

ABSTRACT

Title of Dissertation: ADAPTIVE AGENT MODELING IN A
POLICY CONTEXT

Timothy R. Gulden, Doctor of Philosophy, 2004

Dissertation Directed By: Professor Herman E. Daly
School of Public Policy

This dissertation examines the utility of adaptive agent modeling (also referred to as agent-based modeling or individual based modeling) as a tool in public policy research. It uses the adaptive agent technique to produce useful results in three diverse areas.

It demonstrates that the adaptive agent framework can be used to extend traditional models of comparative advantage in international trade, showing that the presence of increasing returns to scale in some industries shifts the basis of comparative advantage arguments, making room for industrial policy and the regulation of trade.

Next, the dissertation demonstrates that the size distribution of cities within nations, generally thought to approximate the “Zipf” distribution, can be reproduced using a simple agent-based model. This model produces insights into the evolution of the distribution as well as departures from it – especially in France and Russia. This understanding of urban dynamics has implications for easing the structural transition of the Russian economy and for designing policies to reduce the size of megacities in the developing world.

The dissertation goes on to examine individual level data from the Guatemalan civil war from an adaptive agent modeling perspective. It finds several novel patterns in the data which may serve as benchmarks for adaptive agent modeling efforts and suggests avenues by which existing conflict models might be brought into closer accord with the data.

The dissertation concludes that adaptive agent modeling is useful in a policy context because it allows quantitative work to be done while relaxing some of the unrealistic assumptions which are often required to gain analytical traction using traditional methods. The method is found to be particularly useful in situations where path dependence, heterogeneity of actors, bounded rationality, and imperfect information are significant features of the system under examination. The individual based nature of the method is also found to be well suited to assessing distributional impacts of changes in process or policy.

ADAPTIVE AGENT MODELING IN A POLICY CONTEXT

By

Timothy R. Gulden

Dissertation submitted to the Faculty of the Graduate School of the
University of Maryland, College Park, in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
2004

Advisory Committee:
Professor Herman E. Daly, Chair
Professor Robert Axtell
Professor Catherine Dibble
Professor Thomas C. Schelling
Professor John D. Steinbruner
Professor Charles Christian, Dean's Representative

© Copyright by
Timothy R. Gulden
2004

Acknowledgements

I would like to begin by thanking the members of my committee including my advisor, Herman Daly, School of Public Policy faculty members John Steinbruner and Thomas Schelling, Brookings Senior Fellow Robert Axtell, Geography faculty member Catherine Dibble and Dean's Representative Charles Christian. They are a remarkable collection of people who have taught me a great deal and have been extraordinarily generous to me over the past five years. Their support and guidance, along with a fellowship provided by the Charles Stuart Mott foundation, have made this work possible.

I am also grateful to the Brookings Institution's Center on Social and Economic Dynamics (CSED) which has funded parts of this work and provided me with a uniquely stimulating environment in which to explore these ideas. I am indebted to Dr. Steinbruner for putting me in touch with the Center. Carol Graham has helped me to focus both the conflict and cities chapters and has inspired me to push these trains of thought forward when I did not know how to proceed. Robert Axtell and Joshua Epstein have been selfless with sharing their insights into agent modeling and have provided powerful examples of what the method can do. Clifford Gaddy and Fiona Hill provided me with a great deal of information on the political and economic framework in which the Russian city size distribution evolved and have spent many hours helping me think about internal migration in the Soviet and post-Soviet eras. Miles Parker, Ross Hammond, Shubha Chakravarty and John Parker provided technical, intellectual, and moral support at various points.

I want, particularly, to acknowledge Ross Hammond's contribution to the cities chapter. He and I developed the "jars and beans" model jointly in the spring of 2001 in an attempt to get a better grip on the origins of power-law distributions. Most of the work to extend the model and apply it to cities is my own, but I expect that it will be improved through future collaboration with Ross on the subject.

I am similarly indebted to Anindya Sen with whom I developed a first version of the international trade model as a project at the Santa Fe Institute Complex Systems Summer School in 2002. Though the model has evolved significantly since that time, his work in helping me frame the problems was invaluable.

Thanks are also due to Dr. Robert Hunt Sprinkle who provided a great deal of good advice (much of it invited) on matters academic. He also provided editorial guidance and oversaw the publication of "Spatial and Temporal Patterns in Civil Violence: Guatemala, 1977-1986" in *Politics and the Life Sciences* in 2002.

Egor Kraev provided much useful mathematical advice and helped me to clarify many of the ideas presented here.

I am deeply thankful for the unwavering and multi-faceted support of my parents, of Bruce and Donna Wilshire, and of my wife Sarah. Finally, I acknowledge the contribution of my new daughter Isabel, who provided the ultimate motivation to finish the project.

Table of Contents

Acknowledgements.....	ii
Table of Contents.....	iv
List of Figures.....	vi
Chapter 1: Adaptive Agents as a Tool for Policy Research.....	1
Types of Numerical Models.....	1
Systems Dynamics Models.....	2
Adaptive Agent Modeling as a Form of Systems Dynamics Modeling.....	4
Uses of Adaptive Agent Modeling.....	7
Adaptive Agent Modeling in a Policy Context.....	9
Chapter 2: The Importance of Assumptions: Adaptive Agent Modeling as a Tool for Trade and Development Theory.....	12
Gomory and Baumol’s Model of International Trade.....	13
A Place for Policy.....	13
Multiple Equilibria.....	16
Cooperation and Conflict.....	18
Revisiting the Infant Industries Argument.....	20
Policy Implications.....	22
An Adaptive Agent Model of International Trade.....	22
Model Specification.....	23
Samuelson’s Analysis of Outsourcing.....	26
Verifying the Agent Model.....	29
Verifying Gomory and Baumol’s Retainable Industries.....	32
Discussion.....	37
Next Steps.....	38
Chapter 3: Beyond Zipf: An Agent Based Understanding of City Size Distributions.....	44
The Zipf Distribution.....	45
Deviations from Zipf.....	47
USA.....	48
France.....	51
Russia.....	55
A Simple, Abstract Model: Jars and Beans.....	59
Model Description.....	59
Results from the Abstract Model.....	61
Limitations of the Abstract model.....	66
A Richer Model: Cities and Citizens.....	67
Model Overview.....	67
Results from the Richer Model.....	72
Limitations.....	81
Discussion.....	82

Is Population Growth Needed to Preserve Stability of the Lower Tail?.....	82
Implications for Developing Nation Megacities.....	83
Implications for Russian Urban Structure.....	86
Next Steps.....	89
More and Better Data on Urban Agglomerations.....	89
Incorporate Florida’s Work on the “Creative Class”.....	90
Spatially Explicit Implementation.....	90
Calibrate Parameters to Historical US Data.....	91
Conclusion.....	92
 Chapter 4: Spatial and Temporal Patterns in Civil Violence: Guatemala 1977 – 1986	
.....	93
The Guatemalan Conflict.....	94
Data.....	95
Methods.....	97
Observations.....	98
Frequency vs. Severity.....	98
Ethnic Mix.....	99
Punctuated Equilibrium.....	104
Distribution of Incident Sizes.....	107
Modeling.....	116
The Brookings Model.....	117
Evaluating the Brookings Model.....	120
Conclusions.....	123
 Chapter 5: Conclusion.....	126
Contributions of the cases.....	126
Trade.....	126
Cities.....	128
Conflict.....	130
The meaning of the Zipf distribution.....	133
General implications of the dissertation for policy research.....	135
 Bibliography.....	139

List of Figures

Figure 2.1: Multiple equilibria in a world with increasing returns to scale. Reproduced from Gomory and Baumol [2000].	17
Figure 2.2: Zones of mutual gain and zone of conflict in bilateral trade. Reproduced from Gomory and Baumol [2000].	19
Figure 2.3: Adaptive Agent Realization of Samuelson Trade Model	31
Figure 2.4: Retainability of Industries with Increasing Returns	36
Figure 3.1 Zipf Distribution ordered histogram on normal and log-log axes.	45
Figure 3.2: United States Core Based Statistical Areas, 2000.	50
Figure 3.3: Four definitions of French city sizes.	54
Figure 3.4: France urban centers, 1999.	55
Figure 3.5: Distribution of Russian City Sizes	56
Figure 3.6: Curvature of the Russian city size distribution.	58
Figure 3.7: Simple model output compared with discretized US data.	63
Figure 3.8: Simple model output compared with discretized France data.	64
Figure 3.9: Simple model output compared with discretized Russia data.	65
Figure 3.10: USA model results.	74
Figure 3.11: France model results.	76
Figure 3.12: Russia model results with constant core sizes.	78
Figure 3.13: Russia model results with larger cores for Moscow and St. Petersburg.	80
Figure 4.1: Map of frequency and severity of killing, by municipality	98
Figure 4.2: Map of Ethnic Distribution in Municipalities	100
Figure 4.3: Histogram of Ethnicity and Killing	101
Figure 4.4: Time Series Graphs: Annual, Monthly, Monthly for a Single Town.	104
Figure 4.5: Killings by month, with power-spectrum and logged power-spectrum plots.	106
Figure 4.6: Rank/Size Plot of Killings per Municipality-Month.	108
Figure 4.7: Confirmed genocidal massacres (dark circles) and massacres for which genocide status was not determined by CEH (light circles), 1981-2.	109
Figure 4.8: Rank-size plots of the nongenocidal-killings subset and the genocidal- killings subset, by municipality-month.	110
Figure 4.9: Trend of the power-law exponent for different levels of aggregation in the nongenocidal subset.	112

Chapter 1: Adaptive Agents as a Tool for Policy Research

The analysis of public policy almost always involves models of some sort. Because the systems involved with real world policy problems are highly complex and often lack clear boundaries, the policy analyst must work from a simplified version of the actual system, i.e. a model. This model may be conceptual and qualitative or it may be rigorously quantitative using a host of statistical and mathematical methods. These models are useful to social science if their abstractions yield insights into the real system. They are useful to policy analysis if they yield insights into how the system might be manipulated in order to generate a socially desired result. Adaptive agent models represent a novel approach to abstracting from real systems. Such models are applicable to a different (though overlapping) set of problems than more traditional quantitative techniques and yield qualitatively different kinds of insights. The object of this dissertation is to contribute to the field's understanding of the adaptive agent approach and to identify some cases where it can be productively used.

Types of Numerical Models

Ruth and Hannon [2001] divide numerical models into three broad classes: static, comparative static, and dynamic. Each of these categories comprises a huge class of models which are suited for different tasks.

Static models seek to explain the state of a system at a single point in time. Many statistical models fall into this category. A hedonic pricing model, for example, uses the statistical technique of linear regression to explain the price of an

asset as a function of its attributes [Rosen, 1971]. A well constructed model of this sort can help a real estate assessor to estimate the value of a house given the recent sale prices of other houses in the area. The model can work even if the house in question has a unique combination of bedrooms, bathrooms, square footage, etc., because it decomposes the price into a function of these attributes. This model allows a price to be computed for combinations of attributes that do not occur in the sample set.

Comparative static models seek to understand a system by calculating its state at two or more points in time. Such models have long been a staple of economic analysis and are particularly useful when the systems that they describe have equilibria which are 1) stable, 2) unique, and 3) reachable. When these conditions are met, it is safe to assume that we will be able to find an equilibrium for a given point in time, and that this equilibrium will give us useful information about the system. The “canonical assumptions” of neoclassical economics (decreasing returns, perfect rationality, instantaneous adjustment, etc.) generally ensure that these conditions will be met, making the analysis of comparative statics a natural tool within this frame of economic reference. The standard ISLM model in macroeconomics is an example of a model which lends itself to comparative static analysis. When we relax assumptions such as decreasing returns and perfect rationality however, we can no longer assume that our models will have stable, unique, or reachable equilibria.

Systems Dynamics Models

In contrast to these various static and comparative static models, dynamic models trace the evolution of a system in time. Most often, a dynamic model

represents the system of interest as a set of differential equations. Where this system is simple, it may prove to be analytically tractable, thus allowing us to produce an equation which predicts the state of the system at any give time in the future. More often, however, a model which is rich enough to provide non-trivial insights contains non-linear terms and other complications which make analytical treatment impossible. In these cases, we must resort to numerical simulation using computers to understand the behavior of the system.

Because numerical simulation is so often needed in order to understand the behavior of a complex dynamic model, various software environments have been developed to aid in the construction and analysis of such models. Examples include Stella, Madonna, and Vensim, among others. These packages are designed to facilitate the development of systems dynamics based models, where the system in question is represented using a visual language of stocks (or state variables) which represent the state of the system in time, flows (or control variables) and transforming variables which represent constants or calculated quantities based on other variables. Once these variables are given initial values and related to one another with appropriate functional forms, the software environment uses various integration techniques to approximate the evolution of the system in continuous time.

While this approach to modeling has been widely used since Ashby [1956], some particularly influential systems dynamics models include the “World 3” model by Meadows et al. [1972], and the model used by Costanza et al. [1997] to estimate the economic value of global ecosystem services. Since the advent of inexpensive computers and user-friendly modeling software, such models have become common

tools in ecology, operations research, climate change assessment, and a host of others areas. While academic economics has been somewhat resistant to the use of numerical models which do not produce elegant analytical proofs, it has increasingly come to recognize that there are important classes of problems for which systems dynamics based dynamic simulation is a useful and necessary tool [Hannon & Ruth, 1997; Sterman, 2000].

Though traditional analytical techniques are extremely powerful for analyzing systems for which they are well suited, many systems have features (particularly nonlinearities) which make analytical treatment infeasible. Systems dynamics modeling is often an excellent tool in these situations. While some systems approach a static equilibrium over time, others never settle down to a constant state. The classic Lotka-Volterra model of the relationship between predator and prey [Lotka, 1925; Volterra, 1926], for example, exhibits periodic or even chaotic long-term behavior – never settling down to a constant level for either predator or prey species. Other systems may have equilibria which would be stable if they were ever reached, but conditions may change too quickly for the system to ever reach them [Epstein and Axtell, 1996]. In each of these cases, a systems dynamics approach allows a researcher to understand the behavior of the system in ways that would be impossible with a static approach.

Adaptive Agent Modeling as a Form of Systems Dynamics Modeling

While Ruth and Hannon's high level taxonomy of models as static, comparative static, and dynamic is undoubtedly useful, it is also very broad. Each of these three categories contains many clearly distinguishable species of model, in the

same way that the categories animal, vegetable and mineral are useful, but far from definitive. Generally speaking, however, systems dynamics models have been used so much more broadly than the other types of dynamic model that they are often equated with this whole category of models and referred to simply as “dynamic models”.

While the systems of differential equations used by systems dynamics models are ideally suited to describing the way that many systems evolve in time, they are not appropriate to all situations. In systems which include strategic actors (i.e. people) the future of the system often depends less on its current state or past trajectory than on inferences about the behavior of others and the anticipated future of the system. Game theory (and evolutionary game theory) provides a mathematically rigorous way of exploring such systems. However, as with the formal analysis of systems of differential equations, many non-trivial systems in game theory prove impossible to analyze in any meaningful way.

Much as systems dynamics modeling provides a less formal but more flexible way of handling complex systems of differential equations, adaptive agent modeling provides a less formal way of dealing with the issues of imperfect information, bounded rationality, and strategic inference which would be formally modeled using evolutionary game theory.

The parallelism implied above is, however, not exact. Systems dynamics software provides a user friendly means of constructing differential equations and a numerical engine for integrating them. Such software is a means of constructing and exploring systems of differential equations. Adaptive agent modeling, however, does

not bear the same close relationship to evolutionary game theory. Though the agent approach makes it relatively easy to handle evolutionary game theoretic problems that would be extremely awkward within the systems dynamics paradigm, adaptive agent modeling is a broad way of thinking about modeling.

The more precise parallel is between adaptive agent modeling and the approach to modeling embodied by systems dynamics modeling – which is often referred to as “systems thinking” [Sterman, 2000]. These both represent general approaches to decomposing a complicated system into meaningful parts which can be recombined in way that contributes to understanding of the system.

Both approaches have spawned a host of software environments which facilitate the development of models. Leading systems dynamics packages include Stella, Madonna, and Vensim. Commonly used adaptive agent packages include Swarm, Repast, and Ascape. Within each class, these packages look reasonably similar, whereas between classes they look quite different. Systems dynamics packages generally build their interfaces on the visual language of general systems theory [von Bertalanffy, 1968], and build their analytical tools around the integration of differential equations. Adaptive agent packages, in contrast, generally facilitate the use of object oriented programming techniques and provide tools for managing the activation, interaction, and behavior of agents.

While the functions provided by these software tools are very different, the underlying goals of both approaches to modeling are essentially the same – to track the behavior of a system through time, and in so doing to develop an understanding of

which parts of the system, and which relationships among these parts, are most important to this behavior.

A major way (perhaps the major way) in which adaptive agent thinking differs from systems thinking is that it takes the physical parts of the system (the agents) as its basic units of analysis. This is a contrast to systems thinking which takes the stocks and flows of aggregate quantities as its basic units of analysis. This means that adaptive agent models involve a collection of similar but in some way heterogeneous parts. While these parts may or may not have identical internal structures, they always have heterogeneous internal states. The strength of the adaptive agent modeling paradigm is, fundamentally, its ability to retain the heterogeneity of system parts while developing a rigorously defined numerical model.

Uses of Adaptive Agent Modeling

Axtell [2000] identifies three distinct types of situations where adaptive agent modeling is of use. First, there are cases where equations describing the system of interest can be written down and solved either analytically or numerically. While more traditional simulation techniques (systems dynamics, etc.) are capable of dealing with such systems, adaptive agent modeling provides a novel way of approaching these problems which may be clearer and more flexible in some cases. Second, there are cases where the equations describing the system can be written down, but can not be solved either analytically or through numerical integration. In these cases, the agent approach can make unique contributions to understanding the problem. Third, there are cases where writing down equations is simply not useful –

where the analysis of these equations would not give us the insights that we seek even if we could do it.

In models of the first class, the agent approach is not strictly necessary, but it is often helpful. The agent model can be used to verify the results of a model which has been solved analytically or numerically. It can also be used to present the result of a more complex mathematical model in a way that is more accessible to a lay audience. Because an agent model can often be specified with simpler equations than an equivalent analytical model, and because the output of an agent model generally lends itself to presentation in a graphical form, the agent approach can be a useful complement to more rigorous mathematical models for the purposes of demonstrating results and building confidence in an analytically or numerically tractable model.

Axtell provides a careful taxonomy of models of the second sort, which can be described mathematically, but where these descriptions are difficult or impossible to characterize completely using either analytical or numerical methods. These include models with badly behaved equilibria, particularly models where the features of interest are not equilibrium states, but rather the fluctuations that the system goes through on its path toward equilibrium. Systems of this sort are often impossible to handle analytically. Systems dynamics simulations are often of great use for systems that can be written down clearly but which resist analysis, however, in cases where heterogeneity of agents, spatial location of agents, or complex internal state of agents contributes significantly to the dynamics of the system, the structure of the system will lend itself poorly to the types of numerical integration on which systems dynamics modeling packages rely.

Axtell defines a third category of systems for which writing down and solving equations is not a productive activity. Because he is writing for a highly technical modeling audience, he defines this category quite narrowly: these systems are ones where writing and solving equations is not productive even in theory. In thinking about policy, however, it is useful to relax this definition a bit to include systems for which writing down and solving equations would be so complex, and the insights gained so hard to fathom, that such approaches are of no practical use. Many systems which rely on agent heterogeneity for critical parts of their dynamics fall into this category. This is particularly true when this heterogeneity is spatial in nature as when an agent's rationality is bounded by the information that it can gather using vision with limited range [Dibble, 2001].

Adaptive Agent Modeling in a Policy Context

The adaptive agent approach to modeling has its roots in Schelling's neighborhood segregation model [Schelling, 1969]. This model had reasonably direct implications for housing policy – providing novel insight into the dynamics of segregation and informing the debate about the kinds of policies which might alleviate it. Many later applications of the method, however, have been geared more toward establishing basic principles in social science, rather than the direct guidance of policy. Influential models of this sort include Epstein and Axtell's [1996] "Sugarscape", Robert Axelrod's work with the iterated prisoner's dilemma [Axelrod, 1984]. Work of this sort has made significant contributions to social science, but has generally yielded results which are too conceptual and qualitative to have definite implications for policy.

Because of the foundation laid by these basic investigations, however, the field of adaptive agent modeling seems poised to emerge as a tool for public policy analysis. This dissertation presents three cases where adaptive agent modeling stands to contribute significantly to the world's understanding of contemporary policy issues. Each of these cases illustrates one of Axtell's categories for the use of agent models: one that clarifies mechanisms and presents results where full analysis is possible, one which uses agents to conduct a numerical simulation in a case where the system can be stated but is both analytically intractable and ill suited to numerical simulation using more traditional techniques, and one which produces insights into a system which is not well suited to traditional mathematical analysis.

Chapter II presents an analysis of international trade, using an agent model to explore the impact of relaxing the assumption of decreasing returns to scale on the "infant industries" argument in development theory. It provides support for the notion that the presence of increasing returns to scale in the early stages of industrial development justifies certain types of protectionism in some cases.

Chapter III produces insight into distribution of city sizes within countries using a model which is simple, but intractable. It generates insight into the most commonly observed distribution of city sizes as well as various departures from it by using a simple adaptive agent model which relies on bounded rationality and lagged adjustment for its dynamics. This model contributes to the understanding of a longstanding puzzle in economic geography and provides policy relevant suggestions for the management of third world megacities and for easing economic transition in Russia.

Chapter IV of this dissertation presents an analysis of data from the Guatemalan civil war which indicates that civil violence is an example of the kind of complex dynamic system for which agent based modeling is uniquely suited among quantitative methods. It compares this data with results from an agent model, providing some insight into the nature of conflict and important directions for further research in this area.

These chapters serve to demonstrate the potential for the policy relevant application adaptive agent modeling by showing how the method can contribute conceptual clarity, produce novel results, and allow for rigorous, quantitative work to be done in areas which have often been thought to be too messy for quantitative approaches.

Chapter 2: The Importance of Assumptions: Adaptive Agent Modeling as a Tool for Trade and Development Theory

Of the many beautiful results which have emerged from economic theory over its long history, few are as elegant or have been as influential as Ricardo's principle of comparative advantage in international trade. This principle is often taken to prove that all nations, regardless of their level of development or productivity, can only benefit from increased international trade. Indeed, this argument is so counterintuitive on its face, but so convincing on further thought that it has come to dominate the thinking of those concerned with international trade, often leading them to overlook the assumptions on which the argument rests.

Every model rests on a set of assumptions. When modeling is conducted in the service of policy analysis, it is particularly important that these assumptions be made plain and that the result be recognized as the result of those assumptions. One critical assumption on which the comparative advantage argument depends is that there are constant or decreasing returns to scale in all industries. The relaxation of this assumption complicates analysis somewhat, leading to multiple equilibria and destroying the market's ability to deliver a unique outcome which can be considered to be "optimal" in some objective sense.

While an adaptive agent model is not strictly needed to explore the implications of relaxing this assumption, the adaptive agent approach can be used to build confidence in the insights generated through analysis and to communicate them to policymakers with limited background in economics. In this chapter, I will review

two models which seek to realign the generalizations from trade theory with their underlying assumptions. I will then proceed to demonstrate how an adaptive agent model can be used to illustrate these points in a way that clearly shows how the results follow from the assumptions about the behavior of the people and nations involved.

Gomory and Baumol's Model of International Trade

In their book, *Global Trade and Conflicting National Interests*, Ralph E. Gomory and William J. Baumol persuasively show that relaxing the assumption of decreasing returns to scale for national industries dramatically changes Ricardo's policy conclusions based on comparative advantage. With the introduction of startup costs and increasing returns, the situation goes from one of always coincident national interests in favor of openness, to a more nuanced picture where interests sometimes coincide and sometimes conflict.

A Place for Policy

A major result of their analysis is to move international trade theory out of the realm of pure efficiency analysis, making way for discussions of equity and the application of policy. In their analysis, it becomes clear that the market can not be expected to deliver a single, "optimal" pattern of production which allows each country to make the most of what God has given it. Rather, the market can produce myriad stable patterns of production. Some of these patterns are more efficient, some less, some distribute income relatively evenly among nations, some distribute income very unevenly. Gomory and Baumol argue convincingly that which one of these

equilibria the market produces depends, to a great degree, on history and therefore on temporary policy measures such as the protection of infant industries.

Under the traditional assumption of decreasing returns, the market can be expected to produce a unique allocation of production and income based on each country's natural endowments, which are given. This equilibrium is independent of history in that over the long run, the system can be expected to allocate production in the same way regardless of the order in which nations develop. Barring market failures, this also results in global production at the maximum scale which demand and technology allow at any given time.

If we relax the assumption of decreasing returns and allow some industries to display increasing returns over at least part of their range of production scale, natural endowments come to matter much less and have little to do with the distribution of productive capacity. Those who are first to enter an industry face falling costs as they increase production, making entry difficult even when the entrants have a lower wage bill. This means that it is often the first county – not always the best suited one – which ends up producing a given product.

In Ricardo's day, the assumption of decreasing returns was a reasonable one. Agriculture made up the largest share of even the most highly developed nation's utility. In many agricultural sectors decreasing returns still dominate: the best land is used first with production increases requiring the use of increasingly marginal lands and more intensive (and expensive) management techniques. Before the industrial revolution, this principle held even in manufactured goods: a hat maker could make

only so many hats in a day, and there quickly came a point where supervising more apprentices became uneconomical.

During the industrial era, however, agriculture and hand crafts became relatively minor economic sectors while large scale manufacturing and high-skill services became the driving force behind the rapid growth of economic activity. These sectors, however, display a different type of productivity curve. While the first tomato may be the cheapest to grow, the first automobile is far from the least expensive to manufacture. In many modern industries, economical production requires huge scale, and that huge scale requires tremendous investment, a high level of skill, and the reputation required to bring the resulting products to market. Gomory and Baumol refer to industries characterized by high startup costs due to significant economies of scale (like automobile manufacture), as “retainable” industries, because once a nation has developed such an industry and realized the resulting cost reductions it becomes very difficult for another nation – even one with lower labor costs and more plentiful raw materials – to take that industry away through competition.

For the sake of simplicity in the models that follow, we will use production functions which exhibit increasing returns throughout their range of production. This is, however, not essential to the argument. An industry is retainable so long as enough of the early part of its production cost curve is characterized by increasing returns that an entrant would be unable to coordinate sufficient capital to reach the later phases of constant or decreasing returns.

Multiple Equilibria

A world with retainable industries has the potential for a great many equilibria (in the two country case, there can be 2^n stable equilibria; where n is the number of industries). Gomory and Baumol observe that these equilibria are not arranged at random, but fall into definite patterns. In the extreme case, one nation may have all of the retainable industries and a high standard of living, while the other nation subsists in poverty. The poor nation is unable to purchase many of the goods produced in the rich country, and it is also unable to develop its own industries because its costs of production are still higher than those in the rich country – so the products of its infant industries would not be competitive, even if they were produced. Because manufactures are less expensive to import than they are to make, the best that the poor nation can do (in the short run) is to produce its low-margin agricultural goods and trade them for small quantities of high value added manufactures from abroad.

Because one country with a high standard of living is making all of the industrial products in this scenario, its labor costs are high and its workforce is fragmented between many industries. Meanwhile, the labor force of the poor country sits in idle poverty, producing next to nothing. In this situation, world output is lower than it would be if the retainable industries were divided between the two countries, employing their combined labor force to produce tradable goods. On a graph with income share on the x axis and world output on the y axis, the various mixes of production form an inverted “U”, with low output associated with a high

concentration in either country and higher output associated with a more balanced division of industries.

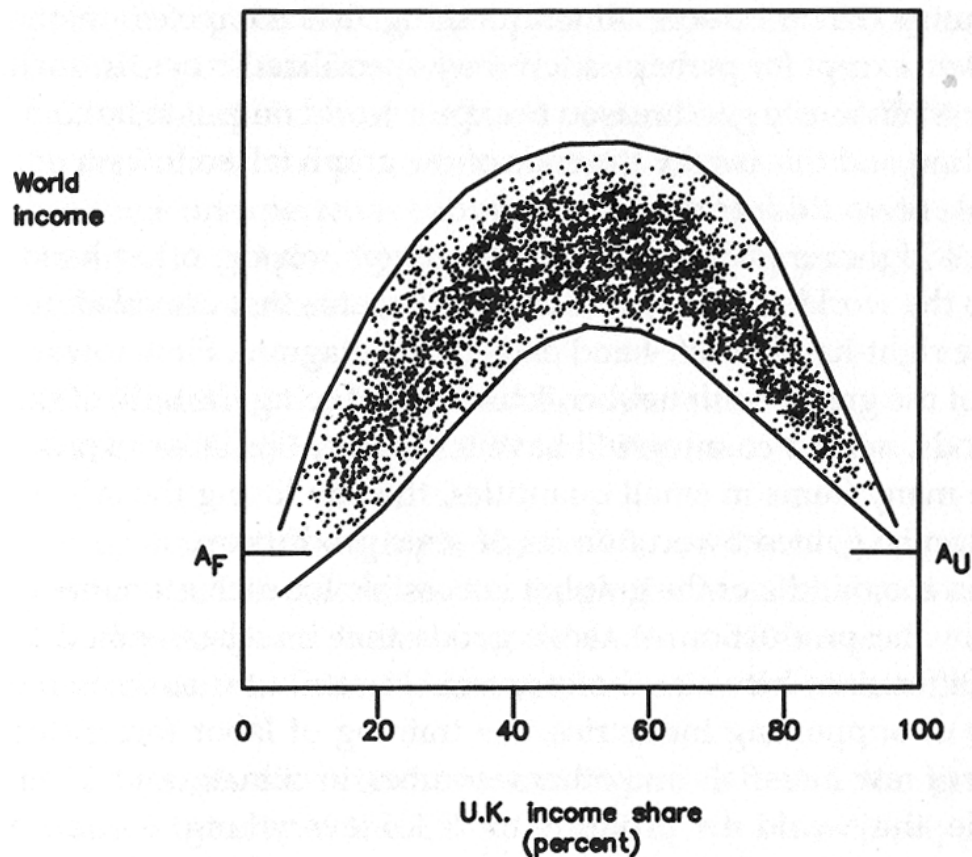


Figure 2.1: Multiple equilibria in a world with increasing returns to scale. Reproduced from Gomory and Baumol [2000].

Gomory and Baumol further point out that this possibility space is actually slightly more complex than a simple inverted “U” because of both natural advantage and synergies between industries. While natural advantage does not play the large role that it did in Ricardo’s theory, there is still a place for it in the world of retainable industries. Some countries are simply better suited to produce some things. If, by accident of history, industries develop in countries where they are not particularly well suited, it is possible to produce an even division of industries between countries

which produces less than the maximum possible because the industries are located in the “wrong” countries.

