981—2011. Source: Author

The data are taken from the Indian census that is carried out once in a decade. Figure 7 plots the Indian city sizes against their respective ranks, in logarithmic scale, for four census years, 1981, 1991, 2001 and 2011. One cannot fail to notice the linearity in the upper tail of the distribution. Indeed the minimum size may be changing for all the different census years which we have estimated (20, 26, 32). Table 2 presents the estimates over time. They are quite stable and Statistically indistinguishable from the value of 2. The estimates of the Tsallis distribution (20) are no different, even quantitatively speaking.

For examining the stability of the Zipf's law, validity of Gibrat's Law is important. The scatter plots of growth rate against size. Figures 8 and 9, is not demonstrative of any size dependency. The mean and variance of the growth rates contingent on size were calculated using the method of Kernel density regression, a non-parametric way of computing statistical highlights. (32) Figures 10 and 11 illustrate the mean growth rate against city sizes; whereas Figures 12 and 13 illustrate the growth rate of variance against city sizes. All are suggestive of empirical validity of the Gibrat's Law.

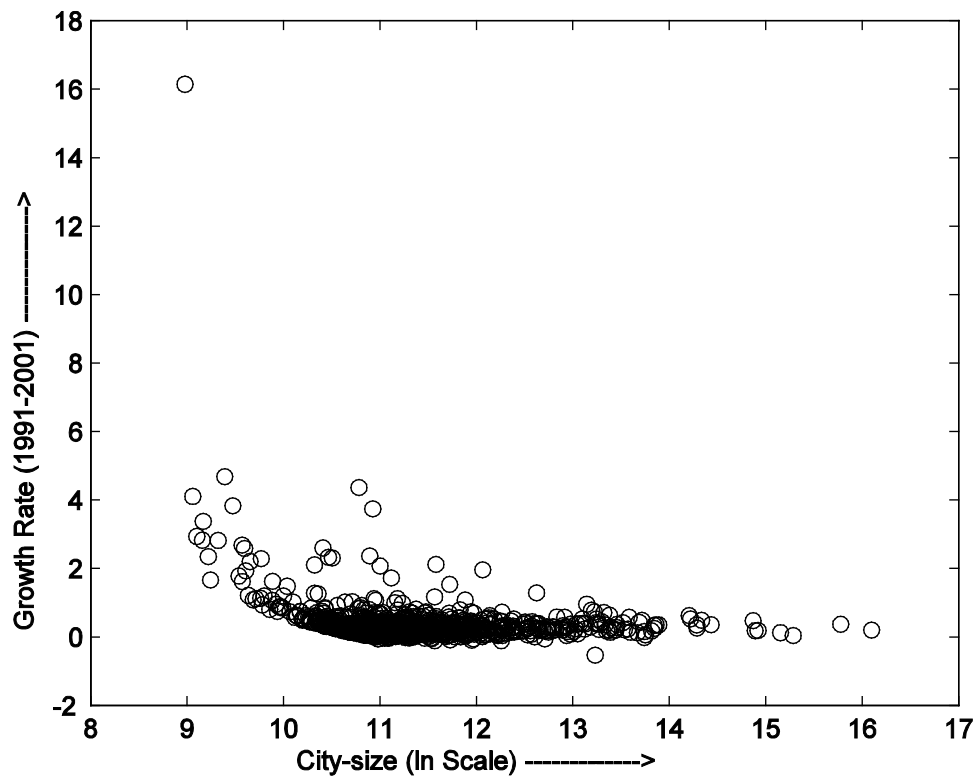


Figure 8: Scatter Plot for Growth Rate of Urban Agglomerations against Size: 1991–2001. Source: Author