Synergy between industries (or the lack thereof) can also lead to different levels of output given the same percentage division of industries between nations. Some industries work well together (e.g. steel making and automobile manufacture) while others do not (e.g. paper making and destination tourism). A division that keeps synergistic industries together while separating those that clash will be more productive than one that does the reverse.

Natural advantage and industrial synergy both lead to a range of possible outcomes for each division of industries between countries. The curve of possibilities, therefore spreads from an inverted “U” to an inverted boomerang which is thin at its tips (because there is only one way for the industries to be packed into a single country) and thicker in the middle, where the industries can be divided in many ways, some more efficient than others (figure 2.1).

Cooperation and Conflict

Gomory and Baumol proceed to unpack this distribution, analyzing the implications of this way of looking at things for the output of each country individually. Using essentially the same logic with which they produced the inverted boomerang for world output, but changing the y axis to reflect national output, they now produce a crossing pair of skewed boomerangs, one for each country. These shapes resemble the shape for world output, but are asymmetrical, with a higher peak on the side of the graph which reflects the larger share of industries for the nation in question.

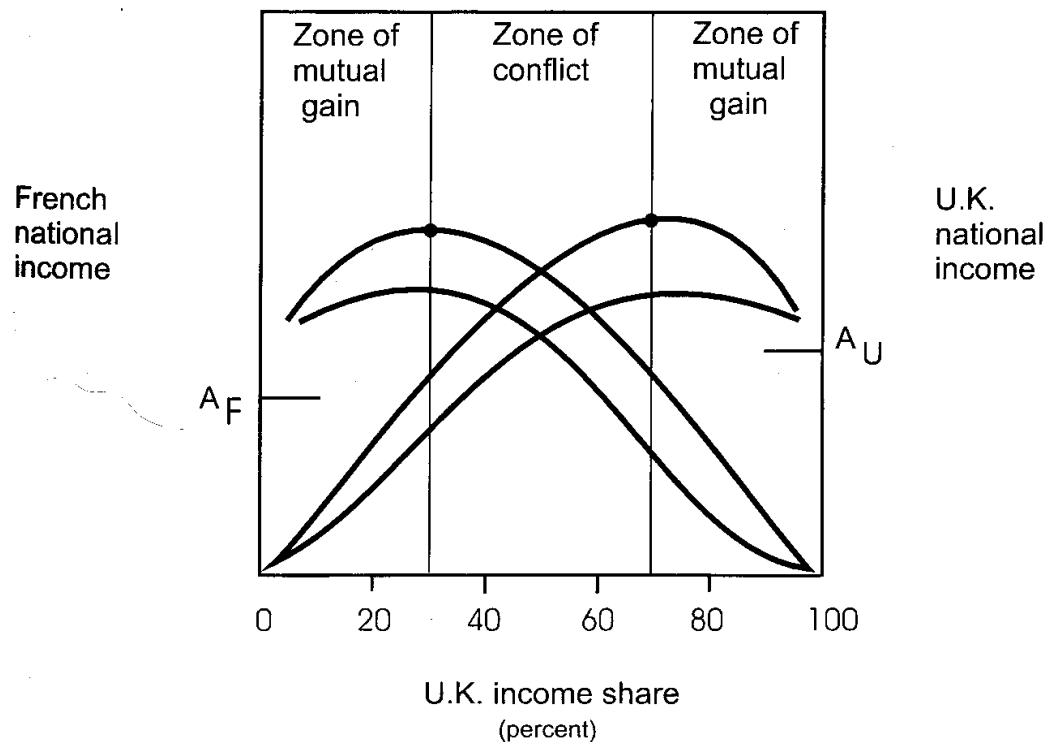


Figure 2.2: Zones of mutual gain and zone of conflict in bilateral trade. Reproduced from Gomory and Baumol [2000].

From this graph (figure 2.2) one can see that there are zones where the interests of the countries either coincide or conflict. In the zones of mutual gain, the curves of both countries slope in the same direction. This indicates coinciding interests. If one partner has a great many industries while the other has very few, both can improve their position by transferring some industries from the richer to the poorer country. This benefits the poorer country by allowing it to produce goods for export and to enjoy the resulting increase in income. It also benefits the richer country by creating a market for its exports and allowing it to purchase low priced goods from its trading partner. In these zones of mutual gain, both partners benefit from increased trade. There is also, however, a zone of conflict where the curves

slope in opposite directions. This indicates that one partner benefits from increased trade at the expense of the other. In this central region of the graph, any movement toward more balanced development leads to greater income for the poorer partner, but less income for the richer one.

It is important to remember here that all of the points within the curves are stable equilibria. If the system finds itself outside of these curved areas, it can be expected to work its way back into them. However, once the market is within these areas, it can not be expected to move the balance in any particular direction, or even to find the maximum output position for a given balance. Instead of market forces, movement within these areas is due to policy decisions: trade policy, development policy, industrial policy, etc.

Revisiting the Infant Industries Argument

Having developed this model of trade in a world with retainable industries which exhibit increasing returns to scale over at least the early part of their development cycle, Gomory and Baumol go on to develop a similar model for industries with linear returns to scale, but where productivity improves with experience. Though some of the details of the analysis differ, the upshot is the same: first movers have a substantial advantage and the market can produce myriad stable outcomes that differ greatly in their equity and efficiency. This conception indicates that the often maligned “infant industries” argument for protectionism in underdeveloped nations has a good deal of merit. Once a country with low wages attains a competitive position in such a skill based industry, its low wage bill will

keep it competitive. However, such entry is only possible once the industry has become efficient enough (through experience) to compete.

This way of looking at development and trade puts the plight of underdeveloped nations in new perspective. Under the traditional assumption of decreasing returns, capital would be expected to flow from wealthy nations to poor ones, eventually equalizing incomes all around and producing high level of world output. To the extent that differences in income remain, in the traditional view, these should be due to differences in the natural endowments of the nations. This world view absolves market participants from any concerns about equity in trade or development because the market is basically egalitarian. Though the developed world may have gained its wealth by having the good fortune to develop first and by exploiting other areas during the colonial era, the market is always working to erase these historical flukes and iniquities. If the market is only allowed to function without impediments, it will eventually allow every nation to produce at the highest level at which its land and people are capable.

Gomory and Baumol make it clear that over a broad range of industries – particularly those which drive the modern economy – this picture is extremely misleading. Underdeveloped countries are not underdeveloped because they are somehow inferior in terms of either land or people. Rather, the operations of the modern international economy work to lock them into their historical patterns of poverty.

Policy Implications

While this finding would seem to be bad news for the developing world, the analysis also offers hope for the most underdeveloped places. While the analysis makes it clear that the market will not automatically improve the lot of Sub-Saharan Africa (for example), it also makes it clear that it is in the interest of wealthy nations to assist the poorest nations to gain a foothold in industries where they have the potential to succeed. Any job transferred from the US to Liberia can be expected not only to make Liberia better off, but to generate more than one job in the US because the reduction in aggregate demand in the US (from the lost job) will be more than offset by an increase in aggregate demand for US imports in Liberia, as well as a reduction in price in the good that is now manufactured abroad. This should result in a more jobs and more consumption in both countries.

They estimate that the ideal trading partner for a wealthy nation is one which has a GDP per capita of about one quarter of its own. This makes Mexico something close to an ideal trading partner for the United States in the sense that the US could not improve its lot by seizing industries from or conceding industries to Mexico. If this analysis is correct and Mexico defines the border between the zone of cooperation and the zone of conflict for the US, then those nations with per capita GDP lower than Mexico (approximately two thirds of the world's nations) fall into the zone of cooperation, where the US could only benefit by helping them.

An Adaptive Agent Model of International Trade

In an effort to gain insight into the mechanisms involved with international trade and development, we can construct a simple adaptive agent model of production

and trade. This model will follow the basic outline of the classic Heckscher-Ohlin trade model, but will further disaggregate the model, resting it on the behavior of individuals and firms. The model is capable of reproducing a contemporary analysis of trade from Paul Samuelson as well as verifying the retainability of industries as described by Gomory and Baumol and demonstrating how recognition of this retainability has important implications for the long discredited infant industries argument for protection of developing markets.

Model Specification

We begin by defining the agents. We define two types of agents: citizens and nations. Citizens are each associated with one nation and possess one unit each of labor and capital, which they choose to deploy in one of two national industries depending on which pays the higher wage or higher return to capital (they may choose to work in one industry and invest in the other). They use these wages and returns to demand goods.

Nations possess national industries (we can follow convention by thinking of them as wine and cloth) which produce goods according to Cobb-Douglas production functions using the labor and capital which the citizen agents provide. They calculate wages and returns to capital along with prices for each of the goods produced. When trade is enabled, they also engage in trade, importing more of a good if its price is lower in the other country and paying for these imports by bartering with goods from the industry where their price is lower.

More specifically, the citizen agents have three basic state variables: a job, an investment, and a demand function. In each round, each agent does these things:

- Asks the nation for the current price of both wine and cloth.
- Asks the nation for the current wage in the industry where the agent works.
- Asks the nation for the current return on capital in the industry where the agent has invested.
- Calculates its demand for both wine and cloth based on its income (from wages and investments) and the prices of the two goods using the simple hyperbolic demand function $D_w = Y/2P_w$. This amounts to saying that each agent spends half of its income on each good – buying less and more of the good as the price goes up and down.
- With a probability of one percent, the agent reexamines its job and investment choice, changing jobs or shifting its investment to the industry which provides the higher wage or return to capital. The low rate of turnover in employment and investment insures that the model is able to adjust to each change, thus avoiding stampedes from one industry to another which dramatically overshoot the required correction in the employment or investment level.

The nation agent also has several state variables. The structure of the nation's two industries is given by a pair of Cobb-Douglas production functions of the form $Q_w = A * L_w^\alpha * K_w^\beta$, where the quantity of wine produced Q_w is the product of an efficiency A , the amount of labor devoted to wine L_w to some exponent α and the amount of capital K_w devoted to wine to some exponent β . These parameters (A , α , and β) are state variables.

Because the model relies on barter rather than money, the price of one good (wine) is fixed at 1, while the price of the other good (cloth) adjusts to reflect its relative scarcity. The price of cloth is adjusted upward by a small amount when demand for cloth exceeds its supply and down by a similar amount when supply exceeds demand. Because wages and returns on investment are calculated as shares of current production, Walras' law ensures that if the cloth market clears, the wine market will also clear. The price of cloth is a state variable.

Finally, when trade is opened, the nations barter goods. Cloth flows from the country in which its price (relative to wine) is lower to that where its price is higher, with compensation being made in wine according to the current price of cloth. When the international market is out of equilibrium (i.e. when the price of cloth differs between the two countries) the trade price of cloth is taken to be the average price between the two countries. The amount of cloth exported is increased by a small amount when the nation's partner has a higher relative price for cloth and is decreased by a small amount when the partner has a lower relative price for cloth. This level of trade is the nation's final state variable.

In each round, each nation does these things:

- Counts the number of citizens working and investing in each industry.
- Determines the quantity of each good which it will produce using each industry's production function and the current level of employment and investment in each industry.
- Determines the wage for each industry by calculating the marginal product of labor in that industry by subtracting the current level of production from the production that would result from the addition of one additional unit of labor.
- Determines the return to capital for each industry by subtracting the wage bill for that industry from the total output of the industry (at current prices) and dividing by the number of investors in the industry.
- Adjusts the price of cloth as described above.
- Adjusts the level of trade to reflect the new price level in both countries as described above.

These straightforward behavioral rules are adequate to reproduce the primary features of the Heckscher-Ohlin trade model in a dynamic context. This model is implemented in Java using the Ascape (Parker 2000) modeling framework. The agents are represented by Java object classes, while Ascape handles the randomized

agent activation regime (i.e. agents activate in a changing, randomized order) while also facilitating the collection of statistics and the production of graphical output.

Samuelson's Analysis of Outsourcing

Paul Samuelson, who is widely considered to be the Dean of neoclassical trade theory, has recently published a paper [Samuelson, 2004] which takes mainstream trade theorists to task for over generalizing the benefits of free trade by demonstrating that there are situations where the gains from trade for one nation can be undone by technological developments in a second nation. Because Samuelson sets up his simple analytical model in a way that is compatible with our agent analysis, it serves nicely to validate our model. If the model is correctly specified, it should be able to produce results which agree with Samuelson's mathematically rigorous analysis.

Samuelson asks us to consider two countries designed to look something like the US and China. His stylized US has 100 citizens while his stylized China has ten times that population with 1000 citizens. For the sake of symmetry, he further assumes that the US average productivity is ten times as high as Chinese productivity, thus producing equal amounts of total production in the two countries (though Chinese per capita productivity is only $1/10^{\text{th}}$ that of the US). These productivities are asymmetrically distributed between industries, however, with the US having Ricardian productivity parameters of 2 and $1/2$, while China has parameters of $1/20$ and $2/10$.

One problem with models of this sort, which represent the economy in barter terms, is that it has traditionally been difficult to compare outcomes in absolute terms.

Samuelson overcomes this problem by pointing out that there is a definite relationship between demand and utility functions. He assumes a J. S. Mill style pair of hyperbolic demand functions: $D_c = Y/2P_c$ and $D_w = Y/2P_w$. These demand functions imply that consumers spend half of their income on each good. He then shows that these are the logical outgrowth of a utility function $U = (C*W)^{0.5}$ which takes the geometric mean of the consumption of the two goods as a measure of welfare. This relationship allows us to measure the total utility of each nation. In the absence of money, this utility measure allows us to assess the value of the nation's consumption. It can thus be used as a fair measure of the nation's utility.

Samuelson refers to this measure as a proxy for GDP, but this is not necessary or entirely correct. Generally, GDP is taken as a proxy for total utility, which is difficult to measure. GDP is, however, a poor proxy for a variety of reasons [Daly, 2003]. Because we are working with a theoretical system, it is possible for us to work directly with utility rather than resorting to the poor proxy of GDP. In the current specification of the model, we would assume that GDP and utility would be highly correlated but other interesting formulations would weaken this link. To avoid confusing the end (utility) with its means (GDP), we will break from Samuelson's usage and refer to the geometric mean of consumption as utility rather than GDP.

Using these production and demand functions, Samuelson demonstrates that there are substantial gains to be had when the countries specialize and trade the product in which they are relatively strong for that in which they are relatively weak. In autarky, the US can produce 100 units of cloth and 25 units of wine. This gives a utility of $(100*25)^{0.5}$ or 50. China, similarly, can produce 25 units of cloth and 100

units of wine to achieve the same utility level of 50. US utility per capita is therefore $50/100$ or 0.5, while China's is $50/1000$ or 0.05.

Samuelson then demonstrates that, under free trade, the US is able to specialize in cloth, producing 200 units of cloth, whereas China is able to specialize in wine, also producing 200 units. Because of the symmetry of the example, each country is able to trade and consume 100 units of each good, thus raising total utility in each country to $(100*100)^{0.5}$ or 100 units. Both countries have thus doubled their real utility by specializing and trading.

Finally, Samuelson demonstrates that not all technological changes need be beneficial for both nations. For the sake of this example, he posits a tremendous technological improvement in China's cloth sector (where the US had previously been stronger) from 0.05 to 0.8. This leaves cloth productivity substantially below the US level of 2, but much higher than it had been. This change serves to equalize the factor prices in both countries (the ratio of the efficiencies in both nations is now 4). This equalization removes all incentive to trade, reducing the problem to calculating the output of each country in autarky.

The result is a boon for China and a plague for the US. China is now capable of producing 400 units of cloth and 100 units of wine for a total utility of $(400*100)^{0.5}$ or 200 (0.2 per capita), while US once again can produce $(100*25)^{0.5}$ or 50 (0.5 per capita). Chinese consumption thus expands by a factor of four while US consumption is halved.

Samuelson uses this model to argue that outsourcing of high technology jobs from the US to India and China is not automatically good for both nations. Indeed

the transfer of jobs in a sector where the US was once a leader to countries which did not previously participate heavily in such industries has the potential to make the economies of various nations more alike in their productivity, thus eroding gains from trade to which the US has become accustomed.

Verifying the Agent Model

We can gain some confidence in both the agent model and in the soundness of Samuelson's analysis by verifying that they both produce the same result. Because our modeling approach is compatible with Samuelson's analysis, it is easy to translate his numbers into parameters which can be plugged into the agent model.

The "US" nation agent begins with 100 citizens. It has two industries specified by these production functions which (following Samuelson) exhibit constant returns to scale:

- $Q_c = 2 * L_c^{0.5} * K_c^{0.5}$
- $Q_w = 0.5 * L_w^{0.5} * K_w^{0.5}$

The "China" nation agent begins with 1000 citizens. Its industries are similarly specified with these production functions:

- $Q_c = 0.05 * L_c^{0.5} * K_c^{0.5}$
- $Q_w = 0.2 * L_w^{0.5} * K_w^{0.5}$

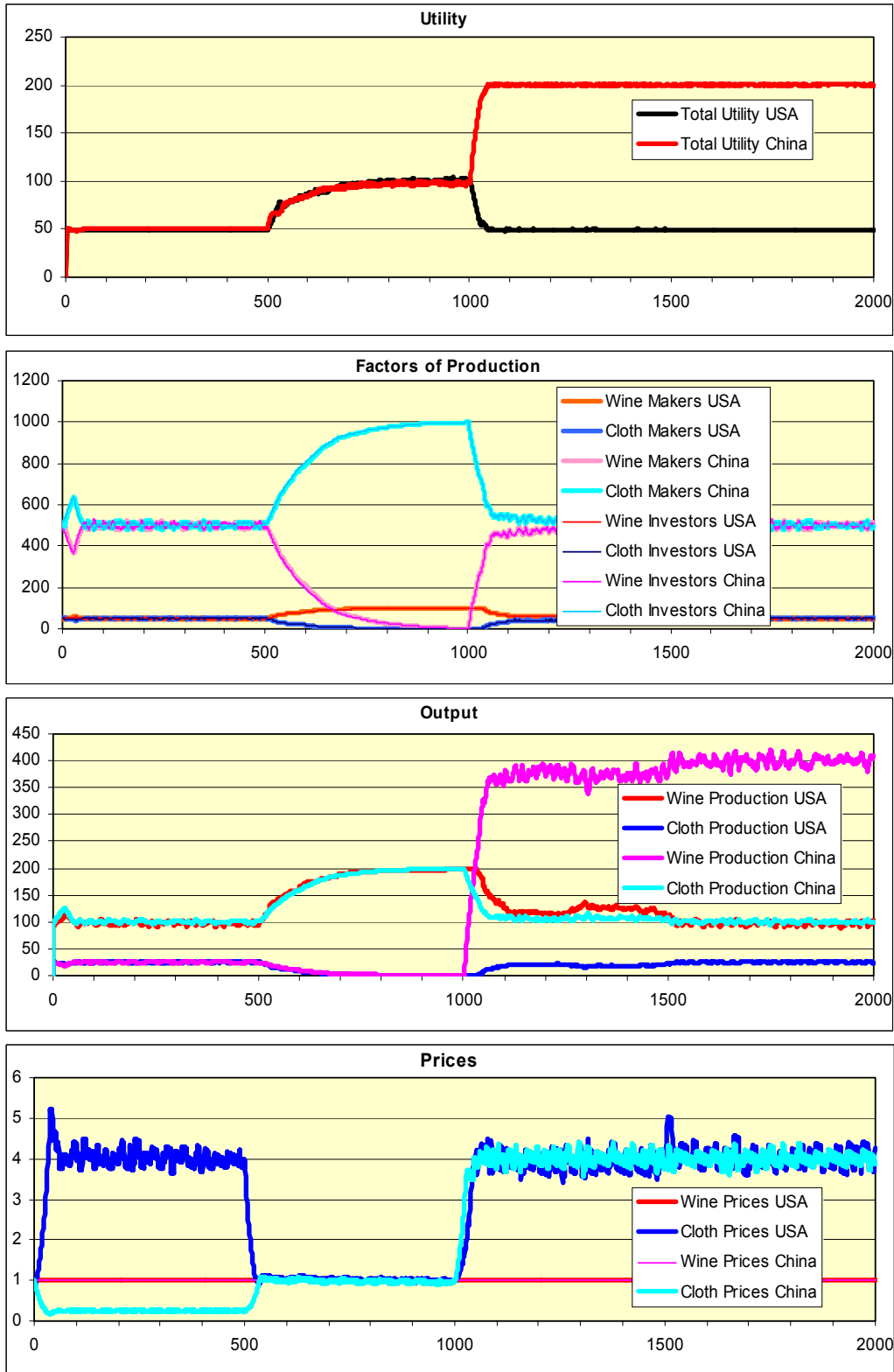
The citizen agents of each country are initially randomly assigned a job, an investment and a demand function as described above. This demand function is identical for each agent.

We begin the model run in autarky. After 500 rounds, both nations have established equilibrium production at 50 units of utility. At 500 rounds, we open

trading which allows the nations to import a good if its relative price is lower in the other country. This results in a major restructuring of each economy.

After another 500 rounds, at round 1000, China undergoes its remarkable invention in the cloth industry, raising its productivity there from $1/20$ to $8/10$. As Samuelson's analysis indicates, Chinese utility jumps to 200, while US utility falls back to its previous autarkic level of 50. After yet another 500 rounds, trade is stopped and the model shows no major difference, thus demonstrating that these productivity levels produce trade terms which are functionally equivalent to autarky.

Figure 2.3: Adaptive Agent Realization of Samuelson Trade Model



Verifying Gomory and Baumol's Retainable Industries

Now that we have established the basic functioning of the model, we can use it to look at what happens when we explore the more interesting case where we relax the assumption of constant returns to scale, shifting instead to the combination of increasing and decreasing returns examined by Gomory and Baumol.

As discussed above, one of the foundations of Gomory and Baumol's argument is that relaxing the standard assumption of constant or decreasing returns to scale to allow for increasing returns to scale in some industries changes the complexion of trade theory dramatically. With constant or decreasing returns, the Heckscher-Ohlin (along with its various Ricardian cousins) indicates that the market will always deliver a better result for each country with trade than it will without. Though the standard model is not dynamic, it also implies that changes in productive capacity will be reflected in the market – as we saw in Samuelson's stylized treatment of the US and China.

Gomory and Baumol observe, however, that in a world where some industries produce increasing returns to scale, these industries can be “retainable” by a nation which develops them early. Because costs fall as more units are produced, it may be possible for a nation with a less efficient production function to retain an industry over a later entry which would be able to produce the good more cheaply if only it could attain the required scale of production. As we will see shortly, a late developing country may, under some circumstances, be able to do better in the long run by abandoning trade in some industries all together.

The recognition of the importance of increasing returns is not entirely new, having been explored by such authors as Kenneth Arrow [1962] by Paul Krugman [1979, 1983], Brian Arthur [1989], among others. It has, however, failed to make a real dent in the policy discourse concerning trade and development

We can illustrate the existence of retainability by running our adaptive agent trade model with an appropriate set of parameters. In this case, we imagine a large (500 citizen), industrialized nation and a smaller (100 citizen) “third world” nation which develops later. Once again we have two industries, but this time they are industries of a specific character. One is a basic agricultural industry which exhibits low productivity and decreasing returns to scale. The other is a high productivity industry – let’s generically call it manufactures – which exhibits increasing returns to scale. We will assume for the moment that this industry exhibits increasing returns over its whole range of production.

With the exception of levels of productivity, these production functions are identical in both countries:

- $Q_a = A * L_a^{0.4} * K_a^{0.4}$
- $Q_m = B * L_m^{0.7} * K_m^{0.7}$

As in Samuelson’s case, the nations differ only in their production efficiency in each industry. The developed nation is more efficient in both industries, having an efficiency in agriculture of $A=0.5$ and an efficiency in manufactures of $B=1.0$. The developing nation begins with equal efficiency in both industries: $A=0.2$ and $B=0.2$. This gives the developing nation a comparative advantage in agriculture and the industrialized nation a comparative advantage in manufactures.

We run the model forward as we did in the Samuelson case. For the first 500 rounds, both countries produce and consume as best they can in autarky. For the next 500 rounds, the nations trade, both realizing gains because they are able to specialize in the area where they are most efficient.

As in the Samuelson case, at round 1000, we introduce a substantial exogenous change in productivity in one of its industries. In this case, the developing country drastically increases its productivity in manufactures from a paltry 0.2 to an impressive 1.5, jumping from 20% of the developed nation's productivity to 150%. At this point, however, we observe a marked contrast to Samuelson's giant increase in productivity: nothing happens.

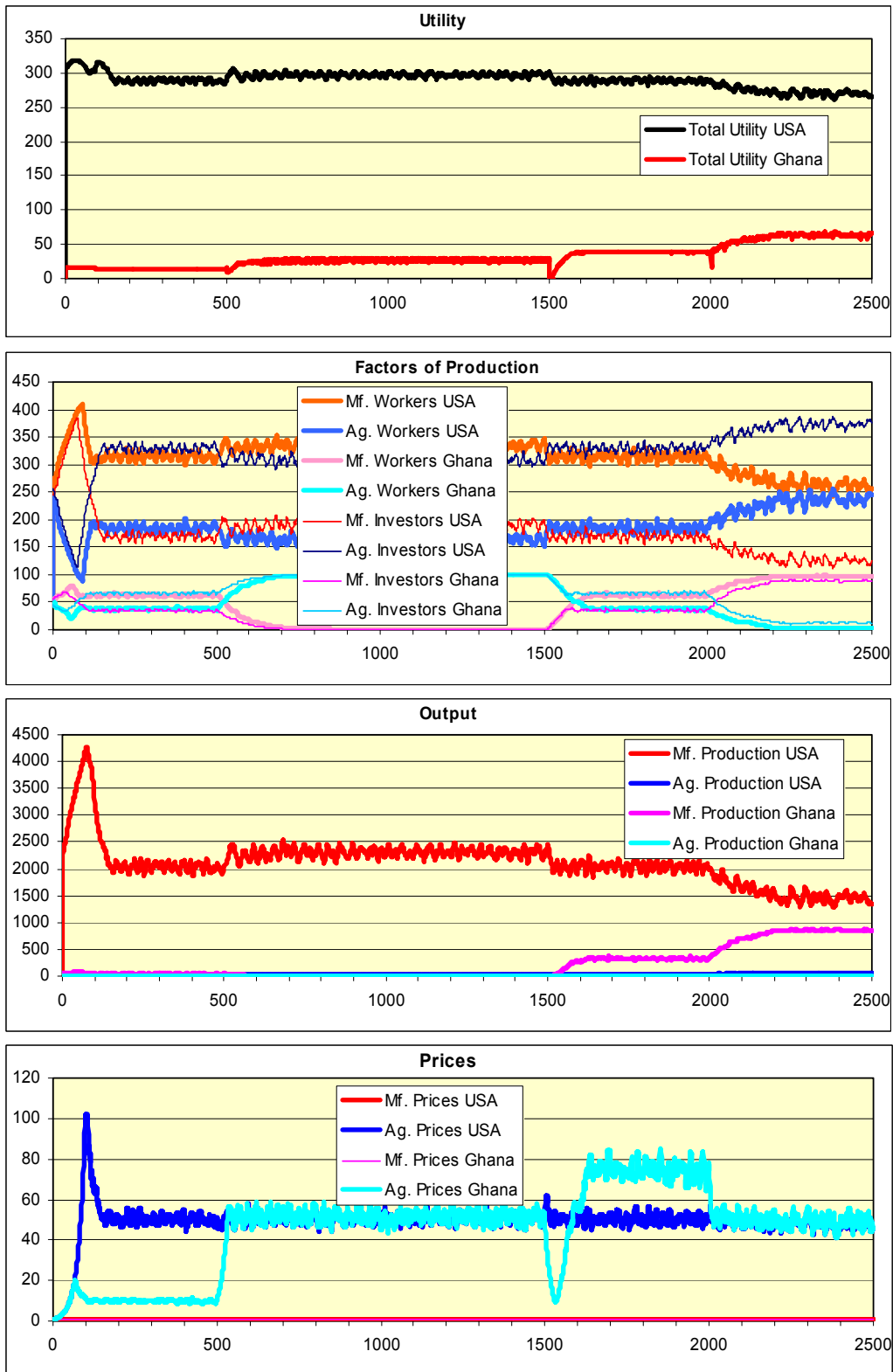
Because the developing nation has specialized in agriculture, it has virtually no industry in manufactures. Any attempt to start such an industry is bound to fail because the industrialized country has attained a scale such that it can produce manufactures more cheaply than the developing nation – even given the developing nation's new, superior productivity at any given point on the production functions. In each round, the citizens and investors of the developing nation examine the feasibility of moving into manufactures, and in each round they find that they can do better by sticking to agriculture. The industrialized nation is thus able to retain the industry despite the fact that, all else being equal, it is no longer the most efficient producer in either absolute or relative terms.

In the Samuelson case, we cut off trade at round 1500 and found that there was no impact on utility in either country because their proportional productivities had become similar. If we cut off trade in this case, something even more surprising

happens. After an initial plunge in utility, the developing country begins to restructure its economy. Where its manufactures had been unable to compete with cheap, mass produced imports in its domestic market, they are now the only game in town. Workers and investors begin to shift away from agriculture and into manufactures. Initially, this sector is not terribly productive, but with experience and scale, it becomes more and more productive. In time, given the parameters we have chosen, the manufacturing sector becomes so productive that the small nation is actually able to do better in autarky than it previously did through trade!

Finally, in round 2000, we reopen trade. The newly industrialized country is now in a much stronger position to compete on the international market and sees a substantial gain. The larger, more established country actually loses more utility as a result of this trade over autarky. It is forced to restructure its economy to produce the lower productivity agricultural good. Because this good has decreasing rather than increasing returns, its productivity erodes as it becomes more specialized, leading to a long term decline in income as compared to autarky.

Figure 2.4: Retainability of Industries with Increasing Returns



Discussion

This case is admittedly highly stylized; however, it makes good sense in terms of development and has important implications for development policy. In the constant or decreasing returns world of neo-classical trade theory, the productivities of nations in different industries determine a unique set of equilibria in trade and utility unless some sort of trade policy intervenes to interfere with trade and lower that utility. A poor country is poor either because it is not very productive, or because it is not making good use of its comparative advantages in productivity through trade.

The policy prescription that comes out of the neo-classical model is simple. Poor nations should try to improve their productivity in areas where they have a chance to compete – keeping wages low and focusing on low skill sectors such as agriculture (the stereotypical example would be bananas). Furthermore, they should seek to increase trade in every situation. The standard set of assumptions about trade indicate that this is the very best they can hope to do. If such a country is unable to compete in any of the more modern industries which are characterized by increasing returns, that is simply because they as a nation are no good at them. Their best strategy for obtaining these high value added goods, in both the short and long terms, is to grow ever more bananas and look for additional markets in which to trade them.

The introduction of increasing returns into this picture changes everything. A poor country no longer faces a simple policy prescription, and the invisible hand can no longer be counted on to deliver the industrial structure which will give the country its highest long-run level of consumption. The multiple equilibrium situation introduced by increasing returns leaves the country with difficult choices. In the short run, protecting a domestic industry will almost certainly hurt them. In the long run,

however, this protection might allow the protected industry to attain sufficient scale that the country would be better off. Even if the long run autarkic equilibrium utility would be lower than the free trade equilibrium, a period of protection and domestic development might allow the protected industry to develop to the point where it could become a competitive producer on the world market, thus allowing the nation to reopen to substantially improved terms of trade and higher consumption. The Asian “tiger” economies come to mind as nations which achieved tremendous development by following this kind of strategy. [UNIDO, 2004]

Next Steps

In this essay, we have used the adaptive agent approach to illustrate a result which can be obtained more simply (but perhaps less convincingly to some) using analysis. This approach, however would lend itself nicely to variations which would be much more difficult to handle analytically.

We have held to the standard economic convention of using consumption as the sole measure of well being. Though this convention is almost universally followed, this probably has more to do with its analytical convenience than it does with any attempt to reflect economic reality. Economic analysis generally assumes a preference curve for goods (as we do in our hyperbolic demand curve and geometric mean welfare function), but assumes that workers are uniformly indifferent about their employment. This adaptive agent modeling framework would make it relatively simple to work with a heterogeneous population of agents who possess different talents for different kinds of work and different preferences for different kinds of work. Not everyone is cut out to be a banana farmer – and not everyone with the

abilities required would want to be one. Such a formulation could reflect not only the efficiencies associated with having a diverse economy which is able to take advantage of people's differing talents, but also reflect the subjective (but very real) welfare gains which would result from people being able to spend their time at jobs which they prefer [Daly, 1996].

Because the current model assumes equal wages and returns within an industry and works to equalize these returns between industries, it can have nothing to say about the impacts of trade on income distribution. While a full scale model capable of reproducing national patterns of income distribution would be more than a minor extension of this model, the ability of the adaptive agent approach to work with heterogeneous agents would make it ideal for this kind of work.

Along these lines, Samuelson [2004] states, "My most important omission, for realism and for policy, is treating all people in each region as different homogeneous Ricardian laborers. That inhibits our grappling with the realistic cases where some Americans (capitalists and skilled computer experts) may be being helped by what is decimating the real free-trade wage rates of the semi-skilled or the blue-collar factory workers." He goes on to discuss ways in which factor price equalization models might predict declining median income even in the face of increasing average income due to increasing inequality. In so doing, he points out that, in a factor price equalization model such as this one, the US unskilled wage would be expected to drop in the face of low wage foreign competition. While it might be possible for the winners in such situation to compensate the losers, he observes that there is no evidence that this has happened or will happen. If citizens were fully aware that this

could happen, a democratic society might well choose to increase median income at the expense of the average (or total) income.

The adaptive agent approach used here would be ideally suited to relaxing the assumption of homogeneous Ricardian laborers. Workers could be endowed with differing abilities in different industries and different levels of effort or energy. Different industries could have various requirements for more and less skilled laborers, with wages reflecting the market for such work. This approach would allow for the rigorous treatment of such issues as offshoring and outsourcing without adding major complexity.

Another way that the adaptive agent approach could contribute to trade modeling would be by providing a natural modeling framework for capturing industrial synergies. A significant part of Gomory and Baumol's analysis rests on the idea that many industries can not operate in isolation, but are dependent on other industries for efficient production. We could further illustrate this point by elaborating production functions to make the output of some industries dependent on the supply of goods produced by others. In the presence of transport costs (which could easily be introduced), this would make some combinations of industries more efficient than others.

It would also be straightforward to generalize this model to include many industries and many nations. This would be useful in evaluating policy issues such as the validity of Gomory and Baumol's claim that it could be in the interest of a wealthy nation to transfer an industry to a poor nation. While their analysis demonstrates that such a transfer would increase global utility, it is not entirely clear,

in a many nation situation, under what circumstances the benefits to the wealthy nation would actually outweigh the costs it incurs. In the two nation case, the wealthy nation sacrifices an industry but is able to reap all of the benefits of lower prices from the lost industry. In the many nation case, the wealthy nation would still incur all of the costs of sacrificing an industry, but the benefits would be distributed among many nations.

This would seem to complicate the self-interest based argument for helping poor nations to take over some of the industries which are currently retained by wealthy nations. While such a move would increase global utility to the point where the winners could, in principle, compensate the losers, this would almost certainly never happen. A multi-nation adaptive agent treatment of this problem could be a useful tool in differentiating the kinds of situations where a pure self interest argument would apply from those which would rely on appeals to the common good (where global welfare would be increased at the expense of national welfare) or to economic justice (where the poor would benefit at the expense of the aggregate).

Finally, the agent framework presented here would be well suited to exploring Daly's [1996] observation (also mentioned by Samuelson [2004]) that the mechanism of the comparative advantage argument depends on internationally immobile capital.

This assumption is explicitly stated by Ricardo [1817], but is generally omitted from modern discussions. Given the realities of early 19th century international travel and communication, Ricardo found this assumption reasonable:

Experience, however, shews, that the fancied or real insecurity of capital, when not under the immediate control of its owner, together with the natural disinclination which every man has to quit the country of his birth and connexions, and intrust himself with all his habits

fixed, to a strange government and new laws, checks the emigration of capital. These feelings, which I should be sorry to see weakened, induce most men of property to be satisfied with a low rate of profits in their own country, rather than seek a more advantageous employment for their wealth in foreign nations.

In the early 21st century, international investment is a much simpler matter and the increasing trend toward globalization continues to make national borders less relevant to investment decisions. Daly points out (following Ricardo closely) that mobile capital shifts the situation from one of comparative advantage – where all nations benefit – to one of absolute advantage. Under absolute advantage total global output can be expected to increase (as capital moves to find its maximum return), but more efficient nations benefit while less efficient nations suffer. In a decreasing returns world, this would lead to equalization of incomes among nations, as capital moved to the places where it was in shortest supply (and thus produced the highest marginal return). In the more complex world that we inhabit, with increasing returns, industrial synergies, critical infrastructure, etc., the effects of relaxing the assumption of international capital immobility are harder to identify with certainty.

An initial exploration of this principle could be conducted by allowing the agents of our model a broader choice of investments. Currently, agents examine the marginal return to capital in the two domestic industries – moving their investments to maximize this return. By allowing the agents to invest in any of the four industries, we should be able to reproduce the basic difference between comparative and absolute advantage.

In its simplest form, the model would pay the return to capital directly to the investor. This would be equivalent to allowing the complete repatriation of revenues (not just profits). Thus, investment abroad would generate considerable demand at

home. The actual fate of revenues from foreign investment is considerably more complex than this [Gomory and Baumol, 2000] and modeling it well enough to make specific policy recommendations would be a non-trivial task. Even a simple model along these lines would, however, make the point that the rosy picture painted by the comparative advantage argument no longer applies. It would make it clear that unless winning nations are prepared to compensate losing nations (which is unlikely), nations would do well to proceed with caution with regard to capital mobility because there is no assurance that each will benefit.

Chapter 3: Beyond Zipf: An Agent Based Understanding of City Size Distributions

George Kinsley Zipf observed in 1949 that the size distribution of cities within nations tends to follow a particular kind of power-law [Zipf, 1949]. This distribution is often described as the “rank size rule” or simply as the Zipf distribution. While Zipf convincingly documented this rule in cities and many other systems (including the frequency of word usage in most languages), he was less successful in explaining its emergence. During the ensuing half century, various theories of city formation and development have emerged, and contributed real insights into the geography and economics of cities. They have, for the most part, however, failed to predict the Zipf distribution of sizes. Another class of theories has been put forward to explain the distribution, but these have tended to rest on unrealistic assumptions, to lack explanatory power, or, at best, to lack the ability to explain the deviations from Zipf which can be observed in many nations. In this paper, we offer a simple, though analytically intractable, adaptive agent model of city size evolution. This model offers substantial insight into the distribution of city sizes in various countries while complementing previous work in the economic geography of cities and offering plausible economic interpretations and logic. The model can also account for several important categories of systematic deviation from Zipf that are observed in empirical data, and offers new insights about how such deviations arise.

The Zipf Distribution

The Zipf distribution is neatly summarized by the expression $S_r = S_0 * r^{-1}$ where S_r is the size of city r , r is the rank of the city (i.e. for the tenth largest city, $r=10$) and S_0 is the size of the largest city. This can be restated as the so called “rank size rule” by saying that the second largest city is half the size of the largest city, the third largest 1/3 as large, the fourth 1/4 as large, etc. One property of this distribution is that when it is plotted as an ordered histogram on log-log axes, it results in a straight line with a slope of -1 (which is the exponent of the power-law).

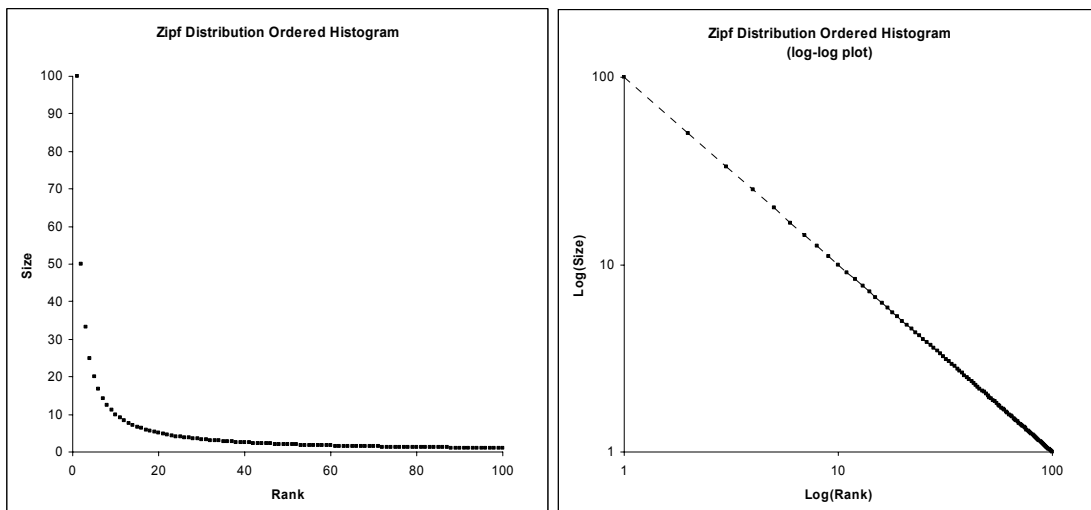


Figure 3.1 Zipf Distribution ordered histogram on normal and log-log axes.

Several Schools of Thought on Why the Regularity Exists

While the Zipf regularity has been observed for some time, it has resisted attempts at theoretical explanation. Fujita, Krugman & Venables (1999) directly address the fit between theory and observation in their chapter entitled “An Empirical Digression: The Sizes of Cities”. They write:

Attempts to match economic theory with data usually face the problem that the theory is excessively neat, that theory gives simple, sharp-edged predictions, whereas the real world throws up complicated and messy outcomes. When it comes to the size distribution of cities,

however, the problem we face is that the data offer a stunningly neat picture, one that is hard to reproduce in any plausible (or even implausible) theoretical model.

The conclusion to this chapter begins by saying, “At this point we have no resolution to the explanation of the striking regularity in city size distributions. We must acknowledge that it poses a real intellectual challenge to our understanding of cities...” Though work in this area has continued in the intervening five years, this remains a valid assessment of the state of the problem.

Attempts to model the dynamics of city size have largely fallen into one of two categories. Models in the first category extend concepts from standard economic theory to apply to city size dynamics. These include externality models, which apply the “Henry George” theorem from urban economics [Marshall 1890, Jacobs 1984, Henderson 1974, Kanemoto et al. 1980], and models which extend Christaller’s [1933] “central place” theory [see Fujita & Mori 1997]. Such models are well integrated with the existing body of economic theory, and are often consistent with other economic evidence about city dynamics. Unfortunately, none of these models convincingly produce the empirical regularity of the Zipf distribution.

Models in the second category apply one or more abstract stochastic processes to represent city size dynamics. Early examples included Simon’s [1957] proportional growth model and Hill’s [1975] application of the Bose-Einstein process. More recently, the most prominent models in this category have focused on descriptions of city growth as a “Gibrat process” [Gibrat, 1931]. Papers applying the Gibrat processes include Gabaix [1999] and Reed [2001]. These processes have all been shown mathematically to successfully generate a stable power-law distribution, and in many cases to closely replicate the Zipf distribution itself. However, such

models have little or no economic content. They demonstrate that the Zipf regularity follows from other statistical regularities, but they do not offer a set of behavioral principles which would produced these regularities. As one recent paper put it: “this collection of models is essentially statistical—they seek to *generate rather than to explain* the regularity” [Overman & Ioannides, 2001]. It is often unclear how the abstract mechanisms represented in many of these models can be useful metaphors for real-world social or economic processes. Indeed, in some cases, closer examination has found strong empirical evidence that mechanisms such as the Gibrat process are *not* good descriptions of real city-size dynamics [see Cuberes 2004]. Abstract stochastic models have also tended to be “brittle”—they can generate the Zipf distribution, but they are “one-process-fits-all” and cannot generally account for the exceptions to or variations in Zipf that are observed in the data.

Deviations from Zipf

While the Zipf distribution offers a remarkably good fit for many nations, the fit is imperfect in many cases. In this paper, we will examine three countries which are particularly interesting with regard to their adherence to and deviations from Zipf. These three countries are: the United States, Russia, and France. All three countries provide excellent data on urban agglomerations. The United States represents a relatively good (though significantly imperfect) fit for Zipf, while France and Russia deviate in different, paradigmatic ways.

Before attempting to analyze the extent to which cities in different countries do or do not deviate from Zipf, we need to address the definition of a city. In this paper, we are interested in the city as a social and economic phenomenon, rather than

as a legal entity. Our unit of analysis is thus not the population within the official city limits, but rather the population of the urban agglomeration of which the legally incorporated city is often only a part.

Consistently defining an urban agglomeration is challenging [Le Gleau et al., 1996], but in the cases we have chosen, it is possible to derive reasonably satisfying definitions of urban agglomerations. The statistical agencies of both the United States and France have addressed this problem directly by developing various functional definitions of urban agglomerations, while Soviet central planning produced Russian cities that are clearly separated, compact and well defined. We will discuss the specifics of each of these cases in turn.

USA

The cities of the United States have generally been regarded as being very nearly Zipf distributed. Because of the sprawling nature of many US cities, and the high daily mobility of the US population, the definition of an urban agglomeration for the US has proven particularly difficult. Over the past several decades, the US Office of Management and Budget has worked with the US Census Bureau to develop a set of Metropolitan Statistical Areas (MSAs) which sought to capture this notion of urban agglomeration. As helpful as this conception was, however, it had significant limitations. For example, the definitions of MSAs depended in part on the desires of local elected officials – thus making them somewhat inconsistent from the standpoint of objective social science.

In 2003, however, the US Office of Management and Budget released a set of carefully, objectively defined data which it terms Core Based Statistical Areas

(CBSAs) or “Metropolitan and Micropolitan” areas [Federal Register, 2000]. This definition attempts to capture spatial and economic integration with a rigor that had not previously been attempted. The result is a consistently defined set of 922 cities. These cities follow the Zipf distribution fairly closely over a tremendous range: from greater New York City with 18.3 million people down to about the 800th city with a population of about thirty five thousand. The largest several cities are significantly smaller than the distribution would predict, yet the distribution generally fits with a power-law exponent of which is very close to -1.

For convenience in the analysis that follows, we will restrict this to a subset of the 250 cities with populations over 150,000. [Figure 3.2] This reduced set of cities looks very much like the full set, displaying a power-law exponent of -1.005.

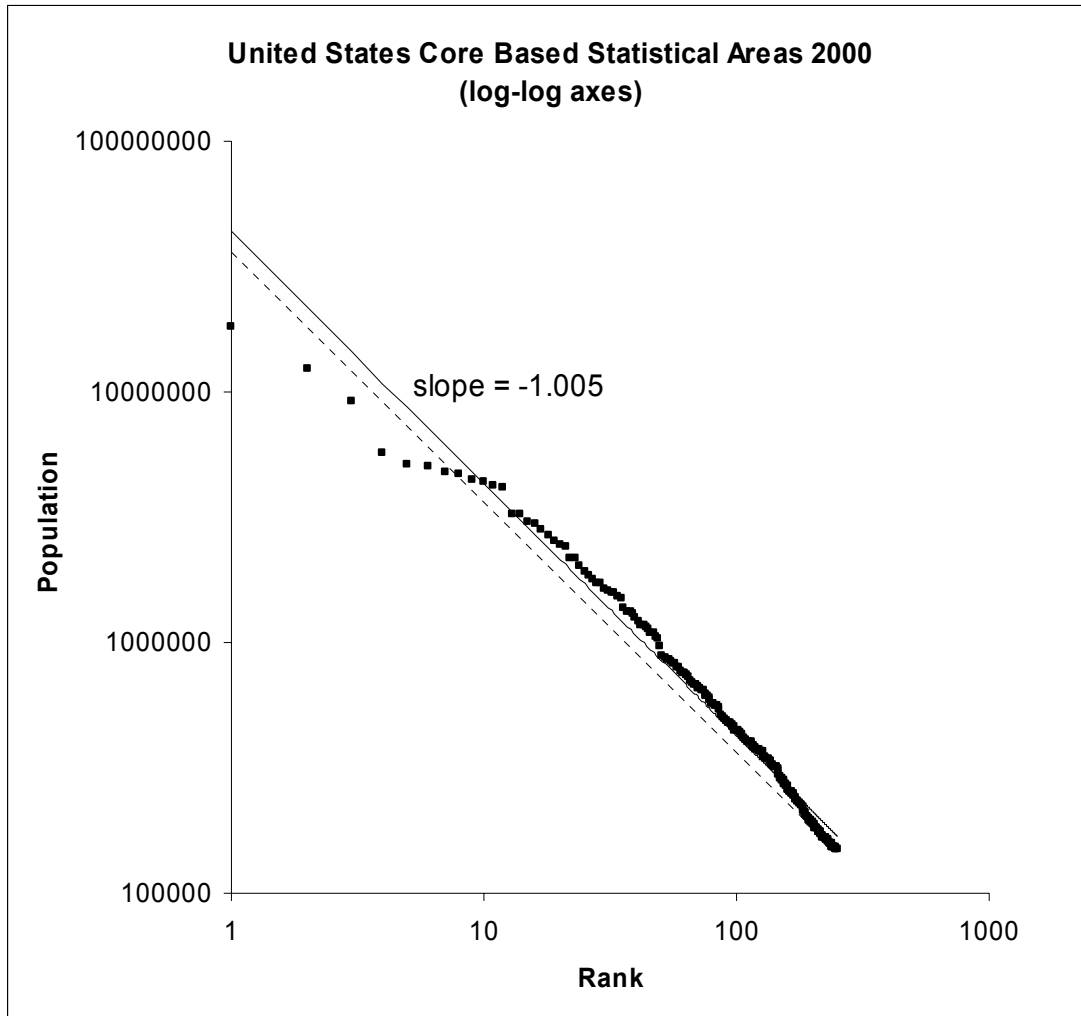


Figure 3.2: United States Core Based Statistical Areas, 2000.

In 2000, these largest 250 US cities collectively contained 220,227,293 people. If we construct a Zipf distribution with this many citizens distributed among 250 cities, it predicts a maximum city size of 36,098,839 (which is significantly larger than the observed 18,323,002 for New York) and a minimum city size of 147,384 (which is quite close to the observed 150,336 for Pottsville, PA which ranks 250th).

We can produce an objective measure of how well this constructed Zipf distribution fits the observed data by dividing the number of people which the Zipf rule misplaces relative to the data (the cumulative error) by the total population of the

cities. The cumulative error is calculated as the sum of the absolute values of the errors for each city divided by two (because each citizen which is in the wrong place is also missing from the right place). This procedure shows that the Zipf distribution misplaces 15% of the population of the largest 250 cities in the United States. We will refer to this measure as the total error.

While the overall error is well reflected by this measure, it does not give a sense of how the error is distributed. A sense of this distribution is given by the error at the median city. This is to say that we measure the error for each individual city $((\text{abs}(\text{Data}_i - \text{Model}_i)/2)/\text{Data}_i)$ and report the median of these values. This indicates whether the error is concentrated in a few large cities which fit poorly or is distributed throughout the range of the cities. We will refer to this measure as the median error. For the United States, the Zipf distribution produces a median error of 9.7%.

France

The French National Institute of Statistics and Economic Studies (INSEE) produces a variety of excellent data on French cities using various definitions. These include the municipality (*commune*); the urban pole (*pôle urbain* or *unité urbaine*); and the urban area (*Aire urbaine*).

Of these three ways of defining a city, the first and third are inappropriate for use in this analysis. The municipality definition is not useful because most major cities are composed of many municipalities. The municipality of Paris, for example, had a population of only about 2.1 million people in 1999. The urban pole of Paris, in contrast, was composed of 396 such municipalities and was home to over 9.6 million people [Chavouet & Fannouillet, 2000]. While the legal definition of a municipality

reflects historical and administrative realities, it tells us little about the urban agglomerations which we are studying.

Where the city as municipality definition is too restrictive, the city as urban area definition seems to be too broad. French urban areas are defined as those areas where at least 40% of the workers commute into an urban center which employs at least 5000 people [INSEE, 2004]. These areas can be very large, often many times the area of the urban pole. A major problem with this definition for our purposes is that this surrounding area mixes people who commute into the city center with people whose social and economic lives are not integrated with the city. This commuting based definition also creates the impression of rapid growth for many cities, not because the cities have changed significantly, but because French commuting patterns have been changing, with workers traveling increasing distances to work [Julien, 2001b]. French cities have therefore been expanding their areas of influence more rapidly than they have been growing in terms of employment, built area, or other measures of city size [Julien, 2001a].

The French definition of an urban pole strikes something of a balance between these two definitions. An urban pole is defined as a collection of contiguous communes in which more than half of the population lives in an area where buildings are separated by no more than 200 meters. This definition is thus a reasonably close approximation of the built up area of the city. However, because this definition includes whole communes which are only partly urbanized, it tends to over count the urban population at the edges of cities. Because the circumference of a circle increases more slowly than its area, this bias tends to inflate the size of smaller cities.

In an effort to avoid this problem, we adopt a slightly more restrictive definition of a French city, which we will call an “urban center”. Our definition follows the spirit of the one described by Le Gleau et al. [1996] while adapting it to better capture the dominance of Paris in the French urban system. Le Gleau defines an urban center such that, if a single commune within an urban pole contains more than half of the pole’s population, then this commune is the urban center. If the central commune contains less than half of the population of the pole, then it is agglomerated with the other communes of the pole which have at least half of the population of the largest commune. This definition has the effect of making the urban centers of France appear very nearly Zipf distributed – but it makes little sense as a definition of a city. Most notably, the central commune of Paris is much larger than any of the other 395 communes which make up the Parisian urban pole. This means that, by Le Gleau’s definition, the urban center of Paris is represented by only this one commune, putting its size at 2.1 million people (as compared to 9.6 million in the urban pole).

We retain Le Gleau’s concept of omitting the fringe areas by changing the criteria for agglomerating secondary communes, but refine it to avoid distorting large cities (particularly Paris). Under our definition, we agglomerate all of the communes in the pole which have a population greater than 20,000 people. Because communes tend to be of roughly uniform size, this is a reasonable proxy for density. We choose the number 20,000 because it is also the minimum size of a city in our dataset. Thus, any commune within an urban pole which would qualify as a city in its own right by virtue of its population of 20,000 is agglomerated into the urban center. This

definition eliminates the inflation of the urban periphery which is present in the urban pole definition while retaining the basic idea of a city as a contiguous built-up area.

The analysis that follows will use this definition of a French urban center.

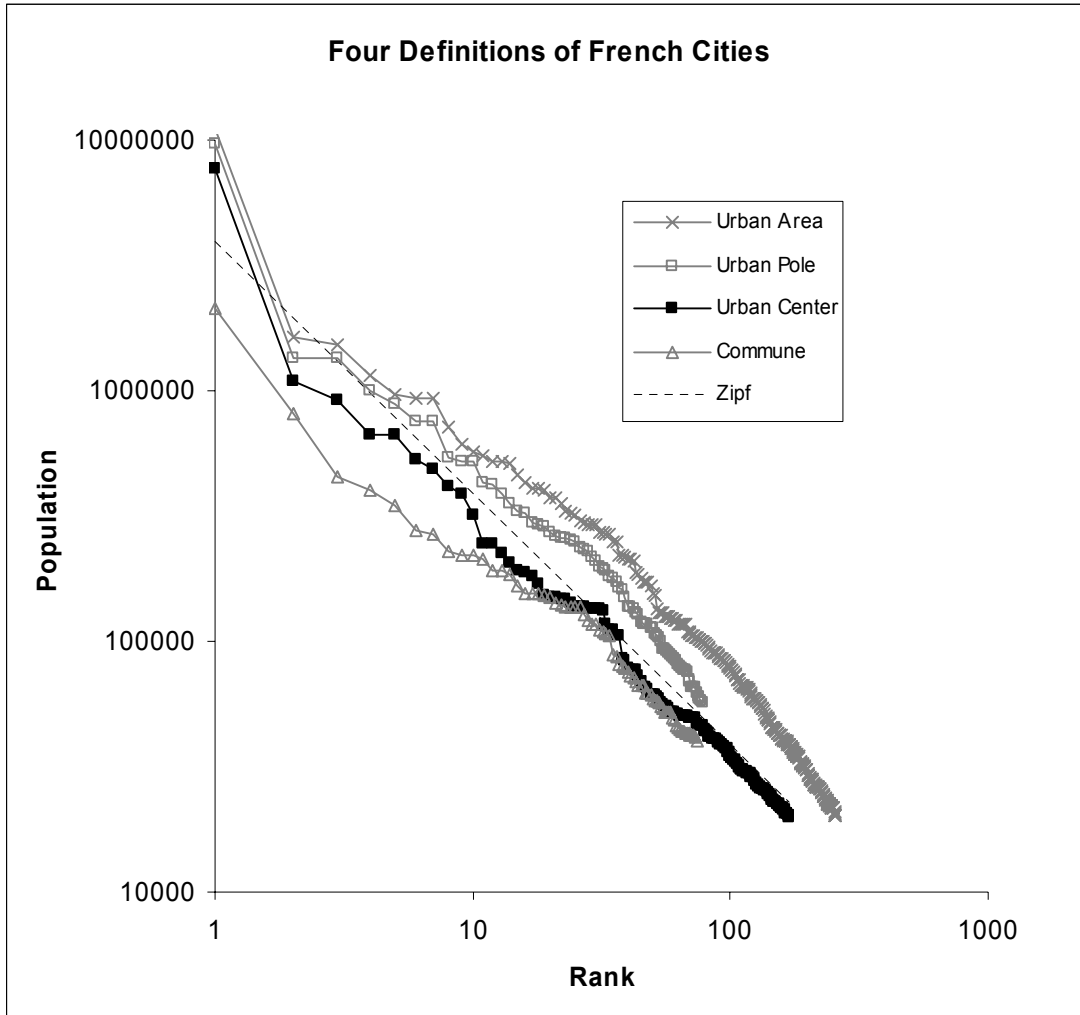


Figure 3.3: Four definitions of French city sizes.

The urban center data conforms fairly closely to Zipf, displaying an overall power-law exponent of -0.98. The primary deviations from Zipf are that Paris is about two and half times the size that the rest of the distribution would predict while the second agglomeration, Marseille-Aix-en-Provence, is about two thirds the size that the distribution would predict. The combination of these two factors makes Paris

about 7 times as large as France's second city – whereas the norm would be twice as large. Overall, the Zipf distribution displaces 17% of the French population, but this is largely due to the very poor fit of Paris. This is pointed up by the fact that the error at the median city is only 7%.

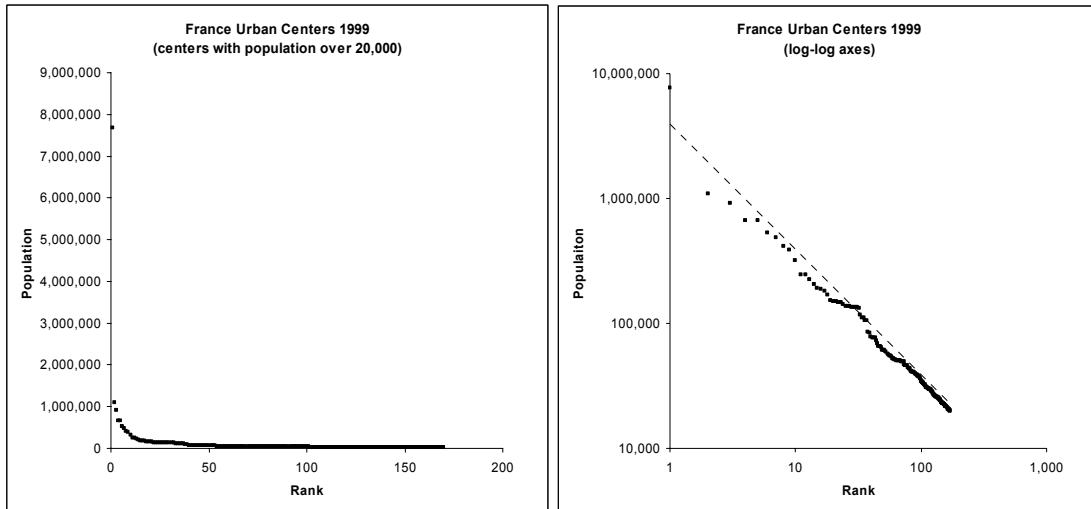


Figure 3.4: France urban centers, 1999

Because France is much less populous than the United States, its urban structure is also much smaller. Whereas the United States has about 900 cities with populations greater than 20,000, France (following the 1999 urban center definition) has only 170 cities above this size.

Russia

Unlike the United States and France, which both adhere closely to the Zipf regularity for all but their largest cities; the Russian city size distribution displays a distinct curvature on log-log axes over the entire range of its urban structure.