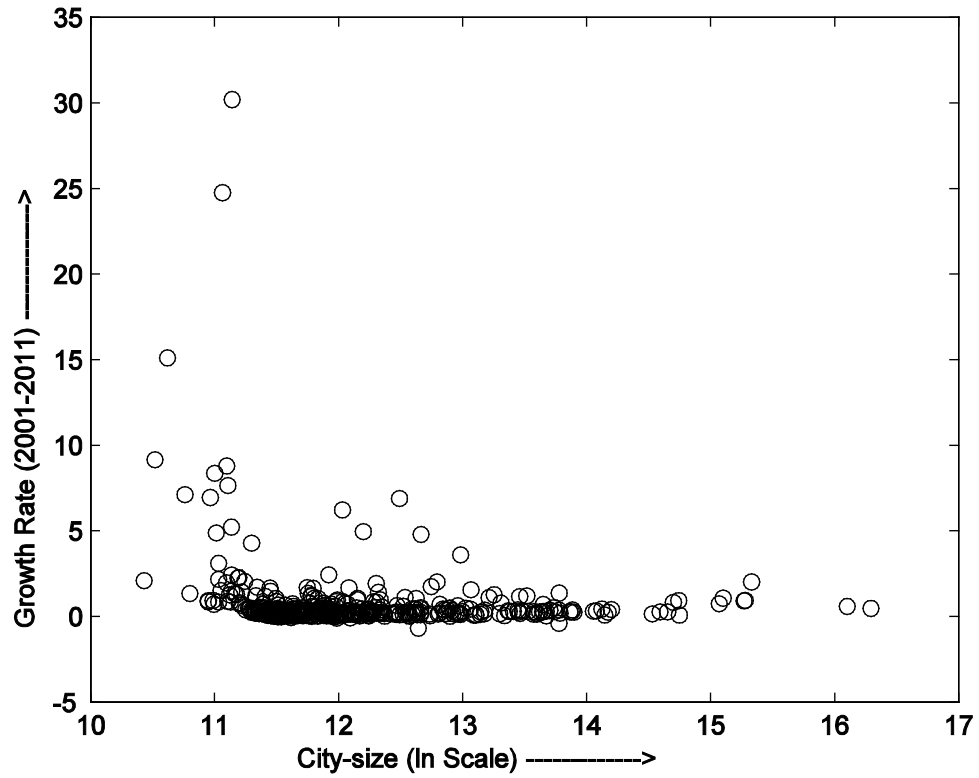


Figure 9: Scatter Plot for Growth Rate of Urban Agglomerations against Size: 2001–2011. Source: Author

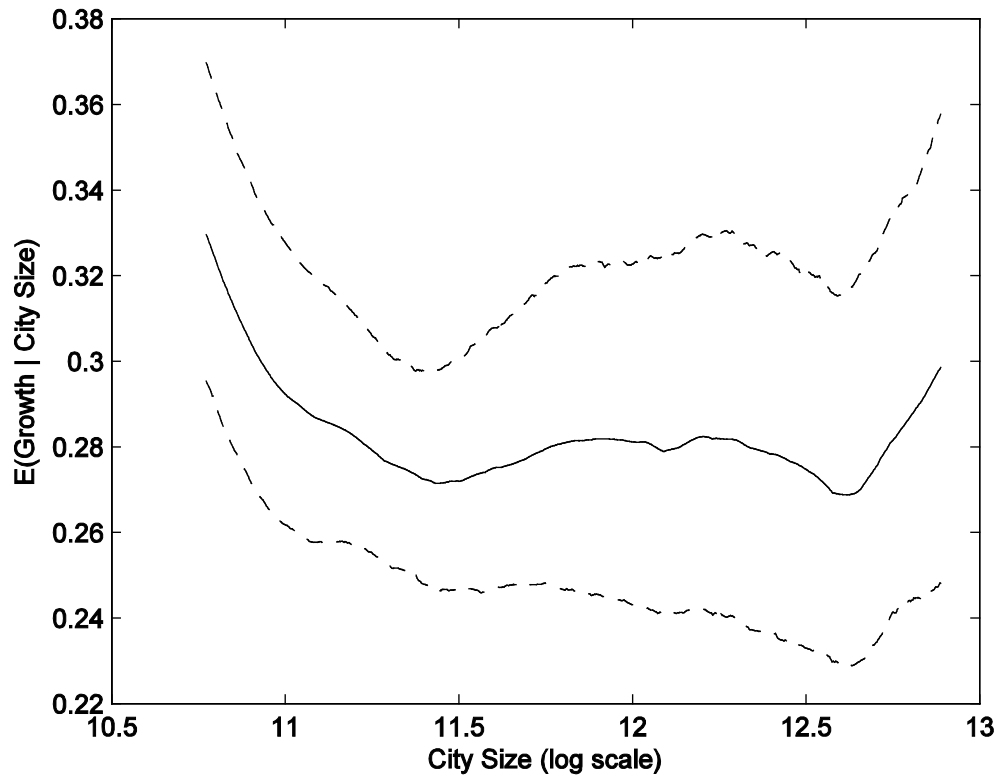


Figure 10: Mean Growth Rate of Urban Agglomerations against Size: 1991–2001.