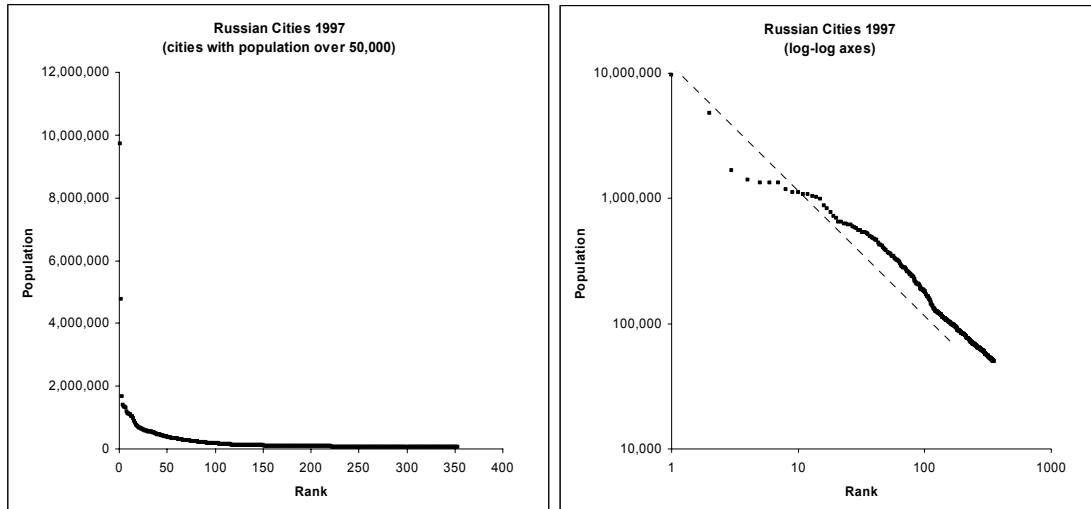


Figure 3.5: Distribution of Russian City Sizes

That the Russian urban structure is substantially different from that of the US and France is not surprising given the radically different physical, social, and economic environment in which it developed. Much of Russia’s urbanization took place during the Soviet period when internal migration was intensely managed by the central government. Soviet planners had various objectives in establishing cities including the extraction of natural resources, the occupation of territory which might be claimed by China, and the movement of industrial production away from the potential front with Western Europe. They pursued these objectives through policies of forced and incentivized migration, costly investments in infrastructure, and intensive subsidies to far flung cities in inhospitable locations [Hill and Gaddy, 2003].

The impact of this managed migration was to increase both the number and the size of cities in far flung parts of the Soviet Union. A basic reality of this system, which we will make use of in the modeling that follows, is that it made it easier to move down the urban hierarchy than it was to move up. A person living in Moscow might be assigned a job in a minor industrial center in Siberia, but a person living in that Siberian city would be unlikely to be assigned to Moscow. Apart from the forced

migration associated with the GULAG prison system and other, less punitively oriented assignments to work, the Soviet system relied on heavy subsidies and incentives to get people to move to smaller places. It also used a system of internal passports to insure that people could not find employment or move about freely in a city other than the one in which they officially resided. These systems insured that the smaller (and often colder and generally less hospitable) industrial cities of Siberia remained populated in spite of Russian citizen's inclinations to move elsewhere [Hill and Gaddy, 2003, Iyer, 2003].

Russian urban agglomerations are easier to define than their US and French counterparts because of the way that Soviet planners designed the Russian urban structure [Hill and Gaddy, 2003]. The desire to spread population over the vast territory of the Russian empire created large distances between cities while the planned nature of these cities reduced or eliminated urban sprawl in most cases. Because Russian cities tend to be distinct and compact, Russian city population numbers and urban agglomeration numbers tend to coincide, requiring the aggregation of suburbs with central cities only for Moscow and St. Petersburg. The data generated by the Russian census are therefore appropriate for our purpose without adjustment beyond the agglomeration of these suburbs.

The overall best fit power-law for this data has an exponent of -0.92 – a number close enough to unity that some authors have failed to remark on it. Our quantitative measure of error indicates that the fit between the Russian distribution and the Zipf distribution is similar to that for the US and France, misplacing 16% of the population (as compared to 15% and 17% respectively), but this apparent

similarity is misleading. This shows up in a median error figure of 17% (as compared to 10% for the US and 7% for France). While the US and France distributions are generally Zipf like, with departures in the largest cities, the Russian distribution is distinctly curved.

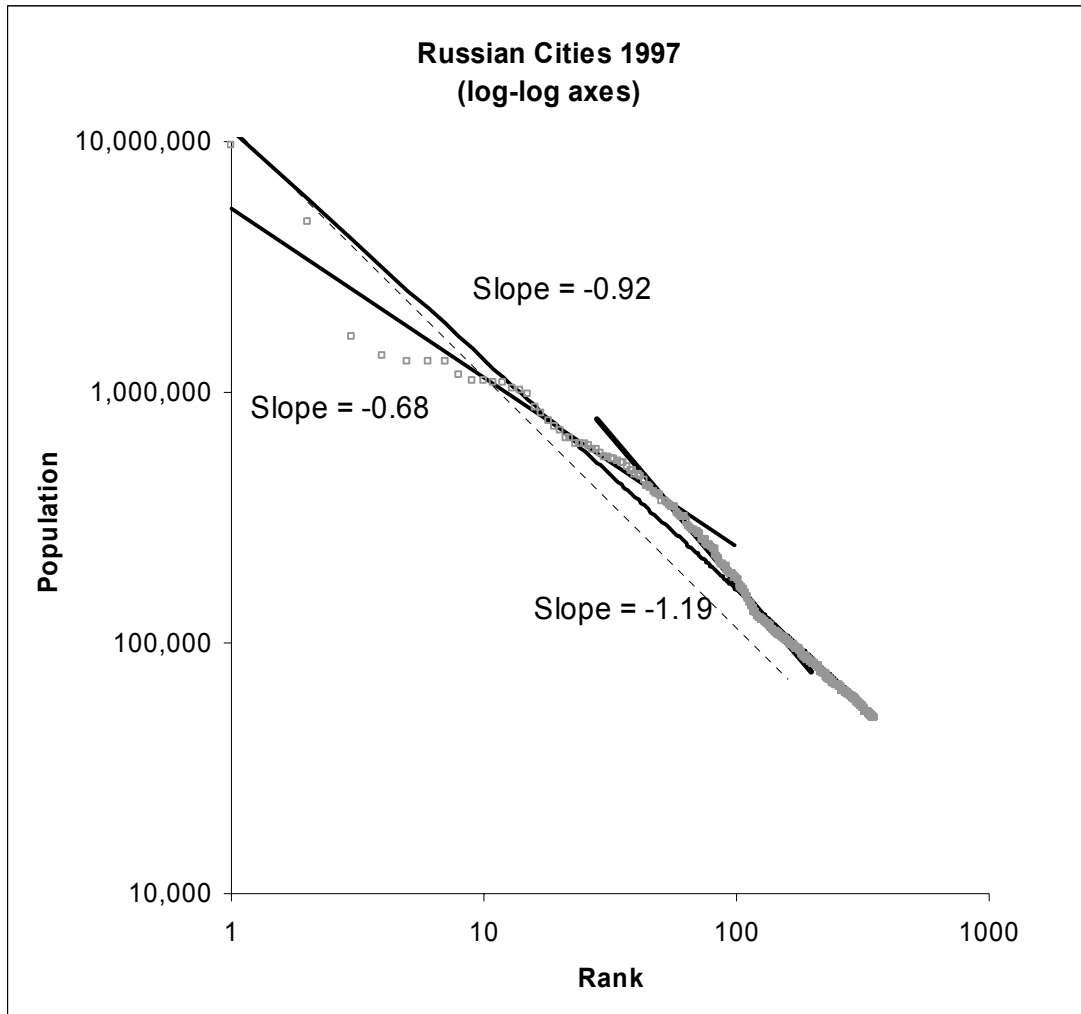


Figure 3.6: Curvature of the Russian city size distribution.

We can demonstrate this curvature by dividing the Russian city distribution into two parts and examining the exponents of the best-fit power-law which describes each part, measuring the power law exponent for cities larger than 500,000 separately from those between 500,000 and 100,000. These sets of cities display two distinct

exponents. The upper part of the curve has a slope of -0.68 while the lower part has a slope of -1.19. These slopes are significantly different with $p \ll 0.001$. Similar tests on data from the US and France yields slopes that are not significantly different.

By inspecting the graph [Figure 3.6] we can see that our cut off of 500,000 between the two groups is an arbitrary one and that the distribution of cities larger than 100,000 is better described by a curve which is concave toward the origin. In this sense, the Russian distribution departs from the Zipf distribution for all of the 161 cities in this range.

A Simple, Abstract Model: Jars and Beans

Model Description

In the sections which follow, we will attempt to explain both the tendency of urban systems to approximate Zipf, and the reasons why the various countries depart from it by constructing a model which is as simple as possible while capturing the essential features of the systems in question.

We begin with an abstract model which can produce remarkably good agreement with real city size distributions. This model is designed to explore the way in which power-law distributions can emerge from systems involving stochastic exchange. Because the abstract model does not itself contain plausible urban dynamics, we describe it in terms of “jars” (rather than cities) exchanging “beans” (rather than citizens). In the next section, we will extend the model in such a way that it demonstrates a plausible relationship to social and economic realities.

The rules of the abstract model are simple. The model begins with some number of jars each of which contains some number of beans. The jars interact in

random pairings. In each interaction, the jars exchange some number of beans ("the bet") equal to half of the beans in the smaller jar. In the base case, both jars have an equal probability of winning the bet. Once the winner is determined, the beans are exchanged and a new random pairing of two different jars is made. Finally, there is a floor size of 1 bean. If a jar of size 1 loses a bet, nothing happens and it remains at size 1. If it wins a bet, it wins a whole bean (rather than half a bean).

An important feature of this model is that it assumes that urban population is conserved. Whereas many others [e.g. Gabaix, 1999a; Fugita et al., 1999] have assumed that people freely enter and leave the urban system, we assume that once people have migrated to a city and have traded their rural skills for urban ones, they tend to remain in the urban system – migrating from one city to another in search of opportunities, but seldom returning to live in the hinterlands. In the simple model, this is reflected in a strict conservation law: beans are neither created nor destroyed, they simply move from jar to jar.

This model differs from other stochastic models typified by Gabaix [1999a] in that the growth rates of cities are not independent. These models generally depend on Gibrat process, wherein cities grow (or shrink) by random amounts. These random amounts are uncorrelated with one another and are drawn from the same distribution. In this model, growth rates are correlated (one city's gain is another city's loss). Also, growth rates depend on city size. When a small city faces a larger city, it faces a gain or loss of half its size, whereas the larger city faces a gain or loss which comprises a smaller fraction of its population. Small cities, therefore, face greater

size volatility than large ones, a fact that coincides with real world observation [Gabaix, 1999b].

Results from the Abstract Model

As simple as this model is, it can produce the Zipf distribution as well as some interesting variations on the distribution. If the model is run with the appropriate number of beans¹ for the given number of jars, it will approach the Zipf distribution regardless of the initial distribution of the beans between jars. Initializing the model with more beans than would be required to fill a Zipf distribution for the given number of jars produces instability in the top of the distribution with large fluctuations in the sizes of the largest jars, with the excess beans tending to float among the top few jars. Radical overfilling of the distribution tends to produce “jamming” at the top, where the largest jar ends up with the majority of the excess beans. Initializing the model with fewer beans than would be required to fill the Zipf distribution produces a curvature of the distribution, maintaining the power-law exponent in the lower tail and progressively lowering it in the upper tail.

Another general property of the abstract model is that the size of the bet is not terribly significant to the dynamics of the model. While it is important that the bet be related to the size of the smaller jar, the size of that fraction generally affects only the

¹ From the definition of the distribution, it follows that a certain number of jars requires a certain number of beans to fill the distribution. When the floor size (the size of the smallest jar) is one bean, the largest jar should contain a number of beans equal to the number of jars. The sizes of all the jars between the largest and the smallest are then given by the rank/size rule, rounding to the nearest whole bean. For example, for 100 jars, 516 beans are required to fill the distribution.

speed with which the system approaches equilibrium, not the nature of that equilibrium. There does, however, come a point where the bet is small enough that the lower tail begins to collapse, with smaller bets leading to faster collapse. This does not occur with a bet size of 50% of the smaller jar, and therefore is not an issue in the runs that follow. We will discuss this property in more detail in the next section where bet sizes are reduced to the point where tail collapse becomes an issue.

We can make a first analogy from this abstract model to urban dynamics by thinking of the jars as cities and the beans as groups of citizens. Each bean represents the number of citizens in the smallest city in the sample. Actual population data can therefore be translated for use in the jars and beans model by dividing the total population of the urban system by the size of the smallest city in the system. This translation means that the units of exchange in the model are the size of the smallest city. This coarse assumption leads to discontinuities in the lower tail of our graphs, but it produces some interesting results and we will subsequently refine it.

Population figures for United States cities can be inserted into this simple model to produce a distribution which bears a noticeable resemblance that which is observed. In the year 2000, according to the Census Bureau data discussed above, the US had 250 cities with population larger than 150,000 and these cities were home to a total of 220,227,293 people. We can translate this for use in the jars and beans model by dividing the total population by the size of the smallest city, giving 1,468 beans. We can then get a first approximation of the US urban distribution by initializing the model with 250 jars and 1,468 beans. Running the model with these parameters gives a fit which is suggestive.

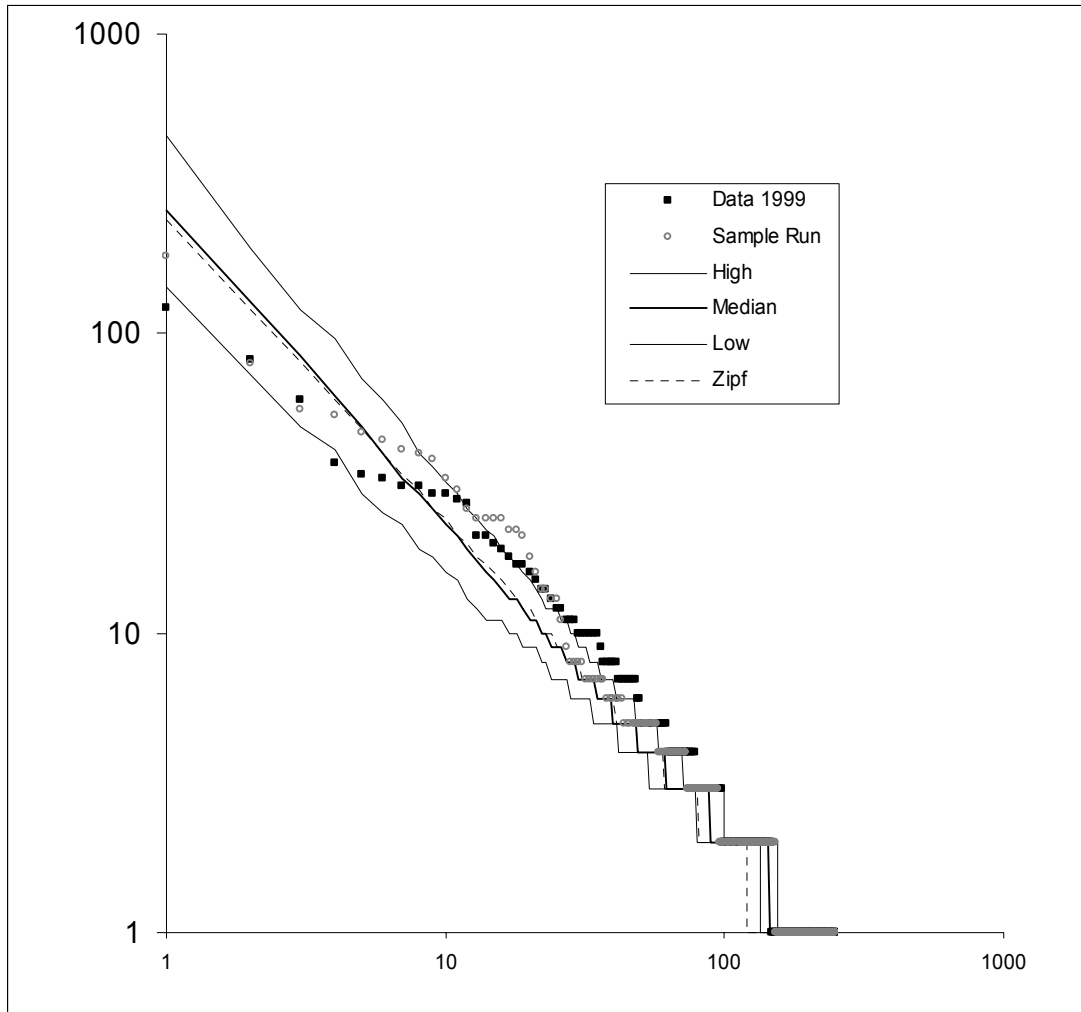


Figure 3.7: *Simple model output compared with discretized US data.*

Figure 3.7 shows the discretized version of the US data compared to output from 100 runs of the simple model using 250 jars and 1,468 beans. The heavier, central line on the graph indicates the median size for the city of each rank across all model runs, whereas the lighter lines represent a 90% confidence interval around this median. This is to say that the lower line represents the fifth largest value for that position over the 100 runs, while the upper line represents the 95th largest value. The US data does not fit precisely within this envelope, but it is not far off. The gray circles in the figure represent one of the hundred sample runs which is comparatively close to the data. Having observed that the model gets the gist of the distribution

right, we will return for a more careful analysis with more complex model in the next section.

Conducting the same exercise for France produces similarly provocative, but not entirely convincing results. Using our definition of an urban center, France has 170 cities with populations larger than 20,000 which collectively contain 22,386,598 people. We thus initialize the model with 170 jars and 1,119 beans.

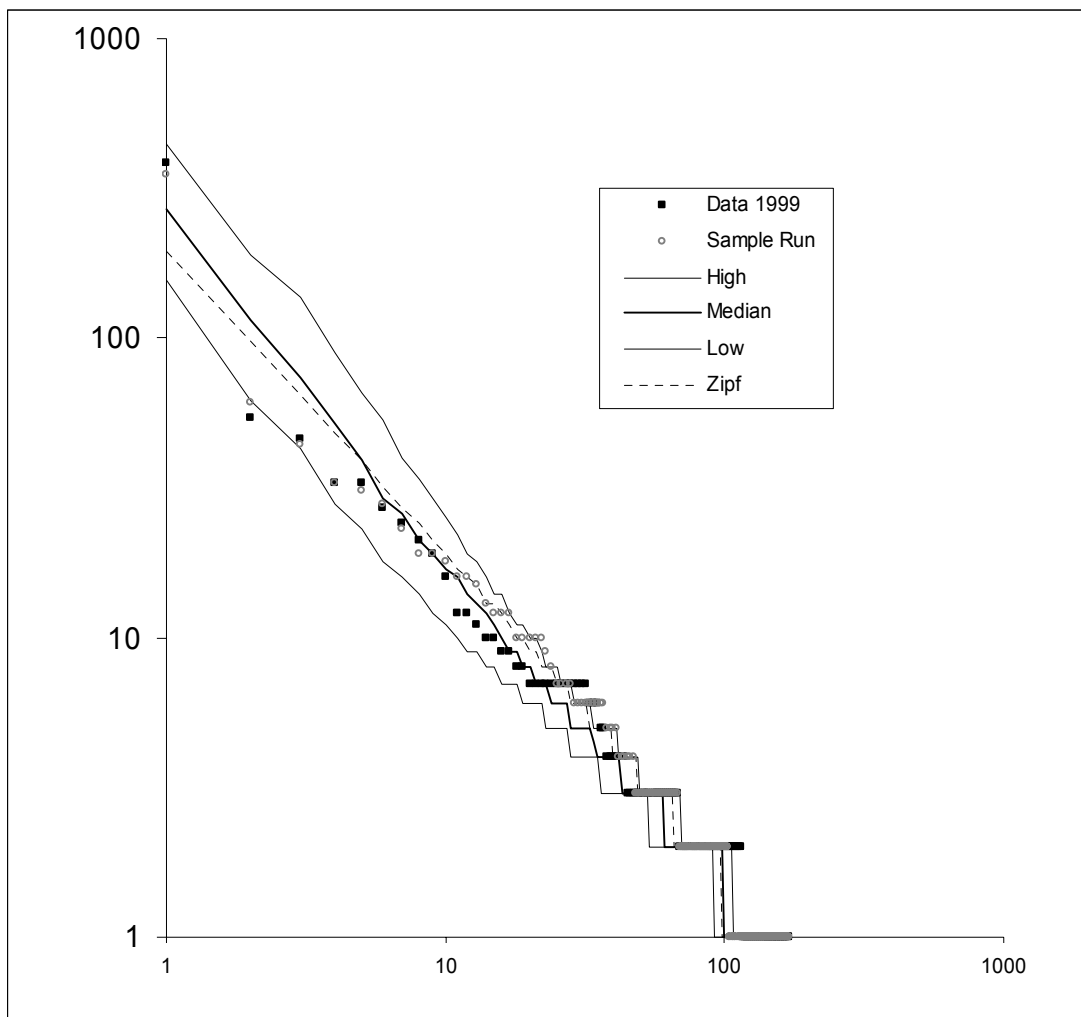


Figure 3.8: Simple model output compared with discretized France data.

Again, we see that the data generally fits within the range of model results.

We can see from the sample run that in a case where the first two cities are of the

proper size, the fit of the rest of the distribution is also very close. Though the simple model does not fully predict the primacy of Paris in the French urban system, the median model run does reflect an increase in slope in the top three or four positions. This is consistent with the notion that a small urban system with a relatively large population will tend to see disproportionately large cities at the top of its range.

Finally, we can obtain intriguing results for Russia by applying the model with a slight variation. In 1997 Russia had 161 cities with populations over 100,000 which collectively contained 70,282,100 people. This yields 703 beans in 161 jars.

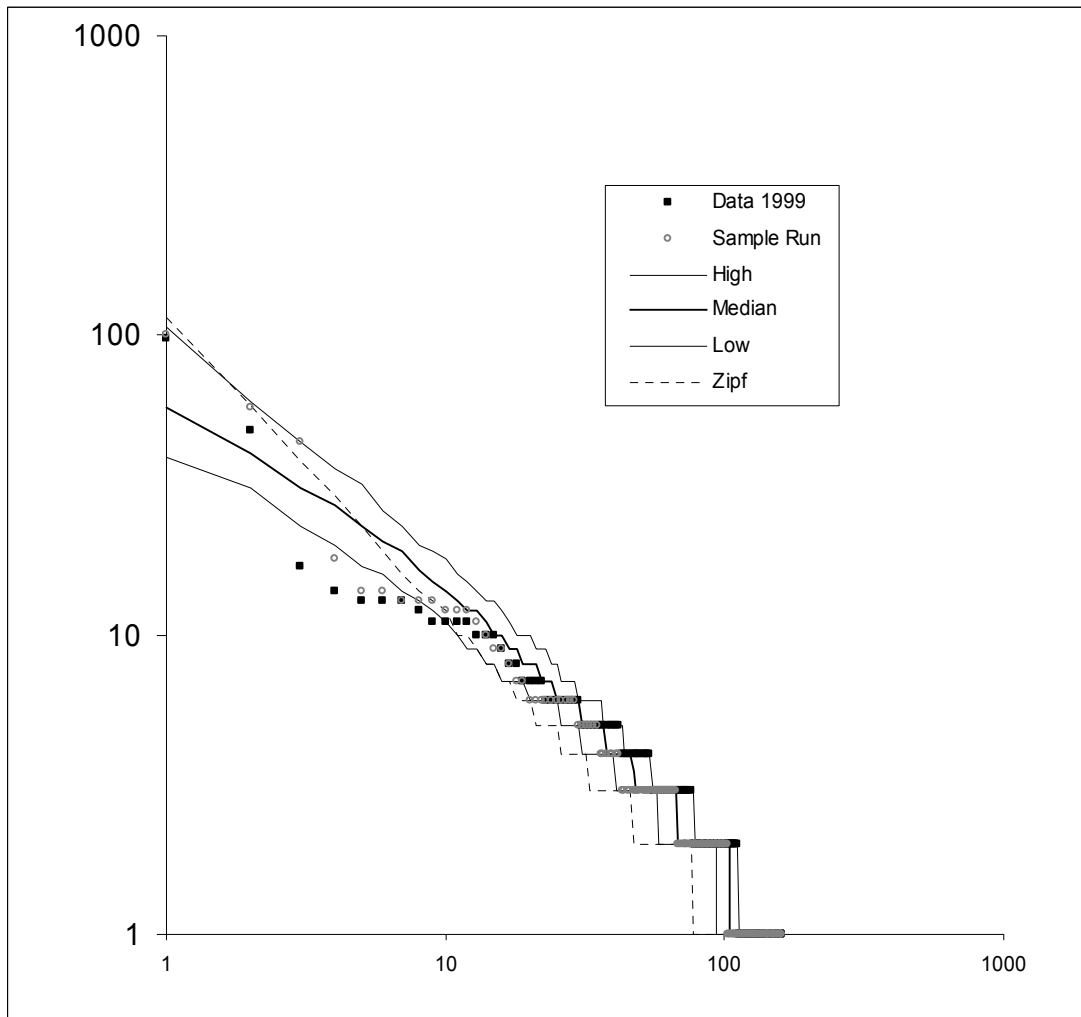


Figure 3.9: Simple model output compared with discretized Russia data.

Initializing the model with these values gives us a distribution which is concave toward the origin on log-log axes, but which has a somewhat different shape than we see in the data from Russia. If, however, we approximate Soviet era restrictions on internal migration by introducing a bias into the process, simulating the asymmetry in difficulty between moving up and moving down the urban hierarchy by giving the smaller city in each pairwise interaction a small advantage, the shape of the distribution comes to match the Russian case more closely.

Limitations of the Abstract model

While the abstract model offers a simple mechanism which creates distributions which look much like real city size distributions, it suffers from several serious limitations. First and most importantly, the dynamics of the model bear little resemblance to those of cities. Cities do not engage in tournaments where they flip coins for half of their citizens. Also, the floor assumption of the abstract model provides a subsidy to the smallest jars – in each interaction they stand to either remain unchanged or to double their number of beans. This mechanism tends to move beans from the upper parts of the distribution into the lower tail in a way that has no clear analog in the dynamics of urban migration.

Also, the simple model implies a highly unstable distribution, where the largest cities vary tremendously in size over time. This also implies a high churn rate, with cities changing rapidly changing their rank within the distribution. In the time scale that is required to achieve the power-law distribution, Chicago might change places with Peoria several times. This point highlights the fact that the abstract model has no place in it for differences in site suitability. Some places (natural ports, for

example) are simply better than others for large cities and any plausible model of urban dynamics should be able to reflect this fact.

A Richer Model: Cities and Citizens

Model Overview

To address these deficiencies, we will now introduce a richer model which comes closer to representing real urban dynamics. This model preserves and improves upon many of the desirable qualities of the abstract model while remedying most of its shortcomings. The richer model relies on the notion that a city has a short-term equilibrium size which balances economies of agglomeration (reasons to move into the city) with diseconomies of congestion (reasons to move out). A city can be thought of as being oversized if it moves above this equilibrium value and undersized if it moves below it. This short term equilibrium is subject to shocks which are based in the bounded rationality of citizens. The equilibrium reacts to these shocks over the longer term according to a lagged adjustment mechanism. Finally, the model introduces the concept of a core size below which it is not economically rational for a city to shrink.

Bounded Rationality

The concept of bounded rationality underlies the betting mechanism in the abstract model and provides us with guidance in refining it in terms of both its size and its "fairness". We can see the centrality of imperfect information in the model by assuming (temporarily) that all cities are at their equilibrium sizes. In this case, with each city is at its optimal size, perfectly informed and rational agents would have no incentive to move from one city to another because any move would leave their home

city underfilled and their new city overfilled – making the mover worse off. Any distribution could therefore become stable over time.

The citizens in our model, however, have imperfect information and bounded rationality. Some citizens, therefore, will move from city to city even at an “equilibrium” distribution of sizes. People are more likely to move from a more crowded city to a less crowded city, but the reverse is also possible. The size of the bet, then, relates to the degree to which the people's rationality is bounded (with a limit at perfect rationality, where the bet is always zero). The expected value of the bet remains at zero (i.e. is varies symmetrically around zero), so the bet can be said to be “fair”.

This principle of fairness is does not obtain in the abstract (jars and beans) model. In that model, the floor mechanism provides a significant subsidy to small jars. With 100 Zipf distributed jars, a bet size of 50% of the smaller jar, and a floor of one, about 1/3 of the jars face positive expected returns -- and the rest face negative expected returns. When the bet is decreased to 1% of the smaller jar, as it is in the model runs that follow, only the single smallest jar can be expected to be within 1% of its floor, and the amount that it stands to win is so small that its effect on the overall distribution can be safely ignored.

The size of the bet, therefore, is a parameter of the model. It represents the degree to which the rationality of the citizens is bounded – the percentage of the citizenry which will move between two equally attractive cities because they mistakenly believe that life will be better in the other city. As with the abstract model, the primary effect of changing the size of the bet is to change the speed with

which the system moves. However, there comes a point where the bet is small enough that very few small cities face positive expected returns. Over the long run, this leads the lower tail of the distribution to sag (i.e. to bend toward the origin) and produces long oscillations in the extent of this sagging. Such bending toward the origin in the lower tail is not observed in real data.

This problem, which appears to be an artifact of the model, can be overcome by introducing a small amount of growth into the system. When all cities grow by a tiny amount each round, the lower tail restabilizes near a slope of -1. The amount of growth does not need to be carefully tuned to achieve this result. The growth rate needs to be enough to keep the tail from sagging yet small enough that the system can “digest” the new citizens. Within that range, the growth rate can vary by an order of magnitude without significant impact on model output.

Lagged Adjustment

Cities adapt to the shocks imposed by the bounded rationality of their citizens through a lagged adjustment mechanism. The basic idea is that if the city grows above its equilibrium size it will become congested. If it remains congested for long enough, however, the city will adapt. Firms will move in to hire idle workers. New housing, roads and facilities will be built. Once these things happen, the city can comfortably accommodate more people than it did before -- its equilibrium size has increased. Similarly, if people move out and stay out for long enough, firms will leave and infrastructure will deteriorate, leaving the city able to comfortably accommodate fewer people than it once could.

Adding an adjustment lag does not change the dynamics of the model, but does impact the rate at which individual cities change size over time and therefore the rate at which the distribution changes. Because the parameters of this mechanism only influence the speed with which the model changes, and we are not attempting to calibrate the model to real time, we will not dwell on the lagged adjustment mechanism. Any mechanism which retains the “fair” quality of the bet from the simple model and does not introduce excessive noise into the model will produce similar results.

Inherent Suitability

A further requirement for the model is to account for the influence of geographic suitability and the persistence of great cities. We accomplish this by positing a more “rationally” determined core size which is only one component of the observed size.

We begin with the assumption that only some fraction of the population of a city is tied to the city’s specific geographic location. Chicago, for instance, is in a unique location to serve as a port for a huge section of the American Midwest. Many of the jobs in Chicago need to be located exactly where they are geographically -- at the base of Lake Michigan. Many other jobs in Chicago, however, do not have to be in that location. But they do have to be somewhere. We thus divide the population of a city into a core population, which is dependent on the city's geographic location and is subject to more or less rational and deterministic microeconomic rules for its size, and a floating population, which is subject to the mechanisms of the model.

A recurring problem for theorists of city sizes has been that models with economic content [Fujita, Krugman, et al., 1999] predict distributions which look quite different from those that are actually observed. The model presented here solves this problem and dovetails nicely with such models by freeing them from the need to predict a Zipf like distribution. A model like Fujita & Krugman's is probably well suited for predicting the core sizes of cities. These core sizes should be much more readily subject to "rational" analysis. The core sizes, however, are not the end of the story and are not the sizes that we see. The sizes we observe are based on the sum of the core size and the size of the floating population that can potentially live elsewhere.

Remarkably, the presence of some cities with higher floors (i.e. larger core sizes) does not change the basic dynamics of the model. It still produces Zipf and the aforementioned departures from Zipf. However, the cities with higher floors tend to stay in the upper part of the distribution, thus reflecting the persistence of major cities which we observe in the real world.

An analogy to a cake with icing is a useful way to visualize the relationship between the core and observed distributions of city sizes. The core distribution is the cake, while the floating population is the icing. All that we observe in city size data is the height of the top of the icing. While the cake of the core size distribution might be rather lumpy and vary depending on economic and geographic structure, the icing of the floating population flows smoothly over the cake and finds its level. In this case, the attractor is not flat, as it is in the case of a physical cake, but rather follows the shape of the Zipf distribution and its related departures as outlined above.

Because this study is concerned with the overall shape of the various city size distributions, it is sufficient that adding heterogeneous core sizes does not change the distributions that emerge from the model. The simulations that follow use uniform core sizes equal to the size of the smallest city in the system, but the results would not be changed if a more complex or dynamic core distribution were used.

Heterogeneous core sizes would have testable implications for the volatility of city sizes over time, but that is beyond the scope of this paper.

While we generally treat core sizes as exogenous to the model, it is easy to imagine variants where they would be endogenized. If such a model involved preferential attachment (i.e. new core firms are likely to locate near existing core firms), then we might expect core sizes to be power-law distributed [Axtell and Florida, 2001]. The important point, however, is that the shape of observed size distribution is independent of the shape of the core size distribution.

Results from the Richer Model

USA

This richer model produces a fit for United States core based statistical area data which is significantly better than the Zipf approximation. The only significant parameters in this model are the number of cities with populations over 150,000 (250), the number of people in these cities (220,227,293), the rate by which each city grows at the end of every round, and the fraction of the smaller city which will serve as the bet. The first two (cities and citizens) are given by the data. The growth rate requires rough, order of magnitude tuning which has little impact on the outcome so long as it is within a reasonable range. In the runs that follow, we use a growth rate

of 0.000005 in each round. The bet size alters the degree of variance between runs, but does not have a noticeable impact on their median outcome. The model thus has no significant free parameters.

We begin the simulation of the United States city size distribution with 250 cities and a reduced population of 50 million citizens (about 1/5 of the actual population) distributed evenly between the cities. The population at this starting point is not significant so long as it is small enough to allow the model to approach equilibrium before the full population is reached. We run the simulation forward with each city growing by a small amount ($1/20,000^{\text{th}}$) at the end of each round, stopping when the population reaches the year 2000 total urban population of 220,227,293.

This growth rate requires some tuning in order to be large enough to prevent the collapse of the lower tail while being small enough to allow the upper tail to assume a mature shape by the time the model population reaches the observed population. The need to tune the growth rate seems, however, to be an artifact of the model. For the sake of simplicity, we begin these simulations with a uniform distribution and with a fixed number of cities, whereas in the real world the urban system is always in the neighborhood of the Zipf distribution, with the number of cities increasing along with their populations. Such a growth pattern is supported by history [Zipf, 1949; Pumain, 2004] and emerges from certain theoretical formulations [Simon, 1957; Gabaix, 1999a; Axtell & Florida, 2001]. When the initial state is close to Zipf, the growth rate becomes much less critical. It needs to be great enough to prevent the collapse of the lower tail, but more rapid growth is not a problem because the system does not need to produce major structural changes.

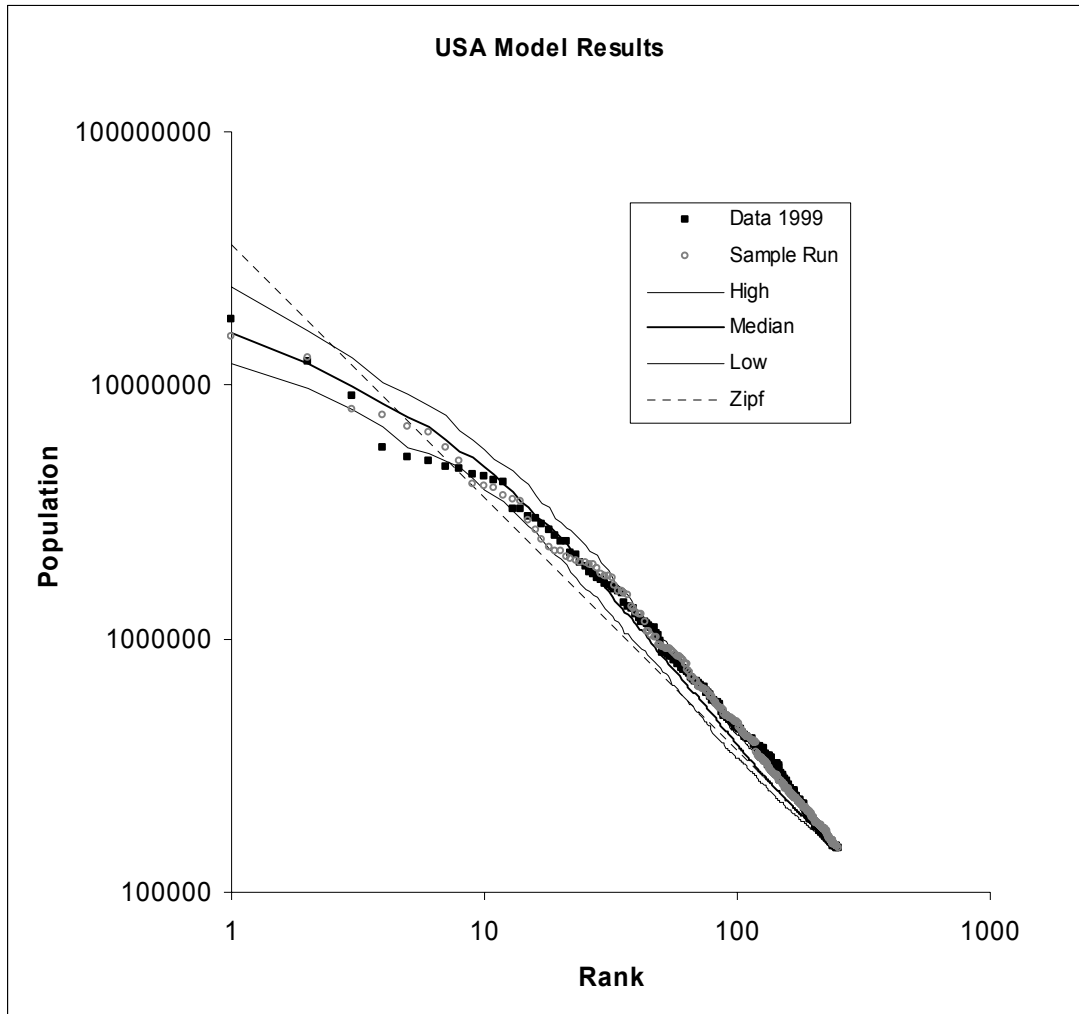


Figure 3.10: USA model results.

When we apply the model to data from the 2000 US census, we can achieve fits that are significantly better than Zipf. Whereas the Zipf distribution predicts US city sizes with an error of 15%, the median run of this model predicts the sizes with an error of 6.0%. The sample run shown in figure 3.10 achieves an error of 4.5%. The median error of the median run also reflects a better fit for the data than Zipf – the model creates an error of 4.5% in the median city compared to 9.7% for Zipf.

France

As discussed above, France is generally characterized by a Zipf distribution with Paris being considerably larger than the rest of the distribution would predict. Though the abstract model produced results which were consistent with French data, the case that we observe – with Paris seven times larger than Marseille – is an unusual one, occurring in less than 5% of model runs. The richer model performs considerably better in this respect.

It should be noted, however, that this vast improvement is in part due to the way that the model is run. In the US run, it is possible to use a constant growth rate, stopping the model when the model population becomes equal to the observed population. This procedure is complicated for France because the model takes considerable time to grow Paris into the prominent position which it occupies in the actual French urban structure. Any growth rate large enough to prevent the collapse of the lower tail causes the total population to be reached before the model has had time to grow Paris to its full size. We therefore begin the model with approximately 90% of the total population of France and run it forward until Paris has reached 90% of its actual population. We then introduce growth at the same $1/20,000^{\text{th}}$ rate used for the US simulation and run until the model population is equal to the French population.

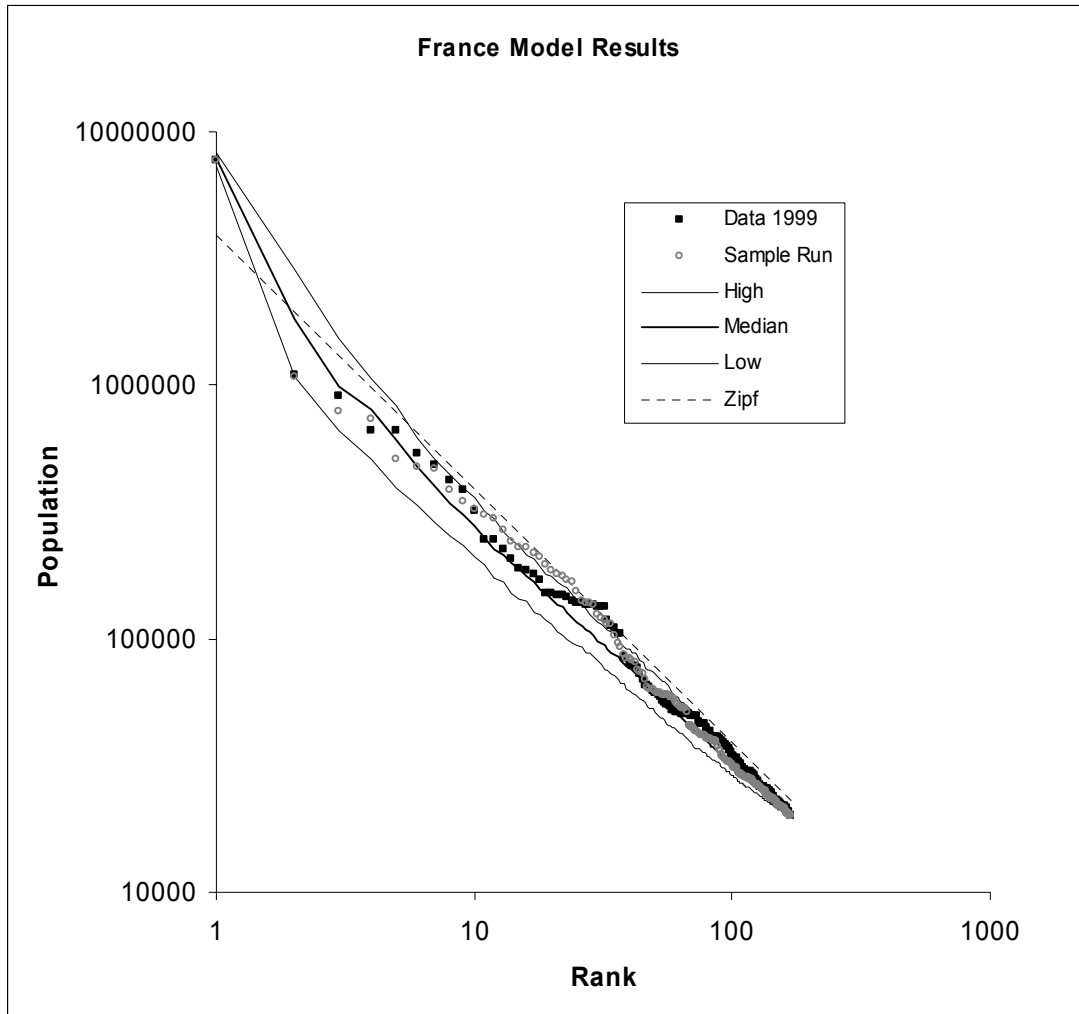


Figure 3.11: France model results.

Whereas Zipf produced a total error of 17% and a median error of 7%, the model produces a total error of 5.3% and a median error of 3.3%.

Russia

To simulate Russia, we initialize the model with 161 cities, a population of 70,282,100 in these cities and a floor of 100,000 (the size of the 161st city). As with the simple model, we introduce a bias into the migration probability to simulate the effects of internal movement restrictions. The degree of this bias is a free parameter of the model.

Given that the model does attempt to represent the urban system in actual space and time, it is not possible to calculate this movement bias using actual data. Because it is the only free parameter in the model, however, we can calibrate it by comparing model results to the observed data. We obtain a good fit by assuming a bias of 0.25% in favor of the smaller city in each pairwise interaction. That is to say that, in each interaction, the probability of the larger city receiving the migration (winning the bet) is 49.75% while the probability of the smaller city receiving the migration is 50.25% ($P_{\text{larger}} = 0.4975$ and $P_{\text{smaller}} = 0.5025$).

A side effect of this bias toward the smaller city is the elimination of the phenomenon (or model artifact) of the collapsing lower tail. In the presence of this bias, the model behaves virtually identically in the presence or absence of population growth. In light of this, we omit growth from the Russian model runs.

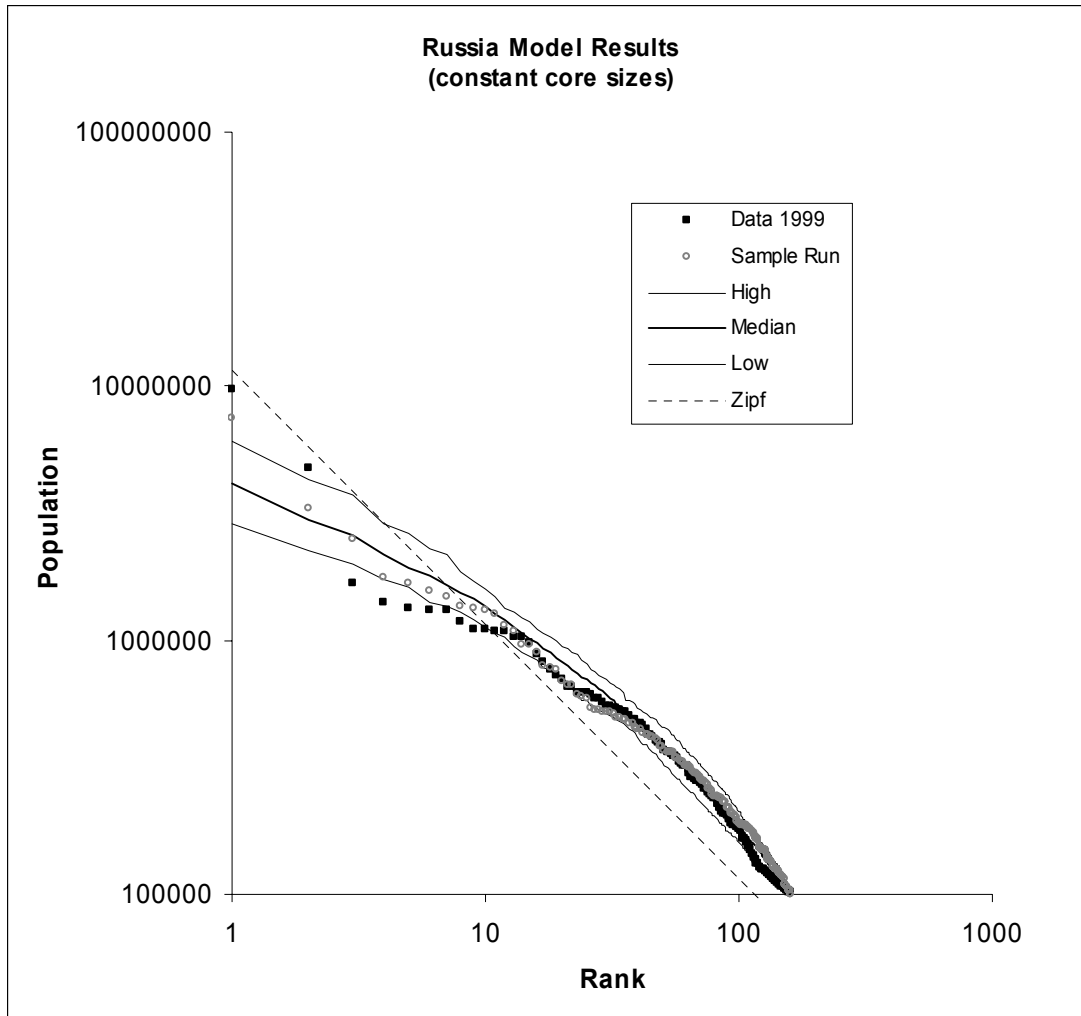


Figure 3.12: Russia model results with constant core sizes.

When run with these parameters, the model captures the basic shape of the Russian city size distribution, but misses the primacy of Moscow and St. Petersburg. These cities each have played unique roles in Russia’s economic and political history serving as capitals of highly centralized political systems under both the Czars and the Soviet system. St. Petersburg is also unique in serving as European Russia’s only ice free port. Given their centrality in the Russian economy, it is perhaps reasonable to treat them differently than the other cities of the nation.

The continuing pressure of internal immigration on these cities – even in the face of falling population in Russia generally [Iyer, 2003], indicates that these cities remain at or below their equilibrium size in the collective mind of the Russian people. In terms of our model, we can say that these factors have led these two cities to have floor sizes which are much larger than the other cities of the system. We can incorporate the unique economic and geographic appeal of these two cities by assigning them a higher floor size than the others. While we leave the floor size of the rest of the system at the size of the smallest city (100,000 people) we will move the floors of Moscow and St. Petersburg to 90% of their 1997 population (90% of 9,735,900 and 4,779,000 respectively).

We observed earlier that introducing heterogeneous floor sizes alters the stability of individual cities but does not change the shape of the overall distribution *unless* floors are set so high as to make a city “protrude” from the distribution. In this case, we are saying that political and geographic forces have caused the core sizes of Moscow and St. Petersburg to protrude from the Russian city size distribution.

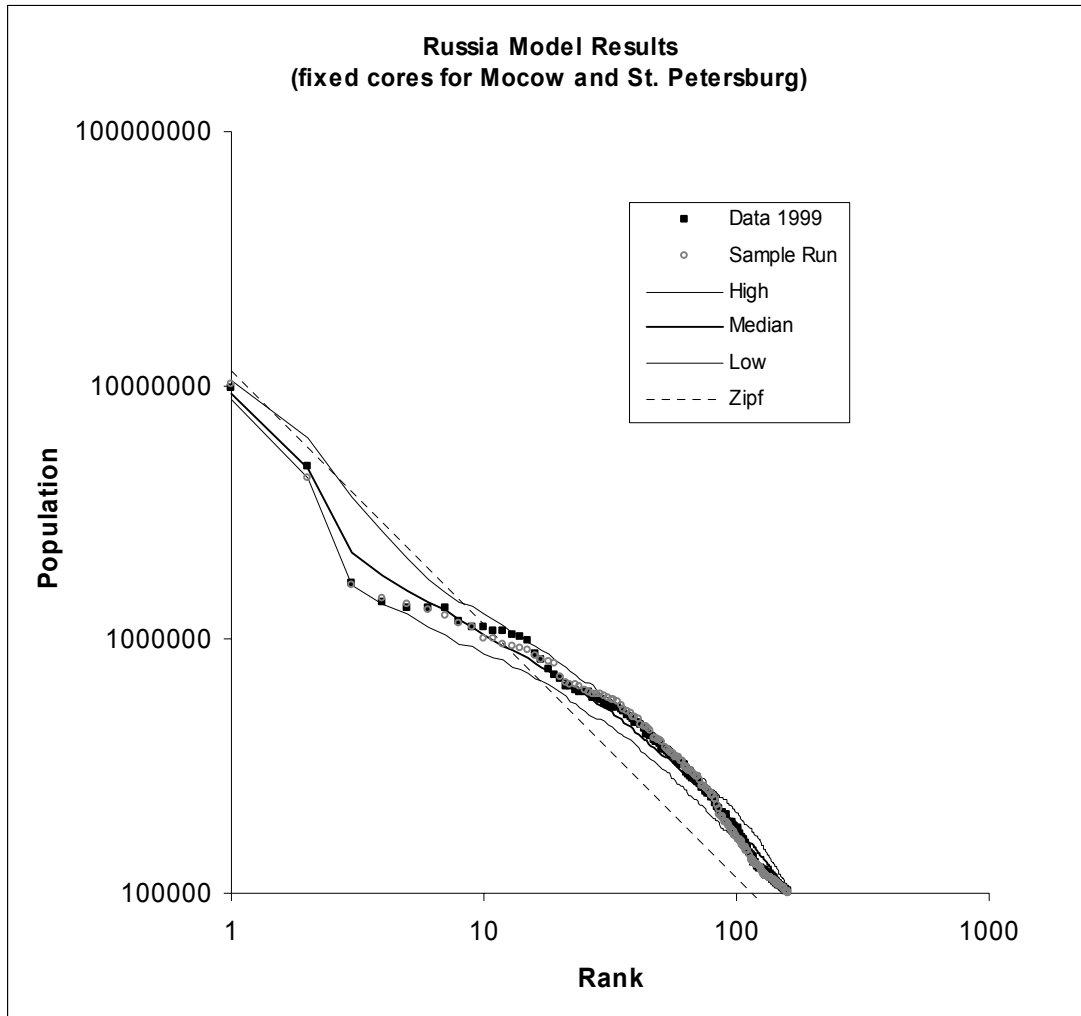


Figure 3.13: Russia model results with larger cores for Moscow and St. Petersburg.

When we incorporate these larger core sizes for Moscow and St. Petersburg into the model, it produces an excellent fit for the data. Overall, the median model run misplaces only 3.25% of the population, which is much better than Zipf, which displaces 12.5%. The error at the median city similarly drops yet further to 2.5% as compared to 16.5% for Zipf.

Limitations

While this model displays a good deal of success in reproducing the distribution of city sizes in the United States, France, and Russia, it is important to recognize several important things which it does not do.

While this model predicts the overall shape of the urban distribution for various countries, it does not predict the movements of particular cities within that distribution. In the simulations presented here, cities display unrealistic levels of volatility. We present the core size mechanism as a method for taming this volatility, however, we do not attempt to model core sizes. The important point here is that the model will produce the same distributions given an extremely broad range of core sizes including, we suspect, core sizes which are compatible with observed levels of volatility.

A second, related, limitation of the model is that its current formulation does not lend itself to calibration to real time. Real urban systems generally expand simultaneously in both population and number of cities, whereas we hold the number of cities fixed. We believe that this assumption, though unrealistic in the long term, can yield insights in the shorter term by keeping the model simple enough that it can be intuitively understood.

A third limitation is that the model uses a simple but highly unrealistic interaction network. Cities in the model interact randomly, regardless of their size or location – indeed, location is not represented in the model at all. The nature and stability of the equilibria are changed by different interaction regimes and it remains to be demonstrated that a realistic interaction regime would produce the same results.

Discussion

Our adaptive agent framework has allowed us to design and explore a simple framework for understanding city size distributions which, in spite of its extreme simplicity, is able to generate close approximations of the actual city size distributions for the US, France, and Russia. Even though the system is simple, it is hard to analyze because of the high degree of interaction among its parts. Previous attempts to explain the distribution in general have gained analytical tractability by assuming independence of the growth rates of cities. While it is possible to generate the Zipf distribution using such assumptions [Gabaix, 1999a] it is hard to imagine how the departures that we have reproduced could be derived without something akin to the interaction regime that we describe.

Is Population Growth Needed to Preserve Stability of the Lower Tail?

Our results for the United States and France depend on population growth to prevent the collapse of the lower part of the distribution's tail. As we have discussed above, this is not unrealistic since the urban population of both nations has been growing since the industrial revolution and continues to grow. Even though the overall population of France has stabilized, the French population continues to urbanize, leading to continued growth in French cities [Julienne, 2001a].

The fact that virtually all nations have growing urban populations makes it impossible to determine empirically whether growth is actually required in order to stabilize the distribution of smaller cities, or whether this is an artifact of the model. This is particularly true because the effect of population growth on the stability of the tail is non-linear. Sufficient growth stabilizes the tail, but there is a broad range of

growth rates that create no additional discernable effect on the behavior of the model. This makes it difficult to generate testable hypotheses which could be used to decide whether the tail collapse that we observe in the model might happen in the real world in the absence of growing urban populations. We believe that a more elaborate version of the model, where space is explicitly treated and the interaction between cities is based on both their sizes and the distance between them may shed light on this issue. We will discuss this elaboration in the succeeding section.

Implications for Developing Nation Megacities

One of the more interesting policy relevant insights generated by the model is that the primacy of Paris (and, by extension, other disproportionately large capitals) might have more to do with the number of small cities than it does with the nature of the large city. Previous efforts to explain urban primacy [e.g. Ades et al., 1995] have tended to focus on the political economy of the capital as the reason that it grows disproportionately large. These theories would attribute the massive size of Paris to the highly centralized nature of the French political system and the fact that it is “the capital of everything” including politics, finance and culture, for the nation. This contrasts with the United States where the political capital (Washington) is different from the finance capital (New York and to some extent Chicago) and the cultural capital (which one might argue is split between New York and Los Angeles). Our model allows for such theories – we invoke this kind of reasoning to explain the size of Moscow and St. Petersburg in Russia – but the model suggests that this kind of explanation may not be strictly required to explain the size of Paris. While the central role that Paris plays in French political, economic, and cultural life undoubtedly does

endow it with a substantial core size, it is not clear that this role requires it to be as large as it actually is.

The stylized result from the model is that a country with a large population and relatively few cities will tend to produce a Zipf distributed population in all but the largest city (or few cities) with the “overflow” population collecting at the top of the distribution. In this interpretation, the centrality of Paris guarantees that it will be that largest city (rather than some other city), and perhaps also ensures that France has a sole primate city, in stead of the small handful that can also emerge from the model.

While these other factors place Paris at the top of the French urban hierarchy, our framework suggests that its actual observed size has more to do with the large number of people and small number of cities in France. We should note at this point that the direction of causality is not entirely clear. Is Paris large because there are so few cities, or are there so few cities because Paris is so large? On the one hand, French towns may be prevented from growing to become integrated into the urban structure by regulations protecting agriculture, by preferences embedded in French culture, by peculiarities of French geography or history, etc. On the other hand the lure of Paris may be depriving small places of the population that they would need to grow up to the point where they could enter the urban hierarchy. While arguments can be made in both directions, it seems likely that both of these factors are at work.

The critical point here is that our framework suggests that there is a definite relationship between number of cities and number of people. It suggests that, in a situation where there are few cities and many people, the “excess” population will

tend to concentrate at the top of the system, forming megacities. This has real implications for urban planning in the developing world.

There is an extensive literature on developing nation megacities and their attendant problems which we will not attempt to survey in depth here. This literature is summarized in various global NGO publications including [UN-Habitat, 2004, UN-Population, 2001]. Megacity related policy challenges involve growth management and the provision of adequate infrastructure for a rapidly growing population. Failure to meet these challenges can create disastrous situations in the areas of environmental protection, public health, and human development and can lead to social unrest, political instability and violence.

Bugliarello [1999], summarizes the problems facing developing nation megacities as follows:

- Explosive population growth.
- Alarming increases in poverty that contradict the reasons why a megacity attracts.
- Massive infrastructure deficits in the delivery of telecommunications services, the availability of transportation, and the presence of congestion.
- Pressures on land and housing.
- Environmental concerns, such as contaminated water, air pollution, unchecked weed growth due to the destruction of original vegetation, and overdrawn aquifers.
- Disease, high death rates, drug-resistant strains of infection, and lethal environmental conditions.
- Economic dependence on federal or state governments that constrains the independence of megacity administrations.
- Capital scarcity, the factor that shapes the economy of the megacity and aggravates its other problems, from infrastructure to environmental deterioration.

Our analysis presents a reason to expect the emergence of megacities such as Sao Paulo in Brazil, Dhaka in Bangladesh, and Jakarta in Indonesia. These countries generally have highly centralized governments and severely constrained capital availability. These factors make it very difficult for their urban systems to expand in terms of number of cities at a rate that bears any resemblance to their rates of population growth and urbanization. Developing nations are therefore left with a small number of cities and a large urban population. While the first move from rural to urban life may be from the countryside to a nearby city (a tendency that would tend toward balanced urban growth) our model suggests that the next step of inter-urban migration will tend to concentrate the urban population.

The model further suggests that efforts to encourage migration from the first tier cities to middle sized cities are not likely to succeed over the long term. A government hoping to stem the growth of a primate city would do better to focus limited resources on providing the infrastructure and economic base which would allow large towns to become full participants in the urban system – thus expanding the number of cities and thereby reducing pressure on the capital.

Implications for Russian Urban Structure

The odd nature of the Russian urban structure appears, in light of this model, to be the result of two factors: a large urban system relative to its population and movement restrictions which have historically biased movements toward smaller cities. Unlike the urban structures of the US and France, the Russian urban structure was not created by free mobility and free markets. Soviet central planning created, instead “a structure of production – location, capital, employment, materials, energy

use, etc. without any regard for economic opportunity costs, in an environment free of economic valuation.” [Ericson, 1999]

The result of this non-market resource allocation was an extensive urban structure with post-Soviet leaders have continued to work hard to preserve through subsidies and other measures. For a host of ideological and security related reasons, Soviet central planners aimed for relatively even dispersal of cities of fairly uniform size while at the same time creating a highly centralized system of power [Demko & Fuchs, 1984] These factors contributed heavily to the creation of the odd urban structure that we see today.

One of the major Soviet era policies used to maintain this sprawling urban structure was a system of permits which were required for one to move from the hinterlands into an industrial center, and from a smaller industrial center to a larger one. This policy may be likened to biasing migration toward the smaller city in our model. While these policies are officially no longer in place since the fall of the Soviet Union, traces of them remain – particularly with regard to migration into Moscow and St. Petersburg. President Putin remains committed to avoiding Siberian “ghost towns” at almost any cost and many subsidies to these towns are in place even at the present time. [Gaddy & Ickes, 2002]

The model’s success at reproducing the Russian urban structure invites speculation as to what Russian planners might expect as the Soviet policies and their aftereffects fade. Because the model is not spatially explicit and is not calibrated to real time, we can offer only broad predictions in this respect, but can offer them with some confidence.

If we assume Zipf as the basic attractor for the urban system (an assumption which is consistent with our results) we would expect to see continued growth of Moscow and St. Petersburg. Though these cities are not far below the line predicted by Zipf, the large core sizes given them by their economic and political centrality mean that they have a relatively small “floating” population. While longer term social and economic forces may reduce the centrality of these cities, we would expect them to grow in the short to medium term.

The second tier of Russian cities, the dozen or so industrial cities with sizes around 1 million, are likely to face a mixed fate. A few of these cities (but only a few) are likely to grow, receiving population from the many cities below them. Others, however, seem likely to shrink. Because some of these cities are in climatically inhospitable places which make them ill suited to support their large populations, the Russian government might do well to recognize that these cities are likely to shrink and to adjust the structure of subsidies appropriately instead of continuing to expend resources in an effort to maintain them at unviable sizes.