Source: Author

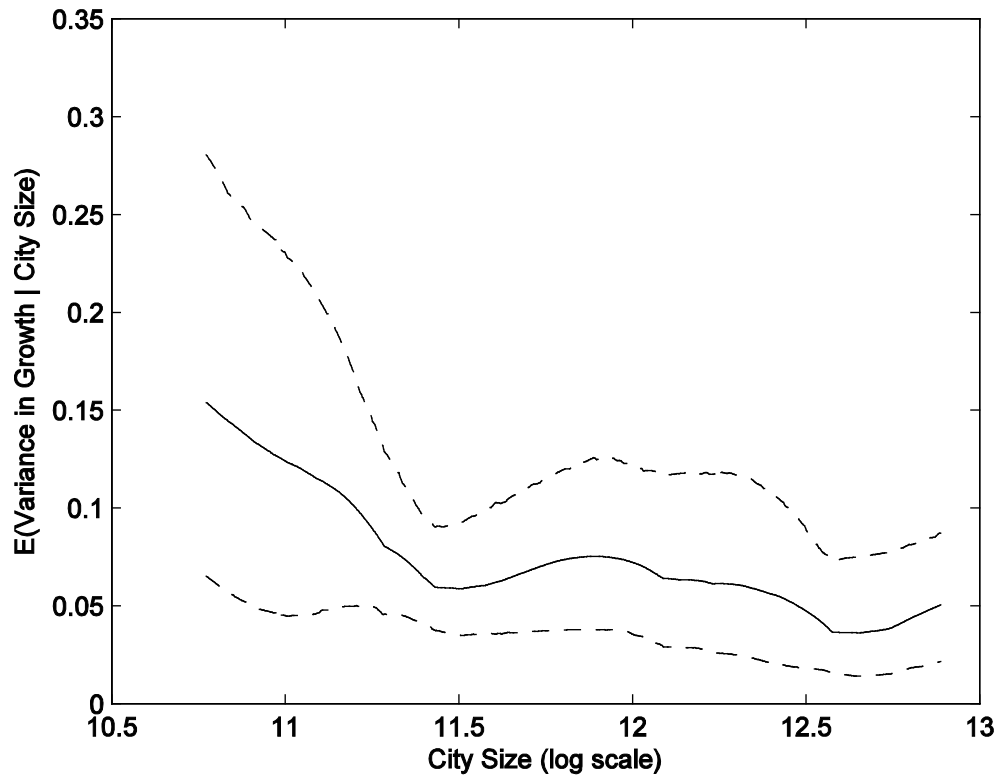


Figure 11: Mean Growth Rate of Urban Agglomerations against Size: 2001–2011.