Finally, the major impact of liberalization of mobility is likely to be a broad shrinking – in both size and number – of cities with populations between 100,000 and 1,000,000. This is the range that came to “bulge” under Soviet policies. Note, in figure 3.6, that the distribution of cities of size below 100,000 is relatively straight. If the transformation of Russia’s social and economic structure leads its migration patterns to become more like those of the other countries we have examined, we would expect the vast majority of cities in this middle size range to lose population to the Moscow and St. Petersburg as well as to the four to six industrial cities which grow.

Next Steps

More and Better Data on Urban Agglomerations

The results presented in this paper are theoretically plausible and provide excellent agreement with the data for the three countries studied. The availability of reliable data played a major part in the selection of these countries and we have attempted to reason from the dynamics of these data rich countries to make inferences about countries where solid data is less available. Our conclusions would be strengthened by a careful examination of data for additional countries. Because the careful definition of a city is critical to success in this work, such a study would require care and effort in order to generate data that are genuinely comparable with the cases presented here. While this is certainly possible for additional OECD countries, which have highly sophisticated statistical services and experienced local analysts, it presents a real challenge in the developing nations where these results are likely to be of most interest. Large developing countries (e.g. India and China) however, do have sophisticated statistical services and a consistent examination of their urban structures is likely to be possible.

We expect results from these countries to be broadly consistent with those observed here, though the fact that these countries are still undergoing massive urbanization is likely to produce urban distributions which are less mature and therefore less likely to accord with our model because their urban systems have not yet had time to approach an equilibrium distribution. Whether this expectation will be met remains to be seen.

Incorporate Florida's Work on the "Creative Class"

The concept of a floating population which is central to this model is akin to Richard Florida's "creative class" [Florida, 2002]. While we do not attempt to integrate Florida's observations about the migration dynamics of creative class workers, we suspect that these concepts are compatible and that Florida's observations could serve as a starting point for further refinement of the model – particularly as it applies to advanced economies like the United States and France.

Spatially Explicit Implementation

It is possible that a more realistic treatment of the urban interaction regime would stabilize the lower tail without the need to assume population growth. The current regime of purely random pairwise interaction has the virtue of simplicity, but is clearly unrealistic. A more realistic interaction regime could be based on the migration model suggested by Lowry [1966] which combines a measure of relative crowdedness (represented by wage levels and unemployment) with distance and size (in a gravity model formulation) to estimate the size and direction of migration.

While an elaboration of the model along these lines will likely be productive, it will introduce substantial complications. Most importantly, it would require the cities to be explicitly placed in space in order to allow for the calculation of distances between cities. In a model like this, the relative placement of cities would almost certainly matter to the model results, so it would probably be wise to begin with the actual arrangement of cities in existing countries. This elaborated model might also require that cities have somewhat realistic core values, which would reflect the insights of central place theory and its descendents [e.g. Fujita, Krugman, et al.,

1999]. We anticipate that a spatially explicit version of the model would be more realistic, but that this realism would come at a considerable expense in terms of simplicity and ability to generate intuitive insight. If however, this modification to the model produced stable power-law behavior without the need to introduce the assumption of population growth (or could help demonstrate that the assumption is warranted) it would be worth the effort.

Calibrate Parameters to Historical US Data

A spatially explicit version of the model with a more realistic interaction network could potentially be calibrated against actual migration and city size data to produce a model that had predictive power in real time. Variables to be calibrated would include: the importance of distance in determining interaction probability; the importance of wage and employment differentials in attracting migrants; the degree of imperfect information among potential migrants (i.e. the “bet” in the model when both cities are in the same position relative to their equilibrium size); the lag with which cities adjust to shocks to their equilibrium size; and the approximate distribution of core sizes among cities.

Performing this exercise would be interesting and, if the model performed well with realistic parameters, would go a long way toward validating the model. Establishing a baseline set of parameters in a country with abundant data like the United States would also provide some confidence when applying the model in countries with less detailed data. A well calibrated model would allow us to estimate the speed of transition in Russia and to gain insight into the magnitude of policy

intervention required to achieve specific goals relative to developing world megacities.

Conclusion

In this chapter we have presented a simple adaptive agent model of interurban migration which is capable of reproducing the city size distributions of the United States, France and Russia with only minimal and theoretically justifiable tuning. This model demonstrates the power of the adaptive agent modeling paradigm in a situation that is defined by simple rules, but is structured in such a way that these rules resist analytical treatment. The adaptive agent model's ability to incorporate bounded rationality (stochasticity) as well as heterogeneity among cities with respect to current size, evolving equilibrium size, and core size makes it a natural approach for this kind of modeling.

The use of this approach has made it possible for us to make real progress in understanding a phenomenon that has puzzled economists, geographers and others for over 50 years. Our model establishes a basis for moving beyond the assignment of mystical significance to the Zipf distribution of city sizes and allows us to see city size distributions as the result of straightforward behavioral rules. We can further understand Zipf as only a special case of city size distributions and see deviations from Zipf not as noise or error of some sort, but as the products of differing policies and situations.

Chapter 4: Spatial and Temporal Patterns in Civil Violence: Guatemala 1977 – 1986

Much of the existing literature examining quantitative aspects of civil violence concentrates on risk factors and searches for correlation between these factors and various indicators of violence. [Bates, 1983; Doyle and Sambanis, 2000; Fearon and Laitin, 1996] The foundation of these studies is generally annual, country level data on conflict deaths [Gurr and Harff, 1996]. While certain types of inferences can legitimately be drawn from such data, it does not lend itself to the study of internal conflict dynamics. This paper examines a substantially more detailed dataset covering the conflict in Guatemala during the ten year period 1977 to 1986. By shifting the basic unit of analysis from the country-year to the municipality-month, many intriguing patterns emerge. These patterns are generally indicative of "complex systems" behavior and point toward the use of adaptive agent modeling as a tool for exploring the dynamics of civil violence.

In our previous applications, we have used the adaptive agent modeling technique's ability to handle large numbers of interacting, heterogeneous parts and to incorporate the effects of bounded rationality. This has given us simplicity and flexibility in looking at international trade and analytical traction in looking at city sizes. An understanding of the phenomenon of civil violence is bound to make heavy use of these tools, but goes farther: the actors in any but the simplest riot think strategically and organize in ways that have a great deal of influence on the dynamics of violence. While the adaptive agent framework is better suited to this kind of work

than other modeling approaches (e.g. econometrics or systems dynamics) the use of the method does not make the problem of understanding conflict dynamics an easy one.

The primary aim of this chapter, therefore, is to demonstrate that the method is well suited to the task and to demonstrate some early work in the area. Though we will conduct some statistical analysis, we do not seek to present a comprehensive statistical, political, or historical portrait of the Guatemalan conflict – a task which has been ably undertaken by others [Ball, Kobrak and Spierer, 1999; CEH 1999]. We will also present some efforts at modeling civil conflict which show promise in beginning to understand the phenomenon. Our main objective, however, is to uncover patterns in the data which illuminate spatial and temporal dynamics in the conflict and which might be used to guide future efforts in the modeling of civil and state violence.

The Guatemalan Conflict

The history of state repression in Guatemala is, in many respects, particular to Guatemala and the victims of this repression were and are particular people with unique histories of their own. The unique history and personalities at the heart of this conflict mean that we can have little hope of generalizing many aspects of it. We can, however, observe certain patterns in detailed data derived from it which may be of use in understanding such conflicts in general. This understanding, in turn, may be of use in predicting, preventing and controlling conflicts in the future.

The Guatemalan conflict lasted from 1960 to 1996 with a period of greatly heightened violence in the early 1980's. The state carried out most of the killing

during the conflict in an ongoing campaign of repressive terror involving the military, the police, semi-autonomous “death squads” and state organized civilian “civil patrols” [Ball, Kobrak and Spierer, 1999]. The CEH estimates that over 93% of the killing was undertaken by agents of the state [CEH 1999].

Ethnicity played a significant role in the conflict. In the early parts of the conflict, the violence was typically between middle class people of the non-indigenous Ladino group struggling for control of the government. As the conflict progressed, it moved from an urban conflict focused on Guatemala City to a rural counter-insurgency campaign. The victims of state repression shifted, in about 1981, from middle class Ladino dissidents to indigenous Mayan peasants who were suspected of aiding rebel groups in the northwestern highlands. The scale and nature of the conflict changed as well, becoming vastly more deadly and including many acts which have been found to meet the formal definition of genocide. In the conflict as whole, about 83% of the victims were Mayans [CEH, 1999]. It should be noted that the dichotomous division of ethnicity into Ladino and Mayan is probably more clear to the Ladino controlled government than to members of the various Mayan groups, who speak many different languages and do not always consider themselves to be of the same ethnic group.

Data

This work is based on a remarkable data set constructed jointly by the American Association for the Advancement of Science (AAAS) and the International Center for Human Rights Research (CIIDH) under the direction of Dr. Patrick Ball of AAAS. It documents over 40,000 killings and disappearances in Guatemala between

1960 and 1996. Many of these records include the specific time and place where the incident occurred as well as other detailed information. It is based on an extensive review of Guatemalan press sources over the entire 36 year period and over 5,000 interviews with witnesses.

While there exist other data sets of this sort (for El Salvador, for instance) this is the only record of its kind which is published and generally available for research. It thus provides a fertile ground for the formation of hypotheses (because it is new) but can provide nothing in the way of confirmation of these hypotheses (because it is unique). It is hoped that research into the spatial and temporal dynamics of violence will spur interest in this kind of disaggregated data and lead to the creation and publication of additional data sets.

This research uses a subset of this data spanning the ten year period of 1977 to 1986. Data are further restricted to killings and disappearances for which the date was known to at least the nearest month. This subset contains 24,000 cases which probably constitutes about 10% of the killings during this 10 year period. This estimate is uncertain because the number of killings overall has been estimated at anywhere from 80,000 to 400,000.

The analysis that follows assumes that this sample is relatively unbiased. This is, of course, a risky assumption in spite of the rigor with which the data were collected. Davenport and Ball provide an excellent discussion of the biases inherent in various types of human rights data collection and of how this data set avoids many of these problems [Davenport and Ball, 2002].

Methods

Much of the existing quantitative treatment of large scale violence relies on summary statistics which provide information about a conflict over a large span of space (a nation or a conflict zone) and of time (a year or the duration of a conflict). Many of these studies use linear regression and related statistical techniques to correlate violence with other factors in an effort to understand and predict such outbreaks.

Given the richness of this data set, this paper takes a different approach. It tries to preserve the complexity of the data wherever possible and to explore the finer grained data for regularities which might be applicable in other situations. Major tools in this effort included complex queries of the data using Structured Query Language (SQL), spatial analysis and mapping with a geographic information system (GIS), histograms, time series plots, rank/size plots and other, mostly graphical, representations of disaggregated data.

This approach has limitations. In most statistical analyses, one tries to form hypotheses independent of the data and then use the data to test these hypotheses. In this case, an examination of the data was used to construct hypotheses, making it impossible to use the same data to test these hypotheses. The observations which follow are therefore offered not as proven generalizations, but as suggestive patterns with theoretical plausibility. The confirmation of their generality will have to wait for detailed data from other conflicts.

Observations

Frequency vs. Severity

In the data set, the frequency of killing in a municipality is only weakly correlated with the overall quantity of killing in that municipality. The coefficient of correlation between frequency of killing (number of months where at least one person was killed in a town) and quantity of killing (number of people killed in the town over the whole study period) is .65. More tellingly, perhaps, the correlation between the number of people killed individually and the number of people killed in groups larger than one is only .31.

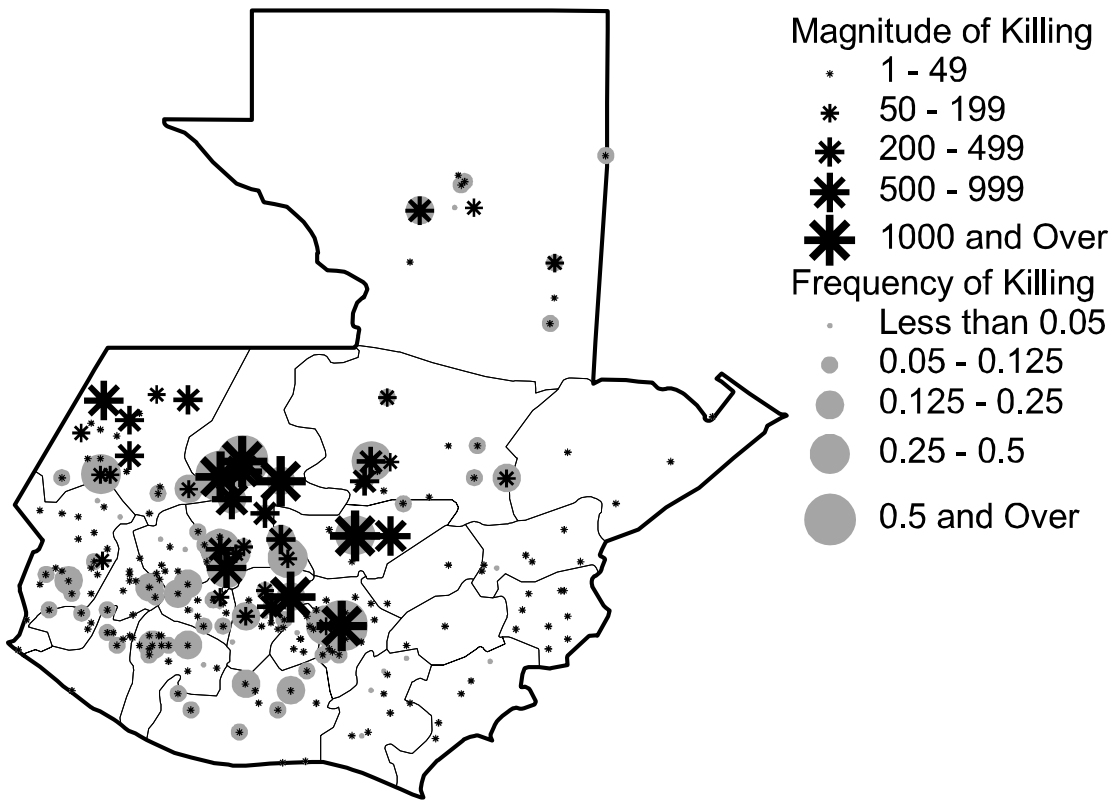


Figure 4.1: Map of frequency and severity of killing, by municipality

Sources indicate that many of those killed individually were personally targeted by the government [CEH, 1999]. This observation supports the modeling observations [Epstein, Steinbruner and Parker, 2000] that the removal of leaders is an effective repression technique. We might assume that the government is well aware of this phenomenon and removes leaders (by killing them) in areas where it knows who these leaders are. These are the areas where we see a large number (and a high frequency) of single assassinations. In areas where the government does not know who the leaders are, we see more people being indiscriminately killed. This may be a result of a combination of two factors. On the one hand, the government may have killed indiscriminately because it did not know how to choose its targets. On the other hand, insurgent activity may have been able to gain a greater base because the government was less able to repress it through assassination. Government perception of this greater insurgent base may have provoked it to more indiscriminant killing.

The hypothesis that less knowledge on the part of the government can lead to more indiscriminant killing is further supported by the tentative observation that violence was more intense in inaccessible areas. While this is hard to quantify precisely, it appears that massacres were more likely to be carried out in the mountains and away from improved roads.

Ethnic Mix

A second observation resulting from the spatial disaggregation of the data is that amount of killing in a municipality has a somewhat complex relationship to the ethnic mix in that municipality. While the population of Guatemala is fairly evenly divided between the Ladino and Mayan ethnic groups, they are generally segregated

at the municipal level. About 76% of the population lives in municipalities which are more than 80% dominated by one group or the other.

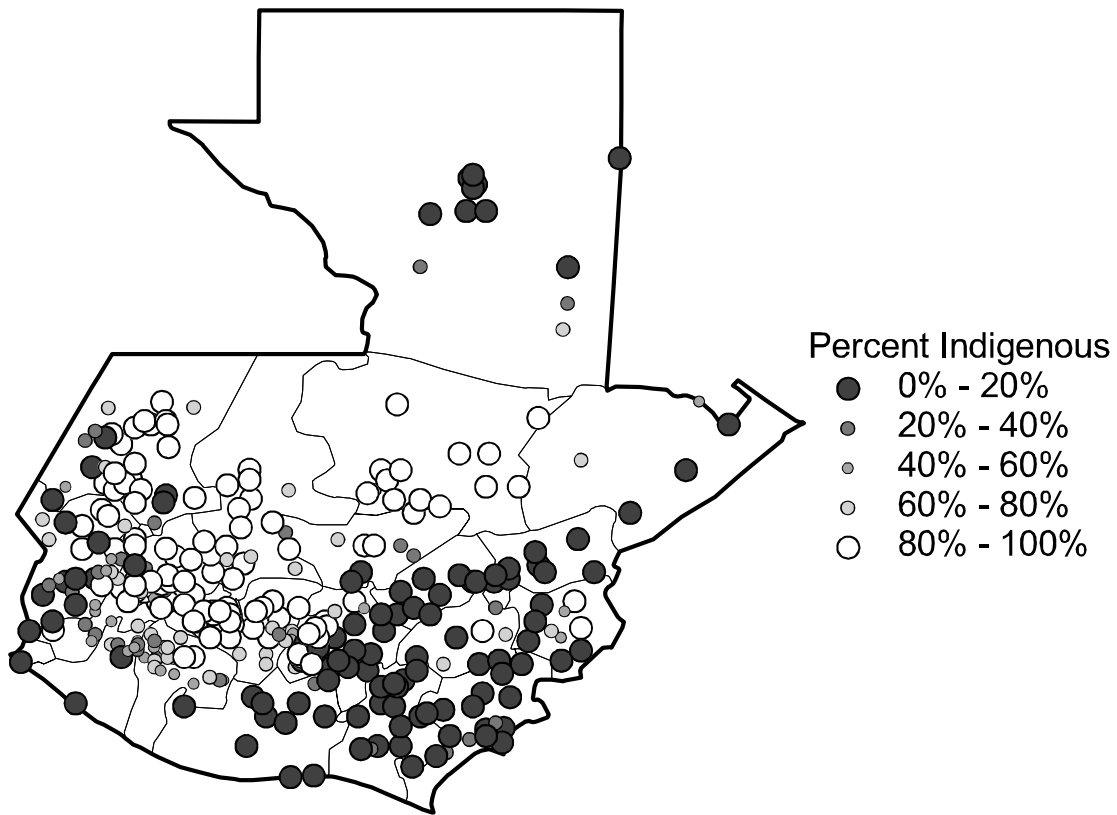


Figure 4.2: Map of Ethnic Distribution in Municipalities

The Mayans live largely in the mountainous northwest section while Ladinos occupy the lower and more agriculturally productive south and east portions of the country. Even within these regions, however, there is significant polarization.

Examination of the quantity of killings within these largely segregated municipalities led to an unexpected finding: the few municipalities where Mayans make up a large, but not overwhelming, majority were the most consistently dangerous. Just over half of the killing took place in municipalities in which the Mayans made up between 80 and 90 percent of the population. This is remarkable because such municipalities make up less than 8% of the municipalities in the country

and house just over 8% of the total population (about 17% of the Mayan population). Many more Mayans (45%) live in municipalities where they constitute upward of 90% of the population. Though these municipalities also saw considerable bloodshed, they did not have as much violence as those that were 80 to 90 percent Mayan. Because the number of municipalities is relatively large (n=345) these variations are unlikely to be a pure statistical artifacts (the differences are significant beyond the 99% level). While we might expect violence to increase monotonically with the percentage of Mayan residents (since it was primarily Mayans who were killed), this proves not to be the case.

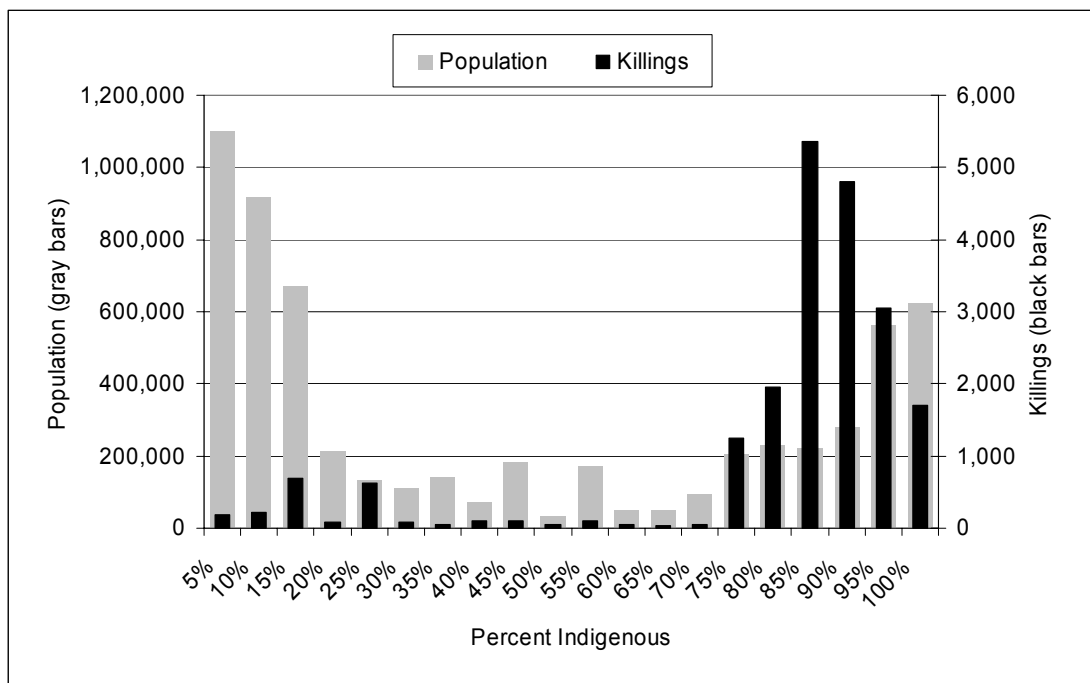


Figure 4.3: Histogram of Ethnicity and Killing.

At least two different mechanisms might explain the fact that more killing took place in municipalities with significant Ladino minorities than in municipalities with almost entirely Mayan populations. One thought is that the rate of killing increases with the percentage of Mayan residents up to a point because as this

percentage increases, the government knows less about the leadership structure of the insurgency, and is thereby inclined to kill indiscriminately as discussed above.

Beyond some point, however, the government may know too little to do anything.

This would be a real world example of a little knowledge being a dangerous thing. A lot of knowledge leads to assassination of leaders, a little knowledge leads to indiscriminant killing, and no knowledge leads to no action.

A second mechanism might be based on group dynamics. There may be a threshold concentration that individuals with a minority trait must reach before they are considered (or consider themselves) a group. It is possible that, in municipalities where the Ladino population constituted less than 10% of the population, tensions between the groups were substantially less because at some basic level, the Ladino population did not constitute a separate ethnic group.

An examination of the opposite end of the histogram provides some support for this interpretation. We see a similar, though much smaller, bump in the number of killings in the range between 10% and 25% indigenous (i.e. 75% to 90% Ladino). The vast majority of the killing in the conflict was directed against Mayans, and these areas had relatively few Mayans. Therefore, it is not surprising that fewer people were killed in these areas. The basic insight remains the same however. In areas where Mayans constituted less than 10% of the population, they may have been perceived more as individuals than as a threatening group.

Thus, at both ends of this histogram where one group or the other is more than 90% dominant, we see less violence. This may be because, in such communities, people relate as individuals rather than ethnic groups. Such communities might be

more tightly knit and better able to avoid government persecution. Also in the middle of the histogram, where neither group is more than 75% dominant, we see relative safety. The area between 75% and 90% dominance, however, seems to be much more volatile. If this observation is born out in the examination of local populations in other conflicts, it could prove to be a useful rule of thumb for peacekeeping operations.

Punctuated Equilibrium

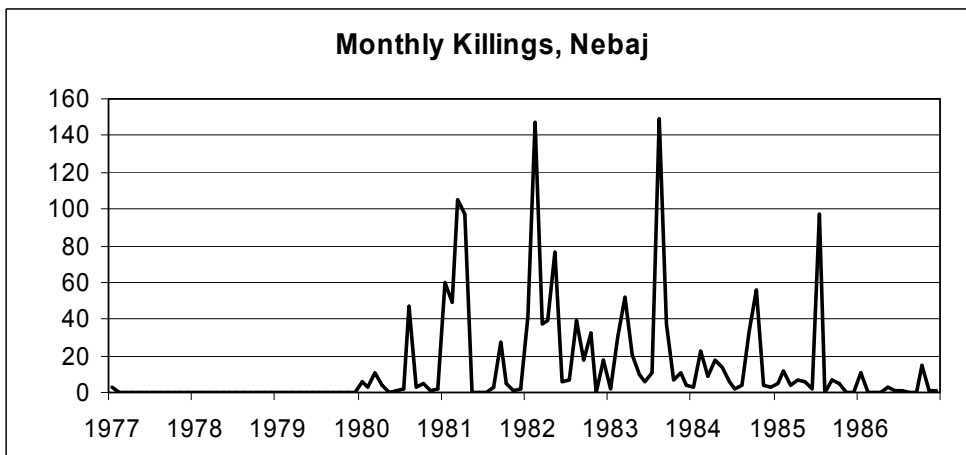
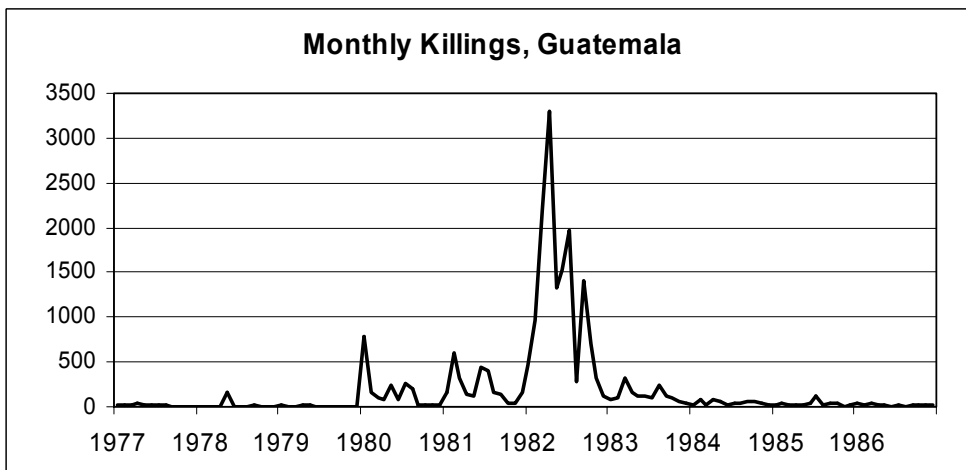
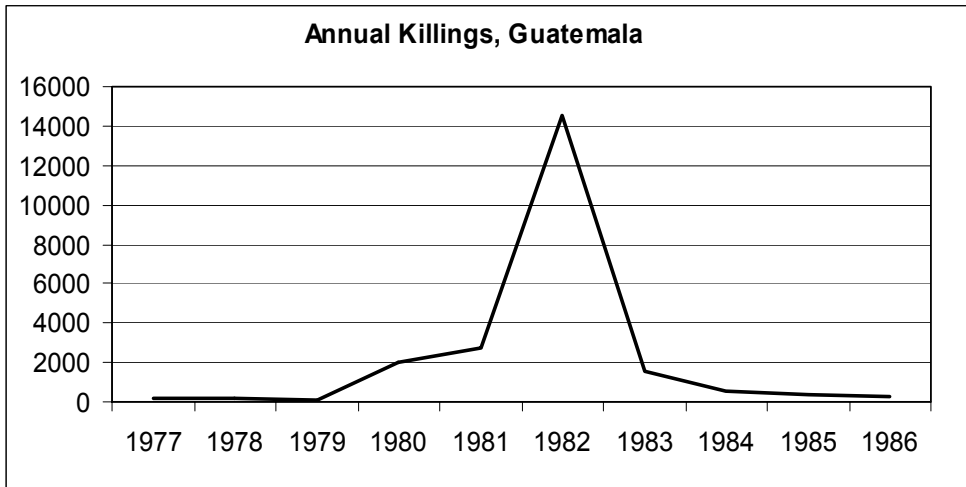


Figure 4.4: Time Series Graphs: Annual, Monthly, Monthly for a Single Town

Another striking observation arises when the data are disaggregated with respect to time as well as space. The violence in a given place does not expand and contract smoothly over time. Rather, the pattern of violence is “spiky”. A municipality may go for some time without an incident and then experience a major incident, or cluster of incidents. This becomes increasingly apparent as we move from aggregate annual numbers to finer resolutions of time and space.

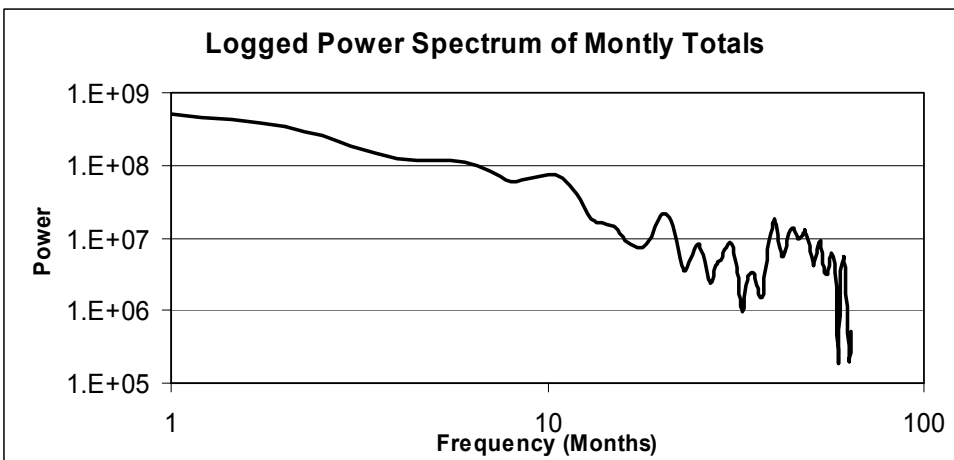
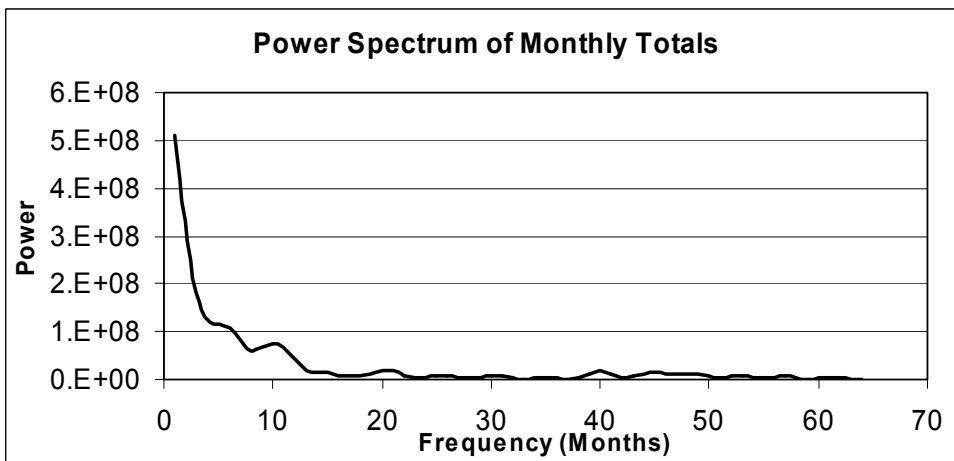
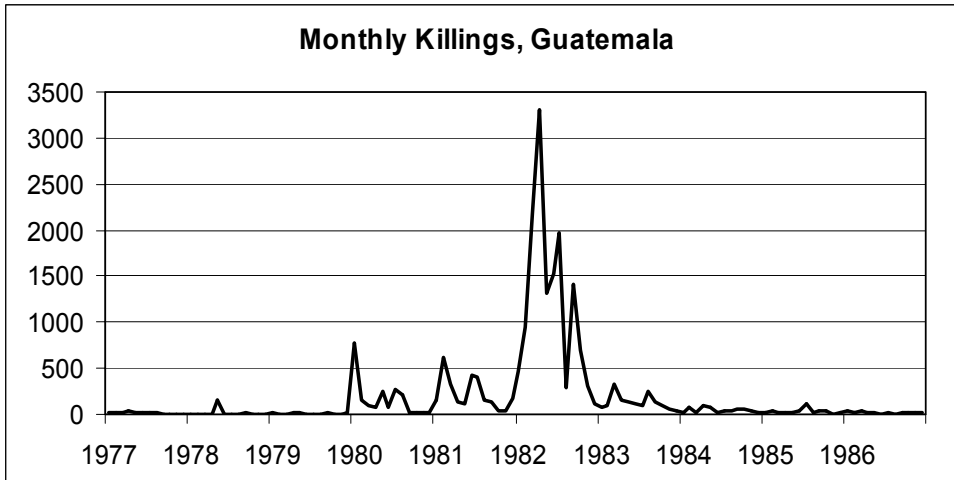


Figure 4.5: Killings by month, with power-spectrum and logged power-spectrum plots

An objective measure of the character of such time series data can be obtained by examining its power spectrum. Power spectrum analysis is a method borrowed

from Physics where it is best known as a way of characterizing sound waves. Mathematical techniques exist, most notably the Fourier transform, to decompose a given signal into the spectrum of sine waves of different amplitudes – different “powers” – which can be combined to reproduce that signal. This technique has been extended to look at a broad range of time series phenomena, ranging from earthquakes and floods to stock prices [Schroeder, 1991].

Purely random noise, also known as white noise, has equal power at all frequencies. Complex systems, however, frequently exhibit "pink" noise, also called "1/f noise", where the power at a given frequency is inversely proportional to the frequency [Schroeder, 1991]. An examination of the time series of monthly killings in the Guatemala data set (using a Fourier transform) reveals this kind of power law spectrum. The presence of a power law in the power spectrum of this time series suggests that the analysis of civil violence might benefit from the application of techniques used to examine other, better understood complex systems.

In the case of this conflict, the exponent of the power law is not precisely -1 (i.e. $1/f = f^{-1}$), but something closer to -1.4. This exponent provides a kind of signature for a process exhibiting pink noise [Bak, 1997]. Examination of other conflicts may reveal that this signature is consistent from one to the next or that it varies in a way that is informative.

Distribution of Incident Sizes

An examination of the distribution of incident sizes within the data set provides some additional insight into the internal dynamics of the conflict. The conflict can be separated into two parts: a "normal" (i.e. non-genocidal)

counterinsurgency and a genocide which was focused in the western highlands in 1981 and 1982. The counterinsurgency is characterized by a "Zipf" distribution of incidents, whereas the genocide follows a different pattern.

A sense of the overall distribution of incident sizes is given by the rank/size (or Pareto) plot presented in Figure 4.6.

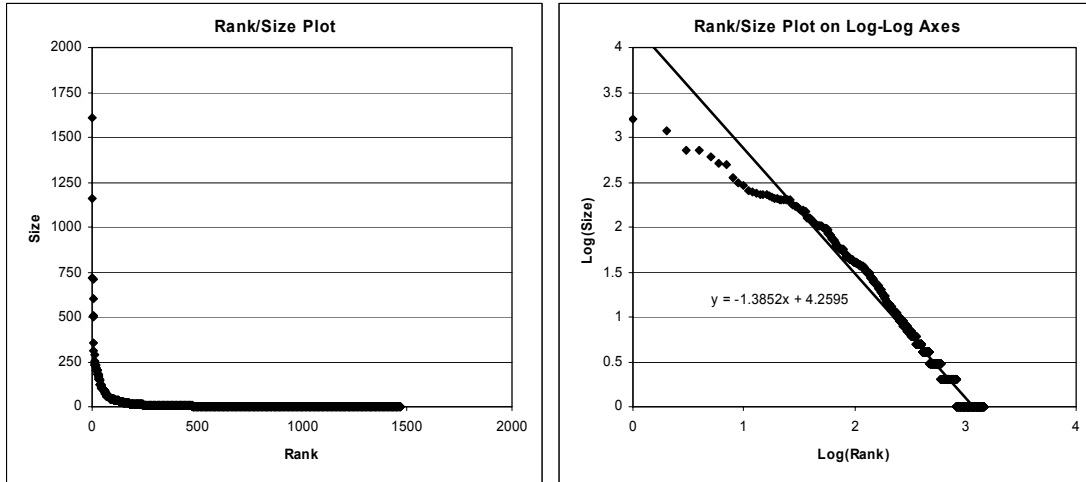


Figure 4.6: Rank/Size Plot of Killings per Municipality-Month.

This ordered histogram (on log-log axis) gives only a rough idea of the real distribution for several reasons. First, it combines regular conflict and genocide -- two processes which, I will argue, follow different dynamics. Second, it does not represent killings per incident directly, but rather killings per municipality per month. This is due to data limitations. Both of these problems can be worked around.

To examine the difference between the regular and genocide parts of the conflict, we need to partition the data with respect to both time and space. We saw above in figure 4 that 1981 and 1982 were years of particularly intensive violence. Figure 4.7 examines this period with respect to spatial distribution by showing the

number of massacres (defined as incidents where five or more people were killed during the same incident) per town in 1981 and 1982.

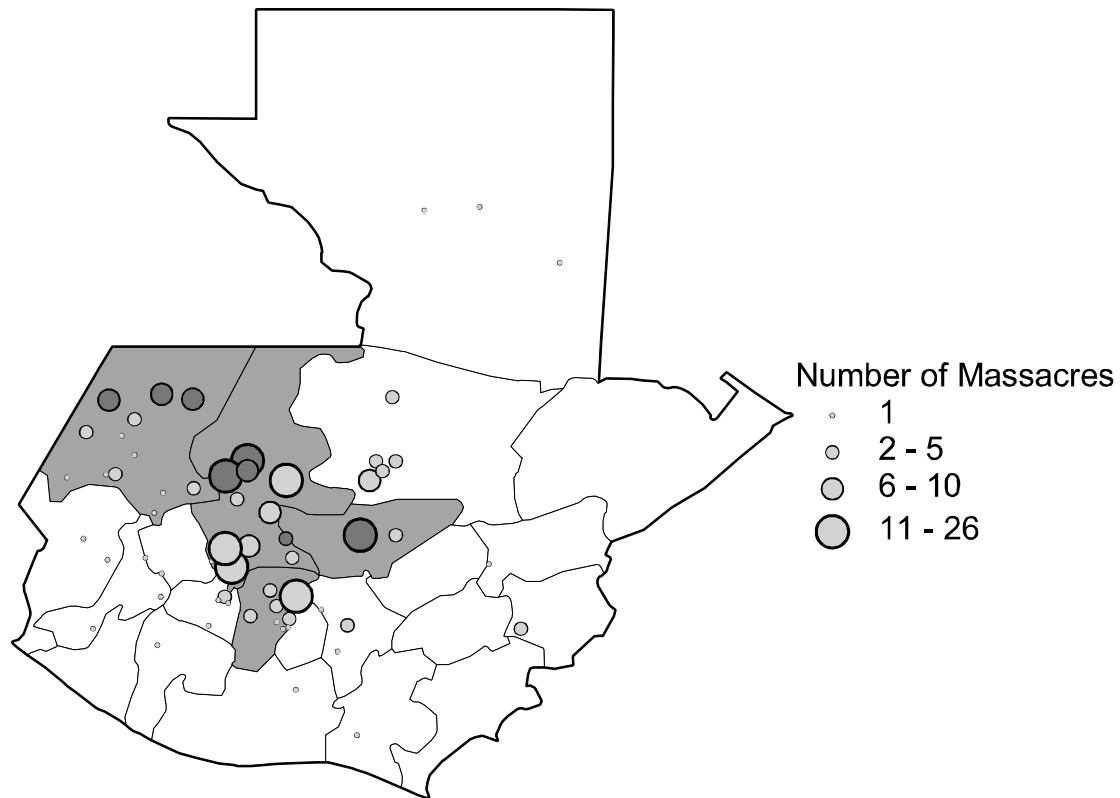


Figure 4.7: Confirmed genocidal massacres (dark circles) and massacres for which genocide status was not determined by CEH (light circles), 1981-2

In its 1999 report, The Guatemalan Commission on Historical Clarification (CEH) documented, with painstaking thoroughness, a number of incidents during which the formal criteria of genocide were met [CEH, 1999].

These criteria are laid out in Article II of the United Nations Convention on the Prevention and Punishment of the Crime of Genocide (1948). The CEH applied the convention using this reasoning:

Considering the series of criminal acts and human rights violations which occurred in the regions and periods indicated and which were analysed for the purpose of determining whether they constituted the crime of genocide, the CEH concludes that the reiteration of destructive acts, directed systematically against groups of the Mayan

population, within which can be mentioned the elimination of leaders and criminal acts against minors who could not possibly have been military targets, demonstrates that the only common denominator for all the victims was the fact that they belonged to a specific ethnic group and makes it evident that these acts were committed “with intent to destroy, in whole or in part” these groups (Article II, first paragraph of the Convention). [CEH 1999]

All of these incidents involved massacres of Mayans in the highlands between 1981 and 1982. The CEH further acknowledges that many additional incidents of genocide took place but were not formally documented. In Figure 7, the municipalities in which the CEH documented genocide are colored black.

By taking the number of massacres in a municipality as a proxy for the level of genocide activity in that municipality, we can roughly identify four departments (Huehuetenango, El Quiche, Baja Verapaz, and Chimaltenango) as the focus of the genocide. In order to look for differences between genocide and regular warfare, we separate records from these four highland departments during 1981 and 1982 (the genocide subset), from the rest of the data set (the non-genocide subset).

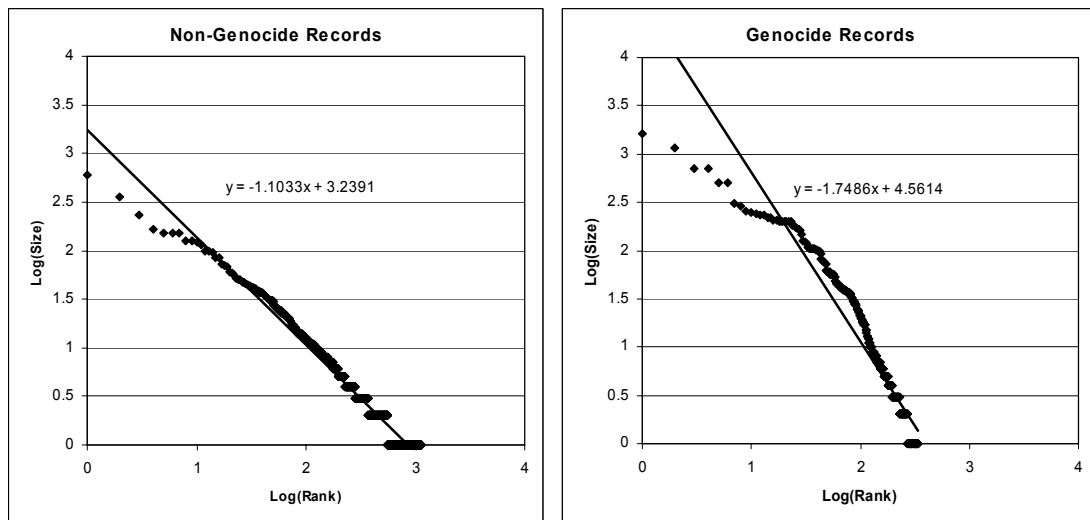


Figure 4.8: Rank-size plots of the nongenocidal-killings subset and the genocidal-killings subset, by municipality-month

The municipality-month rank/size plots for these two subsets look significantly different. The non-genocide subset (n=1133) closely approximates a straight line with slope -1.13 in log-log coordinates. This is to say that the distribution can be described by a power law of the form $S=\alpha R^{-1.13}$ (where S is size and R is rank). The genocide set (n=338), on the other hand, is quite concave toward the origin and has a much higher slope (to the extent that it can be described by a power law at all).

The non-genocide subset is actually even closer to the power law distribution than it might appear. The departure in the upper tail is due to two or three "extra" events with size around 250. It is these few events which leave the distribution short at the top end. This is quite different from the genocide dataset, where the largest 30 or so events describe a curve with slope much lower than the distribution would require.

Once purged of the genocide related records, the regular conflict data adhere more closely to a power law distribution, but still reflect a slope based on the somewhat artificial unit of the municipality-month. While the resolution of the data is not sufficient to examine the exact size distribution at the incident level, it does allow us to estimate the total number of incidents represented by the data – about 3500 in the non-genocide set.

If we think of the municipality-month as an aggregation bin, then the non-genocide, municipality-month set (n=1133) represents an average of 3.05 incidents per bin. By further aggregating the data temporally at the 6 month, 1 year, 2 year, 5 year and 10 year levels, and determining the power law exponent at each of these

levels of aggregation, we are able to establish a linear relationship between incidents per bin and the exponent.

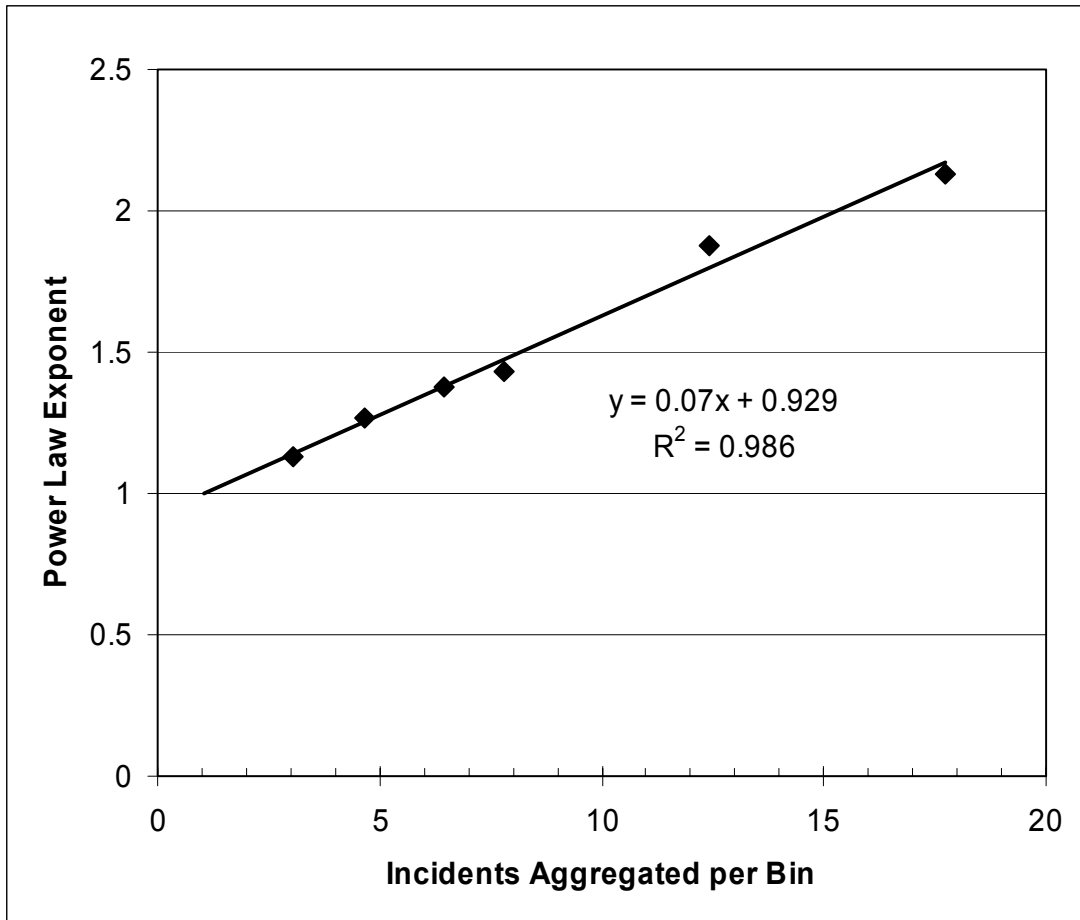


Figure 4.9: Trend of the power-law exponent for different levels of aggregation in the nongenocidal subset

This relationship is nicely described ($R^2=0.986$) by the linear relationship $y = -0.7x - 0.929$, where y is the exponent and x is the average number of incidents per bin. From this empirically derived relationship, we can estimate the exponent for the fully disaggregated case where there is only one incident per bin by simply evaluating the expression at $x=1$. The resulting value of -1.056 is extremely close to -1 , the exponent which defines the so called "Zipf" distribution. The Zipf distribution is characteristic of many processes in the physical and social worlds including city

and firm sizes, earthquake magnitudes, certain aspects of Internet traffic, and a host of other phenomena [Bak, 1997]. Similar results have been demonstrated for the distribution of conflicts (rather than incidents within a conflict) [Richardson, 1960; Cederman, 2002; etc.].

That incidents from the regular part of the conflict follow the Zipf distribution is interesting and invites speculation as to why this might be the case. While this is same distribution that we saw in the previous chapter on the distribution of city sizes, the mechanism that we explored in that chapter does not seem applicable here. In both cases, the distribution emerges from complex system, but it would seem to emerge from quite different rules.

A random growth rate model is a simple way to create a Zipf distribution and the workings of such a model are suggestive here. The model involves an arbitrary number of objects (in this case, potential incidents), each of which has a size greater than or equal to one ($S \geq 1$). The initial distribution of sizes is not important to the long term behavior of the model, so we will start them all at one. In each model iteration, each object grows or shrinks by a random amount ($S_t = g * S_{t-1}$ Where g is a random variable: $-0.1 < g < 0.1$). A final condition of the model is that no object can become smaller than one (If $g * S_{t-1} < 1$ Then $S_t = 1$). The result is a collection of exponential random walks bounded by a floor of one.

If, from this set of exponential random walks, a sample is drawn at any arbitrary time (of size N), the sample will be Zipf distributed [Gibrat, 1931; Gabaix 1999]. The largest object in any given sample can be expected to have a size approximately equal to the size of the sample ($S \approx N$) and the distribution is described

by $S_n = N \cdot n^{-1}$. Thus, for $N = 1000$ the size of the largest object (S_1) could be $1000 \cdot 1^{-1} = 1000$. The size of the next largest object (S_2) would be $1000 \cdot 2^{-1} = 500$. The size of the smallest (S_{1000}) would be $1000 \cdot 1000^{-1} = 1$.

The same distribution arises independent of the initial distribution of sizes and also independent of the range of growth rates. So long as the growth rate is drawn from a range equally distributed around zero, the distribution will converge toward $S_n = N \cdot n^{-1}$, with larger ranges converging faster.

We can make an analogy here to incidents of violence during a conflict. During a "normal" conflict (say a counterinsurgency like Guatemala's), the objective of the repressive force is not directly to kill people. The objective is to put down the rebellion and secure the power of the state. Killing is a means to this end. The state and its agents therefore operate according to heuristic rules under which the level of killing can vary tremendously depending on the situation. Repression according to heuristic rules can be conceived of as similar to the random growth rate model. Since there is no guide to how much killing is the right amount, each incident unfolds according to the goals and perceptions of the two sides.

This is not to say that there was not central control behind the Guatemalan state forces. There undoubtedly was. However, during the "normal" parts of the conflict, the central orders may have taken the form of rules: "Suppress the insurgents", etc. The objective was specified, but the amount of killing required to meet the objective was probably not specified and was thus dependent on the dynamics of the given situation.

The fact that incidents in the non-genocide subset appear to be Zipf distributed is remarkable because it relates the number of incidents to the sizes of incidents in a more direct way than one might think possible. Given the sizes of the largest few events, one can estimate the number of events and the total number killed. Given the number of events, one can estimate the size of the largest events and the total number killed. Given the total number killed, one can estimate the size of the largest events and the number of events. These estimates would be expected to be rough but a rule which would allow even order of magnitude guesswork would be unexpected and might have considerable prognostic power.

Because the distribution is drawn from a single conflict, there is reason to ask whether the Zipf distribution of incident sizes is a common one in conflicts. This is an open question which awaits empirical verification against other data sets. The fact that Zipf distributed events often occur in complex systems phenomena, along with the random growth model analogy gives us some reason to believe that it might be typical.

Examining the genocide subsample, on the other hand, reveals a different pattern. As discussed above, the distribution of incidents in the western highlands in 1981 and 1982 (where the CEH identified acts of genocide) was quite different from that resulting from the rest of the conflict. There are far more "middle sized" events where between 10 and 100 people are killed. This is consistent with a different kind of command, a much more direct order to go to a place and kill people. Where the basic logic of normal conflict is to accomplish the objective while taking as little risk as possible (which means avoiding incidents if possible), the basic logic of genocide

is to kill some fraction (perhaps 100%) of a given population. In normal conflict, killing is a tactic whereas in genocide, it becomes a strategy. It is because of this basic difference in the function of violence that incidents under normal conflict lack a characteristic size and follow the Zipf distribution, while incidents of genocide tend to have a characteristic size that relates to other factors like the size of a military unit or the size of a village.

If this hypothesis proves consistent with data from other conflicts it would provide several potent tools. First, if the conflict was known to be of the “normal” sort, it might be possible to assume that incidents would be Zipf distributed – providing statistical leverage which has previously been unavailable. Second, the distribution of incidents could provide evidence of the nature of the orders and command structure in a conflict, providing a statistical means of differentiating normal and genocidal warfare.

Modeling

Efforts to use adaptive agent modeling to understand civil conflict remain embryonic, though they have recently begun to elucidate parts of the problem. Cederman [2003] updates Richardson’s [1960] work on the magnitude of wars, finding them to be power-law distributed and goes on to present an agent based approach to understanding this distribution. Bhavnani and Backer [2000] use an adaptive agent approach to explore the role of group coherence and information flow in determining the duration and intensity of violent interethnic conflicts including those with genocidal components. Srbljinovic et al. [2003] have applied agent based techniques to understanding the process of ethnic mobilization in the former

Yugoslavia. Taylor et al. [2004] have made progress in developing an agent based platform designed to help intelligence analysts understand complex geopolitical situations. Kewley and Larimer [2003] have taken an agent based approach to quantifying the value of tactical information in a battle situation – finding the technique useful in both analytical and decision support modes.

The Brookings Model

We will explore the utility of this approach in context of Epstein, Steinbruner and Parker's [2001] model of civil violence which was developed at the Brookings Institution. The Brookings model presents a simple, highly generalized framework of civil violence. While its authors do not make any claim of completeness in the model, it does reproduce a number of features observed in civil conflicts and provides a conceptually elegant way of approaching the problem.

In its most basic formulation, the Brookings model is implemented with two types of agents: citizens and cops, which move about and interact on a lattice. The citizens have four state variables:

- Hardship (H) – This measure of perceived hardship is set exogenously and is distributed heterogeneously among agents (i.e. different citizens suffer different levels of hardship).
- Legitimacy (L) – This measure of perceived legitimacy of the central authority is also set exogenously and is equal across citizens. It can be varied over the course of a run.
- Risk aversion (R) – This measure varies across citizens and represents the variation among individuals in their tendency to act on their grievances. Some people can become very angry without acting out, whereas other, more hot-headed, types will express their displeasure under almost any circumstances.

- Vision (v) – The agents do not have global information about the system, but instead act on what is happening around them within their range of vision. Vision is set exogenously and is the same across agents.

These state variables are used to calculate three important quantities for each agent in each round:

- Greivance (G) is calculated as $H(1-L)$. This is to say that the grievance that a citizen feels toward the central authority is a product of the hardship that she experiences and her measure of the “illegitimacy” of the regime. Epstein et al. point out that the high legitimacy that British government enjoyed during World War II ensured that the extreme hardship imposed by the blitz of London did not generate grievance toward the government.
- Arrest Probability (P) is calculated as $1-\exp[-k(C/A)v]$ where k is a constant, and $(C/A)v$ represents the ratio of cops to actively rebellious citizens within the given citizens vision.
- Net Risk (N) is calculated as RP . The agents perceived (or net) risk is the product of her level of risk aversion (R) and her calculated probability of arrest (P).

These seven quantities are sufficient to formulate a rule for acting rebelliously:

- If $G-N > T$, be Active; Otherwise, be Quiet – In other words, if the agents grievance exceeds her net risk of arrest by some threshold value, she will act out against the government, otherwise, she will not.

Finally, citizens have a movement rule: At the start of each round, they move to a random, unoccupied space within their vision.

Cops are much simpler. They have only one state variable: Vision. They have one basic behavioral rule: in each round, they arrest a random active citizen within their range of vision. They also follow the same movement rule as citizens do, moving to a random unoccupied space within their vision at the beginning of each round.

Epstein et al. demonstrate that these simple rules are sufficient to generate a remarkable variety of phenomena which have been observed in civil conflict. They demonstrate that random free assembly can lead to rebellious outbursts. They reproduce the observation that sudden shocks to a regime's legitimacy are much more destabilizing than slower erosion of legitimacy – even if the sudden shock is smaller in magnitude than the slower loss. They contrast the stability of a slow reduction in legitimacy with the explosive potential associated with a slow reduction in repression (simulated by slowly reducing the number of cops), thus illustrating DeTocqueville's comment that "liberalization is the most difficult of political arts."

Epstein et al. then introduce a model of intergroup violence where two types of citizens (dubbed "red" and "blue" are introduced. Legitimacy is redefined in terms of each groups willingness to recognize the other groups right to exist and activation redefined to involve the killing of an member of the outgroup. Cops retain the same behavior as in the single group model. Because they arrest active agents without regard to their group identity, the cops now take on the function of peacekeepers.

This two group model, though highly stylized, is also capable of reproducing features of real conflict. Epstein et al. observe that, when intergroup legitimacy is reduced, peaceful coexistence gives way to localized ethnic cleansing and then to genocide. When peacekeepers are present from the outset, the model can produce a basically stable society with endemic ethnic violence. When they are introduced after violence is underway, they can produce safe havens which allow both groups to exist. At lower levels of intergroup legitimacy, however, they observe that peacekeeping is

a dicey venture – often failing to prevent genocide even with a large number of peacekeepers.

Evaluating the Brookings Model

Our goal in examining the data from Guatemala was to establish benchmarks which an agent model might seek to reproduce. Though the Brookings model does not reproduce our findings in its current form, it reproduces more of the dynamics of civil violence than one might guess based on their simple structure and rules. Given the high degree of abstraction in the current model, it is not surprising that we observe are not precisely reproduced. However, a comparison of the capability of the model relative to our findings demonstrates that the approach is a promising one for gaining better understanding of the dynamics of civil violence.

Our observation that the frequency and the severity of violence were only weakly correlated (i.e. that the places with frequent violence were not the same as those with extreme violence), is consistent with the Brookings group's observation that the elimination of leaders is an effective repression technique. Perhaps more interestingly, it is consistent with Bhavnani and Backer's [2000] model result (backed by data from Burundi and Rwanda) that the conditions that lead to interethnic trust relationships which lead to endemic violence are different from those which lead to extreme violence of shorter duration.

While the Brookings model does not reproduce our finding of the complex relationship between ethnic mix and level of violence, it could be extended to explore this relationship. The gist of our finding was that areas which were between 75% and 90% populated by one ethnic group or the other were more violent than areas that

were either more or less balanced. The Brookings model, as currently specified, treats both ethnic identity and interethnic legitimacy exogenously.

Various studies have demonstrated that adaptive agent methods are useful in understanding the emergence of group identity and the examination of intergroup dynamics [e.g. Epstein and Axtell, 1996, Axelrod, 1997]. The agent approach would make it relatively easy to introduce endogenous dynamics in both the strength of ethnic identity (along the lines of Srbljinovic et al. [2003]), and of intergroup legitimacy (along the lines of Fearon and Latin [1996] as well as Bhavnani and Backer [2000] suggest ways in which these issues could be incorporated into the model without undue complication. Epstein et al. outline steps in this direction in an appendix. If this relationship between ethnic mix and violence can be produced with a plausible theoretical model and proves to be consistent with data from a small number of additional conflicts, it could serve as a rule of thumb for planning peacekeeping operations and prioritizing the deployment of peacekeepers.

The Brookings model also shows promise in being able to reproduce the punctuated equilibrium nature of the Guatemala data. We observed that the incidence of violence was far from smooth, particularly when the data were disaggregated with respect to space and time (this is illustrated in Figure 4, above). Such punctuated equilibria are common in complex systems [Bak, 1997], and are produced by the Brookings model. While the time series data for the Guatemalan conflict are not sufficient to compare with the waiting time analysis conducted by the Brookings group, the “spiky” texture of the violence is qualitatively similar. This is a sharp contrast to the dynamics of econometric or systems dynamics models, which tend to

produce predictions that rise and fall smoothly over time. In this respect, the agent approach would appear to be the only viable way to quantitatively explore the micro structure of civil violence.

The Brookings model can be measured against the Guatemala data in terms of the distribution of incident sizes it produces. The model as specified does not produce Zipf distributed events, however, it does produce a heavy-tailed event distribution. The fact that the distribution does not match the one we observe in Guatemala is not surprising given the highly abstract nature of the model. Epstein et al. point out, “the point to emphasize here is not which distribution is best, but that some macroscopic regularity emerges. A major strength of agent models is that they generate a wealth of data amenable to statistical treatment.”

If the Zipf distribution of incident sizes proves to be a common feature of non-genocidal conflict, then we would expect it to emerge from a relatively stylized model of violence. While the Brookings model achieves a number of striking results with a minimal set of rules, the fact that it does not produce this feature may indicate that significant aspects of the dynamics of violence are missing from it.

One way of thinking about the Zipf distribution in relation to a normal distribution is that the large incidents are very large and the small incidents are very numerous. A Zipf distribution can, in general terms, be produced by a phenomenon which balances positive feedback (making the large events larger) and negative feedback (keeping most events small). The Brookings model includes a mechanism for positive feedback in that the more active citizens a given citizen can see, the more likely that citizen is to become active. It also includes a mechanism for negative

feedback in that arrests reduce the number of relatively aggrieved and/or risk neutral agents in a local area until the violent outburst can no longer be sustained.