Source: Author

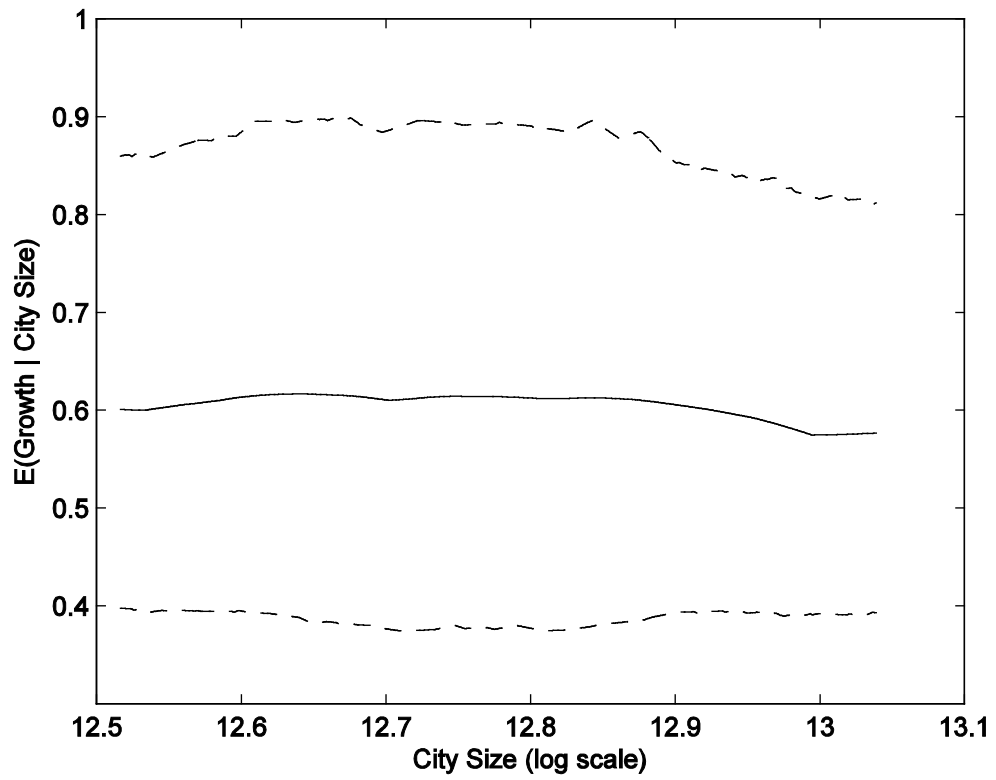


Figure 12: Growth Rate Variance of Urban Agglomerations against Size: 1991–2001.

Source: Author

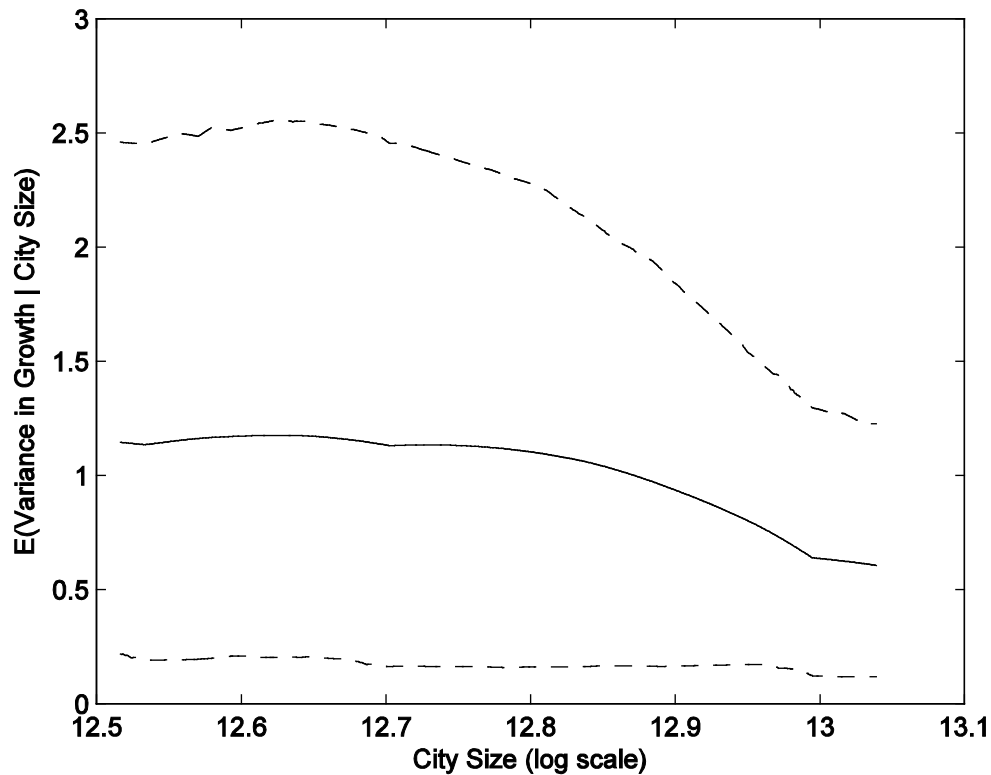


Figure 13: Growth Rate Variance of Urban Agglomerations against Size: 2001–2011.