While these characteristics of the model are sufficient to produce punctuated equilibrium, they do not operate in a way that produces the observed distribution of events. The density based negative feedback mechanism may be a factor here. In the Brookings model, an incident generally dies out because a critical mass of active agents is arrested, thus reducing their local density. This mechanism does not kick in until the incident is under way, with the number of arrests required to bring it to an end being dependent on the vision and density of the agents. A more realistic event distribution might be produced by supplementing this mechanism by endogenizing risk aversion by letting arrests have a deterrent effect on those who witness them.

The point here is not to propose a revised version of the Brookings model, but simply to point out that its adaptive agent structure makes it a flexible tool for exploring the dynamics of civil violence. Some of these explorations (e.g. endogenizing intergroup legitimacy or risk aversion) would involve simple modifications of the model, whereas others (e.g. introducing hierarchical command structures or strategic behavior) would be more challenging. The general approach of adaptive agent modeling, however, seems ideally suited to exploring the complex dynamics of civil violence.

Conclusions

Examination of detailed data from the Guatemalan conflict between 1977 and 1986 reveals a number of novel patterns which support the use of complex systems methods, including adaptive agent modeling, for understanding the dynamics of civil

violence. The lack of strong correlation between individual and larger scale killings within municipalities provides some support for the notion that the removal of leaders is an effective repression technique. A comparison between the amount of killing in municipalities and the ethnic mix in those municipalities reveals a non-linear relationship between ethnic mix and killing; this invites analysis based on group dynamics. The temporal texture of the conflict is far from smooth, with a power spectrum that closely resembles that of other, better understood, complex systems. The distribution of incident sizes within the data seems to fall into two distinct sets, one of which (corresponding to "regular" conflict) is Zipf distributed and lacks a characteristic size, the other of which includes acts of genocide and is distributed quite differently -- possibly reflecting the different role that killing plays in these different types of conflict.

Because of the unique nature of the Guatemala data set, all of the findings in this paper need to be considered as preliminary empirical results. It is hoped, however, that the findings are sufficiently provocative to encourage the compilation and release more data sets of this sort. Many aspects of civil violence seem to depend on the internal dynamics of a conflict, and will not be revealed without a careful examination of detailed data from many conflicts.

Finally, we saw that adaptive agent modeling is a technique which is well suited to exploring the regularities presented here. A brief survey of the literature in the area of agent based modeling of conflict indicated several promising lines of research relating to intergroup dynamics and the general structure of violence. A more detailed examination of the Brookings civil violence model showed that while

the model does not produce the observed regularities in its current form, variants on this model might well be able to explain these observations and thereby contribute to our understanding of the internal structure of civil conflict and to efforts to prevent, predict, and/or control it.

Chapter 5: Conclusion

In the course of this dissertation, we have explored the use of adaptive agent modeling in policy relevant contexts. In three diverse cases, we have argued that the method is capable of generating useful, policy relevant results and that it provides a tool for exploring aspects of social systems which are often overlooked in quantitative analysis because of the inadequacy of traditional tools. In this chapter, we will review the major findings of each of these cases and proceed to examine the meaning of these findings for the use of the adaptive agent method in a more general way.

Contributions of the cases

Trade

In chapter two, this dissertation makes primarily theoretical contributions in the area of international trade. First, we “docked” an adaptive agent model with an analytical model from the recent literature on the subject [Samuelson, 2004], finding that the adaptive agent model produced results which are in line with those produced by analysis. This served to demonstrate that when the model is instituted with the same assumptions as a traditional trade model it produces the same result. We then took an additional step by relaxing the assumption of constant or decreasing returns to scale and demonstrated that the model is capable of reproducing Gomory and Baumol’s [2000] result with regard to increasing returns to scale and the importance of history and policy in development and trade. We went on to discuss the suitability of the model for relaxing other assumptions of traditional trade models, including

those of capital immobility, consumption as the sole determinant of welfare, and the homogeneity of people both as laborers and as consumers.

While we did not attempt to ground these theoretical points in empirical data, we did compare our results with other policy oriented works, particularly those of Samuelson, Gomory and Baumol, and Daly. We argued that the adaptive agent approach provides an additional tool for examining issues in trade policy, that it can be used to lend support to existing arguments that there is a place for policy in trade, and that the method may be uniquely well suited to examine the class of questions where the heterogeneity of workers, consumers, and industries plays a significant role.

We argued that the adaptive agent approach provides a platform for rigorous, quantitative work which is able to relax the standard assumptions of economics. These assumptions (e.g. decreasing returns, perfect rationality, representative agents, etc.) are generally adopted because they produce analytical tractability – not because they are universally applicable. In some cases, these assumptions are harmless abstractions from reality, but in others they lead to genuinely misleading results. With regard to international trade, we demonstrated that the agent method can be used to show that the assumption of no increasing returns to scale leads to important mistakes in thinking about trade and development. We also discussed (but did not demonstrate) how the method could be used to show that the assumption of internationally immobile capital leads to other important problems in thinking about international commerce.

Cities

In chapter three of this dissertation, our theoretical contributions toward understanding the size distribution of cities in the United States, France and Russia were grounded firmly in data. The adaptive agent modeling perspective allowed us to present a very simple model of human migration under bounded rationality which was able to explain not only the tendency of national city size distributions to approximate the Zipf distribution, but also was able to account for the deviations from Zipf which are present in all three countries. While we do not present this model as the last word on the subject, we do see it as a significant contribution to understanding a mystery which has intrigued economic geographers for over fifty years. This model is simpler than most of its predecessors and makes stronger and more accurate predictions than any of them.

By separating core size from observed size, this conceptualization of urban size dynamics is able to account for the observed sizes of cities while remaining compatible with existing work in economic geography which seeks to explain urban agglomerations in terms of central place theory and increasing returns (e.g. Fujita, Krugman, and Venables [1999]).

By recognizing that the process of internal migration involves the conservation of urban population, the model is able to explain considerably more of the shape of the observed distributions than did previous statistical models. These models (e.g. that proposed by Gabaix [1999]) have sought to explain why city distributions approximate the Zipf distribution. They did this by assuming that cities have growth rates which are independent from one another and offering a statistical

process with little behavioral content. We have gone substantially beyond this, offering simple and plausible behavioral rules with logic that leads not only to approximations of Zipf, but to systematic deviations from Zipf. Our results indicate that the true attractor of national city size distributions is not Zipf, but a class of distributions of which Zipf is only one. This allows us to explain the observed deviations from Zipf as signal rather than noise, opening the door for policy insights and interventions with regard to urban size structure.

The two major areas where we offer preliminary policy suggestions are the control of third world megacities and the management of transition in the post-Soviet Russian urban system. The model suggests that policies aimed at reducing pressure on megacities by shifting development to second tier cities are likely to continue to fail. In contrast, policies geared toward providing the physical and social infrastructure that would allow smaller places to become functional parts of the urban system offer more promise.

In Russia, we explained the odd distribution of cities as a result of centrally designed policies of the Soviet era. The model suggests that the Russian urban structure may shift substantially as the last vestiges of these policies are removed. In the course of this transition, Moscow and St. Petersburg are likely to grow, while all but a handful of Russia's mid-sized industrial cities are likely to shrink. The model further suggests that Russian cities with sizes between 1,000,000 and 100,000 are likely to shrink in both size and number. Because the model suggests that there is only a loose coupling between "economically rational" size and observed size, we further suggested that the Russian government might do well to institute policies to

further this apparently inevitable transition. As this transition goes forward, it seems unwise to subsidize the many Russian industrial cities which are in climatically inhospitable locations with poor access to markets at the expense of cities in more viable locations.

Conflict

In the course of chapter four, where we examined data from the civil conflict in Guatemala, our contributions were primarily empirical. We observed that the relationship between ethnicity and violence during the conflict was complex, with more violence taking place in areas where one ethnic group made up between 75% and 90% of the population. This contrasted with significantly lower levels of violence where the population was either more or less balanced between the two ethnic groups. We observed that outbreaks of violence followed a pattern of “punctuated equilibrium” which resembles the progress in time of other, better understood, complex systems. Finally, we observed that the sizes of the incidents of the non-genocidal part of the conflict followed the Zipf distribution, whereas the sizes of incidents from the genocidal part of the conflict were distributed differently. In developing this result, we introduced a novel method for estimating a power-law exponent from grouped data.

These empirical observations have the potential to be important in several ways. First, if they prove to be general results, they could be directly useful for managing and understanding conflicts. The literature on the relationship between ethnic proportions and propensity to civil violence has produced mixed results in large part because such studies are often based on national level data. We demonstrated

that, while Guatemala is nearly evenly divided between Mayans and Ladinos overall, very few of its municipalities are divided this way. This means that very few Guatemalans experience an evenly split population. Instead, the daily experience of most Guatemalans is of being either in the majority or in the minority. We found that the great majority of violence was contained in the 8% of towns where Mayans made up between 80% and 90%. If this observation holds up in looking at other conflicts, it could be useful in targeting efforts to diffuse potential violence and in placing peacekeepers.

The observation that genocidal and non-genocidal conflicts produced differently distributed violent events also has the potential for direct use in post conflict reconstruction as a society seeks to come to terms with a bitter conflict and to deal honestly and justly with those who were involved in it. While it is hard to imagine the distribution of incidents serving as the key evidence against a war criminal, it seems useful to have many different ways of characterizing a genocide. If it proves to be a general result, this technique could provide yet another piece of evidence in establishing a pattern of genocide once a conflict has concluded.

Though the direct uses of the empirical observations developed in this chapter have the potential to be important, they are far from established as general rules. Though further research may prove that they are robust enough to be of direct, practical use, their utility in this regard remains to be established.

A more immediate use of these observations is as benchmarks for evaluating adaptive agent models of conflict. Though this data set has been examined in depth by highly innovative statisticians (e.g. Patrick Ball [1999]), previous analysts did not

observe the patterns noted here. This is largely because our analysis began as an effort to provide an empirical grounding for the development of adaptive agent models of civil violence. The dynamic, disaggregated and bottom-up perspective required by the adaptive agent approach lead us to look for different types of patterns in the data. These patterns have more to do with the statistical “texture” of the dynamics of the conflict than they do with predicting outcomes.

These textures (i.e. the non-linear relationship between ethnic mix and violence, the temporal pattern of incidents, and the distribution of incident sizes) are the kinds of quantitative phenomena to which the adaptive agent method is uniquely suited to address. These would seem to be “emergent phenomena”. Though it would be very hard to predict these patterns by looking at the parts of the system (however these parts might be defined: individuals, political parties, economic forces, etc.), it seems likely that they are logical outgrowths of the way that these parts interact. The patterns are not so much qualities of the individuals involved as they are qualities of logic of the system as a whole.

If these observations have any degree of generality, a successful adaptive agent model of civil conflict (or at least the Guatemalan conflict) should be able to reproduce them. The observations should, therefore, be valuable in guiding efforts to understand the complex interactions which underlie this kind of violence. These observations provide a small, but potentially important step toward establishing an improved set of conflict models which might lead to real progress in the preservation of the global peace.

The meaning of the Zipf distribution

While both city sizes within nations and incidents of violence in a non-genocidal counterinsurgency are distributed approximately according the Zipf rule, it is important not to read too much into this. Various authors have attributed significance to the Zipf distribution as a signature of a complex system and, in some cases, have attributed ill-defined normative significance to it. The fact that this distribution appears in two of the three cases discussed here provides an opportunity to discuss its causes and meaning. While these distributions appear similar, when we looked into the possible origins of the distribution in both cases, we found a great deal of difference.

In describing city size distributions, we gained a great deal of explanatory power by assuming a conserved (or constantly growing) urban population which migrated between cities according to ruled dominated by bounded rationality. In this formulation, Zipf has no normative significance. This is to say that there is nothing in our formulation to suggest that Zipf distributed cities are more efficient or in any way more desirable than alternative distributions. Though we theorized that they are a product of free markets and free mobility, it is important to remember that our process is not one of optimization. In a world of perfect information, where everyone could be expected to accurately seek optimality at all times, our model would predict stability of the existing city structure, whatever that might be, not the emergence of Zipf and its variants.

We do suggest that nations with disproportionately large megacities might have reasons to prefer a more Zipf-like distribution to their current configuration, and that there are policies which might be more or less successful in bringing this about.

This is, however, because of the recognized environmental and social problems associated with megacities, not because a smooth Zipf distribution is somehow more harmonious. We see no particular advantage in the fact that these policies are likely to make the nation's city distribution adhere more closely to Zipf.

Similarly, in the case of Russia, we suggest that the decay of Soviet era policies designed to produce cities that were evenly sized and evenly distributed in space is likely to result in a more Zipf-like distribution – but our assertion here is strictly positive. There are copious reasons why Russia might want to alter the sizes of its cities (as documented by Hill and Gaddy [2003] among others), but these reasons do not have to do with a size distribution that is somehow “bad”. We do, however, observe that the positive prediction that the Russian city size distribution is likely to change in the absence of restrictive Soviet policies provides an opportunity to pursue the normative desire to restructure the Russian urban distribution with a minimum of pain. If the Russian people can predict where their city size distribution is headed, they will be in a better position to get there with fewer missteps.

In the Guatemala case we also identified an important feature, the incidents of killing in a non-genocidal counterinsurgency, as being Zipf distributed. As in the case of city sizes, we theorize that this distribution is the emergent result of a complex process, but at that point the resemblance between the two phenomena ends. With cities, we contend that the migration of urban population between cities under conditions of bounded rationality leads to a range of distributions of which Zipf is a special case. In the case of civil violence, we have less success in modeling the phenomenon, but we do not see a strong analogy between the mechanisms. The

number of dead is not conserved and the dead do not migrate between incidents. It is hard to imagine a parallel to the interurban migration citizens in a model of incidents of deadly violence.

While it is probably fair to say that both the evolution of city size distributions and the evolution of incidents of violence involve complex dynamics involving a balance between positive and negative feedbacks, that fact leaves a great deal unexplained about these phenomena. We have tried to go beyond identifying the similarity in distributions to provide insight into the mechanisms that might have produced them. In doing so, we have found the mechanisms to have little in common beyond the fact that they involve the interaction of heterogeneous parts.

General implications of the dissertation for policy research

The focus of this dissertation has been to demonstrate the utility of the adaptive agent method in examining issues with policy relevance. We have sought to do this by using it to demonstrate useful, novel results in diverse areas of application. Having staked our admittedly modest, but potentially important claims in this regard, we will look now at what the adoption of the method, and the habits of thought which it requires, might mean for the conduct of policy research in general.

Traditional methods of quantitative modeling have tended to be characterized (in broad terms) by:

- Unique, non-path dependent equilibria.
- The use of representative agents as a proxy for homogeneous populations of agents.
- Assumptions of perfect information and rationality which are technically difficult to relax.

- Elaborate analytical techniques to enforce conserved quantities.
- Limited insight into distributional impacts (stemming largely from the use of representative agents).

Through the examples presented here, we have demonstrated that the adaptive agent approach differs fundamentally from more traditional quantitative methods in that it is fundamentally suited to:

- Modeling path dependent processes where the history of the system matters. (Particularly relevant in the chapter on trade)
- Modeling individual based processes where the heterogeneity of actors matters. (Particularly relevant in the chapter on civil violence)
- Modeling situations where bounded rationality and imperfect information are fundamental to the process under study. (Particularly relevant in the chapters on cities and civil violence)
- Managing conserved quantities. (Relevant in all three cases)
- Examining distributional impacts of changes in process or policy. (Relevant in all three cases)

When a problem is to be treated quantitatively, the way that the problem is conceived must be constrained by what can be done with the quantitative methods that are at hand or can be readily devised. The fact that the adaptive agent method allows for these various restrictions to be simultaneously lifted allows for an expanded scope for quantitative work in social science in general and policy analysis in particular.

By expanding the range of applicability of quantitative methods, the introduction of adaptive agent methods expands the way that problems can be conceived. Before a problem can be formally analyzed, the analyst must form a mental picture of the system – a pre-analytic vision. Armed with a set of tools which

produce powerful results when assumptions about homogeneity, linearity, etc. are upheld, but are much less useful when they are not, there is a tendency to build the pre-analytic vision from the parts of the problem which will lend themselves to analysis. The pre-analytic vision of those who work with numbers has therefore tended to be characterized by the same features that characterize the available methods. This has led to a systematic under exploration of path dependent systems, the importance of individual differences, the effects of bounded rationality, the recognition of limits, and the distributional effects of policy.

While the adaptive agent method is still in the early stages of its development and lacks both the completeness of analytical proof and the methodological refinement of modern econometric methods, its has great potential for expanding the range of pre-analytic vision with regard to policy. The method does not allow insight to be derived from muddled questions (no method can hope to do that), but does allow a new class of questions to be asked and answered in a rigorous quantitative context. Because adaptive agent methods allow researchers to explore issues that were previously relatively intractable, they allow this class of problems to be brought into the realm of quantitative analysis.

The existence of adaptive agent modeling allows for a richer pre-analytic vision which takes account of history, social organization, and human diversity. Such thinking has, of course, always been an important part of policy discourse. The fact that many of these problems did not lend themselves to quantitative analysis by traditional methods, however, meant that they had to be treated verbally and could not easily be combined with quantitative results. The rise of adaptive agent modeling

offers the potential to merge these strains of policy discourse, bringing some of the rigor of quantitative analysis to subjects that have only been treated in words, and bringing some of the richness and subtlety of philosophical discourse to quantitative models.

Bibliography

- Ades, Alberto F & Glaeser, Edward L, 1995. "Trade and Circuses: Explaining Urban Giants," *The Quarterly Journal of Economics*, MIT Press, vol. 110(1), pages 195-227
- Armano Srbljinovic, Drazen Penzar, Petra Rodik and Kruno Kardov. 2003. "An Agent Based Model of Ethnic Mobilisation," *Journal of Artificial Societies and Social Simulation* vol 6. no. 1 <http://jasss.soc.surrey.ac.uk/6/1/1.html>
- Arrow, Kenneth. 1962. "The economic implications of learning by doing." *Review of Economic Studies* 29: 155-73.
- Arthur, W. B., "Urban Systems and Historical Path Dependence" in *Urban Systems and Infrastructure*, edited by R. Herman and J. Ausubel, Washington D.C., National Academy of Sciences, 1987.
- Arthur, W. Brian. 1989. "Competing technologies, increasing returns and lock in by historical events." *Economics Journal* 99: 116-31
- Ashby, W. Ross. 1956 *An Introduction to Cybernetics*, Chapman & Hall, London.
- Axelrod, Robert. 1984. *The Evolution of Cooperation*. Basic Books, NY.
- Axtell, R.L. and R. Florida, 2001. "Zipf's Law of City Sizes: A Microeconomic Explanation," (submitted for publication; working paper, The Brookings Institution, www.brook.edu/dynamics/papers/cities).
- Axtell, Robert. 2000. Why Agents? On the Varied Motivations for Agent Computing in the Social Sciences, CSED Working Paper No. 17
- Bak, Per. 1997. *How Nature Works: The Science of Self Organized Criticality*. New York: Copernicus.
- Ball, Patrick. 1999. AAAS/CIIDH database of human rights violations in Guatemala (ATV20.1). <http://hrdata.aaas.org/ciidh/data.html> (July 1, 2000).
- Ball, Patrick, Paul Kobrak, and Herbert F. Spierer. 1999. *State Violence in Guatemala, 1960-1996: A Quantitative Reflection*. <http://hrdata.aaas.org/ciidh/data.html> (July 1, 2000).
- Bates, Robert. 1983. "Modernization, Ethnic Competition, and the Rationality of Politics in Contemporary Africa." in *State Versus Ethnic Claims: African Policy Dilemmas*, ed. Donald Rothchild and Victor A. Olorunsola. Boulder, CO: Westview.

- Bhavnnani, Ravi & David Backer. 2000. "Localized Ethnic Conflict and Genocide: Accounting for Differences in Rwanda and Burundi," *Journal of Conflict Resolution* 44(3):283-306.
- Bugliarello, G. 1999. Megacities and the developing world. *The Bridge* 29(4):19-26.
- Cederman, Lars-Erik. 2003. "Modeling the Size of Wars: From Billiard Balls to Sandpiles." *American Political Science Review* 97: 135-150.
- Chavouet, J. M. and J. C. Fanouillet, 2000. Forte extensions des villes entre 1990 et 1999. INSEE Paper #707.
- Christaller, Walter. *Die zentralen Orte in Süddeutschland*. Jena: Gustav Fischer, 1933. (Translated (in part), by Charlisle W. Baskin, as *Central Places in Southern Germany*. Prentice Hall 1966).
- Clifford G. Gaddy and Barry W. Ickes, 2002, *Russia's Virtual Economy*, *Brookings Institution Press*, Washington DC, p. 235
- Costanza, R., R. d'Arge, R. de Groot, S. Farber, M. Grasso, B. Hannon, S. Naeem, K. Limburg, J. Paruelo, R.V. O'Neill, R. Raskin, P. Sutton, and M. van den Belt. 1997. The value of the world's ecosystem services and natural capital. *Nature* 387:253-260
- Cuberes, David, 2004, The Rise and Decline of Cities, unpublished manuscript, <http://economics.uchicago.edu/download/cities3.pdf>
- Herman E. Daly, 1996. *Beyond Growth*. Beacon Press, Boston
- Daly, Herman and Joshua Farelly, 2003. *Ecological Economics: Principles and Applications*. Island Press, Washington, DC.
- Davenport, Christian and Patrick Ball, 2002 "Views to a Kill: Exploring the Implications of Source Selection in the Case of Guatemalan State Terror, 1977 – 1995" *Journal of Conflict Resolution*, Vol. 46 No. 3, June 2002 427-450.
- Demko, G.J. and R.J. Fuchs. 1984. "Urban Policy and Settlement System Change in the USSR, 1897 – 1979," in *Geographical studies on the Soviet Union: essays in honor of Chauncy D. Harris*, George J. Demko and Roland J. Fuchs, editors. Chicago, Ill. : Dept. of Geography, University of Chicago.
- Dibble, Catherine. 2001. *Theory in a Complex World: GeoGraph Computational Laboratories*. PhD Dissertation, Department of Geography, U.C. Santa Barbara, CA.
- Doyle, Michael W. and Nicholas Sambanis, 2000. "International Peacebuilding: A Theoretical and Quantitative Analysis." Draft (March 7, 2000).

- Epstein, Joshua M., John D. Steinbruner, Miles T. Parker 2001. "Modeling Civil Violence: An Agent-Based Computational Approach." Brookings Institution Center on Social and Economic Dynamics Working Paper No. 20.
- Epstein, Joshua and Robert Axtell. 1996. *Growing Artificial Societies: Social Science from the Bottom Up*. The MIT Press, Cambridge, MA,.
- Ericson, Richard. 1999, "The Structural Barrier to Transition Hidden in Input-Output Tables of Centrally Planned Economies." *Economic Systems* 23(3): 199-244
- Fearon, James D. and David D. Laitin, 1996. "Explaining Interethnic Cooperation," in *American Political Science Review* 90:4, December 1996.
- Federal Register Vol. 65, No. 249, December 27, 2000
- Florida, Richard. 2002. *The Rise of the Creative Class: And How It's Transforming Work, Leisure, Community and Everyday Life*. Basic Books, New York.
- Fujita, Krugman, and Venables, 1999. *The Spatial Economy: Cities, Regions, and International Trade* Cambridge, MA: The MIT Press.
- Fujita, Masahisa and Mori, Tomoya, 1997. "Structural Stability and Evolution of Urban Systems", *Regional Science and Urban Economics*, August 1997, Vol. 24, No. 4-5, pp. 399-442.
- Gabaix, Xavier, 1999b. "Zipf's Law and the Growth of Cities", *American Economic Review Papers and Proceedings*, 89 (2), pp. 129-32.
- Gabaix, Xavier. 1999. "Zipf's law for cities: an explanation", *Quarterly Journal of Economics*, 114:739-767.
- Gibrat, Robert. 1931. *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*. Paris: Librairie du Recueil Sirey.
- Gulden, Timothy R. 2002. "Spatial and Temporal Patterns in Civil Violence: Guatemala 1977-1986" *Politics and the Life Sciences*. Vol 21, No 1. (March).
- Gulden, Timothy R. 2002. "Spatial and Temporal Patterns in Civil Violence: Guatemala 1977-1986" Brookings Institution Center on Social and Economic Dynamics Working Paper No. 26.
- Gomory, Ralph E. and William J. Baumol. 2000. *Global Trade and Conflicting National Interests*. Cambridge, Mass.: MIT Press.
- Guatemalan Commission for Historical Clarification (CEH). 1999. *Guatemala: Memoria del Silencio*. <http://hrdata.aaas.org/ceh/> (July 1, 2000).

- Gurr, Ted Robert and Barbara Harff, 1996. *Early Warning of Communal Conflicts and Genocide: Linking Empirical Research to International Responses*. United Nations University Press.
- Hannon, Bruce and Matthias Ruth. 2001. *Dynamic Modeling: 2nd Edition*. Springer-Verlag, New York
- Henderson J.V., 1974, 'The Sizes and Types of Cities', *American Economic Review* 64, 640-656.
- Hill, B. and Woodrofe, M. 1975. "Stronger forms of Zipf's law." *Journal of the American Statistical Association*, 70(349):212
- Hill, Fiona and Clifford Gaddy. 2003. *The Siberian Curse: How Communist Planners Left Russia Out in the Cold*. Brookings Institution Press, Washington, DC.
- INSEE, 2004, http://www.recensement.insee.fr/EN/RUB_MOT/ACC_MOT.htm
- Iyer, Seema D. 2003. "Increasing Unevenness in the Distribution of City Sizes in Post-Soviet Russia." *Eurasian Geography and Economics*. 44 No. 5 pp. 348-367
- Jacobs, Jane, 1984, *Cities and the Wealth of Nations: Principles of Economic Life*. New York: Random House.
- Julien, Phillippe. 2001a. Les Grandes Villes Francaises Etendent Leur Influence. ISEE Paper No. 766.
- Julien, Phillippe. 2001b. Les Deplacements Domicile-Travail: de plus en plus d'actifs travaillent loin de chez eux. ISEE Paper No. 767 (translated in part as "More and more workers are working far from home")
- Kanemoto, Y., 1980. *Theories of Urban Externalities*, North-Holland.
- Kewley, MAJ Robert and LTC Larry Larimer. 2003. "An Agent Based Modeling Approach to Quantifying the Value of Battlefield Information," *PHALANX*, June 2003, Volume 36 Number 2
- Krugman, Paul R. 1979. "Increasing returns, monopolistic competition and international trade." *Journal of International Economics*. 9:469-79
- Krugman, Paul R. 1983. "Targeted industrial policies: Theory and evidence." *Industrial Change and Public Policy*. Kansas City, MO: Federal Reserve Bank, pp. 123-55.
- Le Gléau, Jean-Pierre, Denise Pumain and Thérèse Saint-Julien. 1996. "Villes d'Europe : à chaque pays sa definition", *Économie et Statistique* , no.294-295,1996-4/5 (in English as "Towns of Europe: to each country its definition")

http://www.tu-chemnitz.de/phil/geographie/material/WS_03_04/BK/villeseuro.pdf

- Lotka, A. J. 1925. *Elements of physical biology*. Baltimore: Williams & Wilkins Co.
- Lowry, Ira S. 1966. *Place to Place Migration Flows. In Migration and Metropolitan Growth: Two Analytical Models*. San Francisco: Chandler Publishing.
- Marshall A., 1890, *The Principles of Economics*, McMillan (N.Y.), re-edition 1925.
- Meadows, D.H., D.L. Meadows, J. Randers, and W.W. Behrens III, 1972. *The Limits to Growth*, Universe Books, New York.
- Overman, Henry G.; Ioannides, Y. M. 2001. "The Cross-sectional Evolution of the US City Size Distribution." *Journal of Urban Economics* 49 (2001), pp. 543-566.
- Pumain D. 2004, "Scaling laws and urban systems." Santa Fe Institute, Working Paper n°04-02-002
- Reed, Bill. 2002, "On the rank-size distribution for human settlements", *J Regional Science*, 41:1-17.
- Ricardo, David. 1817. *On the Principles of Political Economy and Taxation*
- Richardson, Lewis F. 1960. *Statistics of Deadly Quarrels*. Chicago: Quadrangle Books.
- Rosen, Sherwin. 1974. "Hedonic Prices and Implicit Markets: Product Differentiation in Pure Competition," *Journal of Political Economy*, Vol. 82, Jan./Feb. 1974, pp. 34-55.
- Ruth, Matthias and Bruce Hannon. 2001. *Modeling Dynamic Economic Systems*. Springer-Verlag, New York
- Samuelson, Paul A. 2004. "Where Ricardo and Mill Rebut and Confirm Arguments of Mainstream Economists Against Globalization." *Journal of Economic Perspectives*. Vol. 18, No. 3. Summer 2004. pp 161-180.
- Schelling, T. C. 1960, *The Strategy of Conflict*. Cambridge, Mass.: Harvard University Press.
- Schelling, T. C. 1969, Models of Segregation. *American Economic Review. Papers and Proceedings*. 59. 488-493.
- Schroeder, Manfred. 1991. *Fractals, Chaos, Power Laws*. New York: W. H. Freeman and Company.

- Simon, Herbert A., 1957. *Models of Man: Social and Rational*. New York: John Wiley and Sons, Inc.
- Steinbruner, John D. 2000. *Principles of Global Security*. Washington, DC: Brookings.
- Sterman, J. 2000. *Business Dynamics: Systems Thinking for a Complex World*. Irwin/McGraw-Hill
- Taylor, Glenn, Richard Frederiksen, Russell R. Vane III and Edward Waltz. 2004. "Agent-based Simulation of Geo-Political Conflict," in Deborah L. McGuinness, George Ferguson (Eds.): *Proceedings of the Nineteenth National Conference on Artificial Intelligence, Sixteenth Conference on Innovative Applications of Artificial Intelligence*, July 25-29, 2004, San Jose, California, USA. AAAI Press / The MIT Press 2004, 884-891
- UN-Habitat, 2004. *The State of the World's Cities 2004/2005*. ISBN: 92-1-131705-3. Earthscan Publications Ltd, London and Sterling
- UNIDO, 2004. *Industrial Development Report 2004*, Vienna: United Nations Industrial Development Organization.
- UN-Population, 2001. *World Urbanization Prospects: The 2001 Revision*, United Nations Population Division
- Volterra, V. 1926. Variazioni e fluttuazioni del numero d'individui in specie animali conviventi. *Mem. R. Accad. Naz. dei Lincei*. Ser. VI, vol. 2.
- von Bertalanffy, Ludwig. 1968. *General Systems Theory*, George Braziller, New York
- Zipf, George Kingsley. 1949. *Human Behavior and the Principle of Least Effort*. Addison-Wesley.