Source: Author

Robustness Checks

Instead of x_{\min} being an exogenous parameter in estimating the scaling exponent, the following algorithm is useful to derive an endogenous values of the x_{\min} . The distance between the two CDFs, the empirical one obtained from the data and the other one arising out of the best-power law by choice of the proper value for x_{\min} . Kolmogorov-Smirnov (KS) statistic is a popular measure to assess the distance between the two probability distributions, which computes the distance between two statistical distributions with CDFs F_1 and F_2 as the supremum of the absolute value of the difference between each these two CDFs at all points of the support. The other robustness exercises include estimation of the Tsallis distribution as opposed to simple Pareto distribution.

Conclusion

Zipf's Law is not a law in the truest sense of the term but a mathematical formulation that is prominent in quantitative description of the human behavior across countries across different time periods. The underlying social conditions differ from country to country from one time to another. However, those fluctuations are not strong enough to mitigate the phenomenon of Zipf's law. Perhaps the impact of those fluctuations is captured by the Pareto exponent. Paul Krugman, the noted economist, noted this powerful implication of the Zipf's law: "the usual complaint about economic theory is that our models are oversimplified — that they offer excessively neat views of complex, messy reality. In the case of Zipf's law the reverse is true: we have complex, messy models, yet reality is startlingly neat and simple."(17)

We have surveyed studies which demonstrated Zipf's law being followed in the distribution of socio-economic variables. By no means, this is an exhaustive list of variables for which Pareto exponent have been estimated. However, many other variables such as the financial market variables may not have much implication on the socio-economic condition of the country and are not central to socio-economic context. The value of the Pareto exponent differs but assumes values which are not unusual in the particular context.

References:

1. Zipf George Kingsley. Severe depression in old age. Human Behaviour and the principle of Least Effort. Cambridge, MA: Addison-Wesley; 1949.
2. Zipf George Kingsley. The psycho-biology of language. Oxford, England: Houghton-Mifflin; 1935.
3. Sasaki Y, Kuninaka H, Kobayashi N, Matsushita M. Characteristics of population distributions in municipalities. Journal of the Physical Society of Japan. 2007 Jun 25;76(7):074801.
4. Huberman BA, Pirolli PL, Pitkow JE, Lukose RM. Strong regularities in world wide web surfing. Science. 1998 Apr 3;280(5360):95-7.
5. Yakovenko VM, Rosser Jr JB. Colloquium: Statistical mechanics of money, wealth, and income. Reviews of Modern Physics. 2009 Dec 2;81(4):1703.
6. Chatterjee A, Chakrabarti BK. Kinetic exchange models for income and wealth distributions. The European Physical Journal B. 2007 Nov 1;60(2):135-49.

7. Baldwin RB. Lunar crater counts. *The Astronomical Journal*. 1964 Jun;69:377.
8. Malamud BD, Morein G, Turcotte DL. Forest fires: an example of self-organized critical behavior. *Science*. 1998 Sep 18;281(5384):1840-2.
9. Boffetta G, Carbone V, Giuliani P, Veltri P, Vulpiani A. Power laws in solar flares: self-organized criticality or turbulence?. *Physical review letters*. 1999 Nov 29;83(22):4662.
10. L.C. Malacarne, R.S. Mendes, *Physica A* 286 (2000) 391
11. Blasuis Bernd, Tonjes Ralf. Zipf's Law in the Popularity Distribution of Chess Openings. *Physical Review Letters* 2009 Nov 16; 103(218701).
12. Mizuno T, Toriyama M, Terano T, Takayasu M. Pareto law of the expenditure of a person in convenience stores. *Physica A: Statistical mechanics and its applications*. 2008 Jun 15;387(15):3931-5.
13. Ghosh A, Gangopadhyay K, Basu B. Consumer expenditure distribution in India, 1983–2007: Evidence of a long Pareto tail. *Physica A: Statistical Mechanics and its Applications*. 2011 Jan 1;390(1):83-97.
14. Zanette DH, Manrubia SC. Role of intermittency in urban development: a model of large-scale city formation. *Physical Review Letters*. 1997 Jul 21;79(3):523.
15. Moura NJ, Ribeiro MB. Zipf law for Brazilian cities. *Physica A: Statistical Mechanics and its Applications*. 2006 Jul 15;367:441-8.
16. Gabaix X, Ioannides YM. The evolution of city size distributions. *Handbook of regional and urban economics*. 2004 Dec 31;4:2341-78.
17. Krugman P. *The self-organizing economy*. Oxford, UK and Cambridge, MA: Blackwell Publishers; 1996 Feb.
18. Gabaix X. Zipf's law for cities: an explanation. *Quarterly journal of Economics*. 1999 Aug 1:739-67.
19. Soo KT. Zipf's Law for cities: a cross-country investigation. *Regional science and urban Economics*. 2005 May 31;35(3):239-63.
20. Gangopadhyay K, Basu B. City size distributions for India and China. *Physica A: Statistical Mechanics and its Applications*. 2009 Jul 1;388(13):2682-8.
21. Tsallis C. Possible generalization of Boltzmann-Gibbs statistics. *Journal of statistical physics*. 1988 Jul 1;52(1-2):479-87.

22. Malacarne LC, Mendes RS, Lenzi EK. q-Exponential distribution in urban agglomeration. *Physical Review E*. 2001 Dec 21;65(1):017106.
23. Clauset A, Shalizi CR, Newman ME. Power-law distributions in empirical data. *SIAM review*. 2009 Nov 6;51(4):661-703.
24. Rao CR. *Linear statistical inference and its applications*. John Wiley & Sons; 2009 Sep 25.
25. Vaidyanathan A. On the validity of NSS consumption data. *Economic and political Weekly*. 1986 Jan 18:129-37.
26. Gangopadhyay K, Basu B. The Morphology of Urban Agglomerations for Developing Countries: A Case Study with China. In *Econophysics and Economics of Games, Social Choices and Quantitative Techniques 2010* (pp. 90-97). Springer Milan.
27. Minhas BS. Validation of large scale sample survey data Case of NSS estimates of household consumption expenditure. *Sankhyā: The Indian Journal of Statistics, Series B*. 1988 Dec 1:279-326.
28. Chakrabarti AS, Chatterjee A, Nandi TK, Ghosh A, Chakraborti A. Quantifying invariant features of within-group inequality in consumption across groups. Available at SSRN 2727616. 2016 Feb 4.
29. Basu B, Bandyopadhyay S. Zipf's law and distribution of population in Indian cities. *Indian Journal of Physics*. 2009 Nov 1;83(11):1575-82.
30. Gangopadhyay Kausik. Interview with Eugene H. Stanley. *IIMK Society & Management Review* 2013; 2(2): 73-78.
31. Malamud BD, Morein G, Turcotte DL. Forest fires: an example of self-organized critical behavior. *Science*. 1998 Sep 18;281(5384):1840-2.
32. Gangopadhyay K, Basu B. Evolution of Zipf's Law for Indian Urban Agglomerations Vis-à-Vis Chinese Urban Agglomerations. In *Econophysics of Systemic Risk and Network Dynamics 2013* (pp. 119-129). Springer Milan.
33. Sinha S. Evidence for power-law tail of the wealth distribution in India. *Physica A: Statistical Mechanics and its Applications*. 2006 Jan 1;359:555-62.
34. Jayadev A. A power law tail in India's wealth distribution: Evidence from survey data. *Physica A: Statistical Mechanics and its Applications*. 2008 Jan 1;387(1):270-6.

Table 1: Estimates of Power Law Exponents for Consumption during 1983 to 2007

Sample	Range of Pareto Coefficient α
Rural Households	5.3–20.0%
Urban Households	8.7–35.4%
Rural Persons	2.8–18.5%
Urban Persons	7.1–27.7%

Table 2: Estimates of Power Law Exponents for City Size

Year	Minimum Size	Linear Fit estimate of α	Maximum Likelihood Estimate of α
2011	212,523	1.935 (0.007)	2.018 (0.069)
2001	180,355	1.921 (0.007)	2.044 (0.075)
1991	148,272	1.899 (0.008)	1.976 (0.075)
1981	120,000	1.889 (0.009)	1.991 (0.080)

Research Office

Indian Institute of Management Kozhikode

IIMK Campus P. O.,

Kozhikode, Kerala, India,

PIN - 673 570

Phone: +91-495-2809238

Email: research@iimk.ac.in

Web: <https://iimk.ac.in/faculty/publicationmenu.php>

