

Editorial: The Future of Power Law Research

Michael Mitzenmacher

Abstract. I argue that power law research must move from focusing on observation, interpretation, and modeling of power law behavior to instead considering the challenging problems of validation of models and control of systems.

I. The Problem with Power Law Research

To begin, I would like to recall a humorous insight from the paper of Fabrikant, Koutsoupias, and Papadimitriou [Fabrikant et al. 01], consisting of this quote and the following footnote.

“Power laws ... have been termed ‘the signature of human activity’...”¹

The study of power laws, especially in networks, has clearly exploded over the last decade, with seemingly innumerable papers and even popular books, such as Barabási’s *Linked* [Barabási 02] and Watts’ *Six Degrees* [Watts 03]. Power laws are, indeed, everywhere. Despite this remarkable success, I believe that research into power laws in computer networks (and networks more generally) suffers from glaring deficiencies that need to be addressed by the community. Coping with these deficiencies should lead to another great burst of exciting and compelling research.

To explain the problem, I would like to make an analogy to the area of string theory. String theory is incredibly rich and beautiful mathematically,

¹“They are certainly the product of one particular kind of human activity: looking for power laws...” [Fabrikant et al. 01]

with a simple and compelling basic starting assumption: the universe's building blocks do not really correspond to (zero-dimensional) points, but to small (one-dimensional) strings. Many different versions and variations of string theory exist. String theory is an incredibly popular area of physics, both with researchers and with the general populace. The mathematical richness, simple starting points, large number of variations, and popularity are all characteristics also shared by current research in power laws.

There is, however, an unpleasant problem with string theory: it has not been experimentally validated, and it is not currently clear what will be required to validate it. There has not been an actual experiment that shows that string theory successfully predicts something that is not also explained by other existing theories, and current suggestions for potential experiments to pass this hurdle appear well outside the realm of what is possible in the foreseeable future. In an era where evolution and intelligent design have set the stage for the question "What is science?", the issue that there is apparently no adequate way to verify or falsify string theory is rightly seen as a major flaw. This problem may be corrected in time, and I do not wish to suggest that string theory is not a valuable scientific contribution. The current state, however, creates a litany of challenging questions about the fruitfulness of this line of research.

I would argue that we are now facing a similar problem in research on power laws in computer networks, and this signals that we must actively change the research agenda. Specifically, while numerous models that yield power law behavior have been suggested and, in fact, the number of such models continues to grow rapidly, no general mechanisms or approaches have been suggested that allow one to validate that a suggested model is appropriate. Like the string theorists, we have beautiful frameworks, theory, and models—indeed, we have perhaps far too many models—but we have been hesitant in moving to the next steps, which could transform this promising beginning into a truly remarkable new area of science.

To make my case, let me introduce a rough taxonomy of the types of results that one could aim for in studying power laws.

- (a) *Observe*: Gather data on the behavior of a system and demonstrate that a power law distribution appears to fit the relevant data.
- (b) *Interpret*: Explain the significance of the power law behavior to the system.
- (c) *Model*: Propose an underlying model that explains the power law behavior.
- (d) *Validate*: Find data to validate, and if necessary specialize or modify, the model.

- (e) *Control*: Use the understanding from the model to control, modify, and improve the system behavior.

My argument is that most research on power laws has focused on observing, interpreting, and modeling, with a current emphasis on modeling. As a community, we have done almost nothing on validation and control, and we must actively move towards this kind of research. Unlike string theory, where validation and control are hampered by nature and constrained by what is physically possible, I believe that in networks we can, in most cases, successfully validate our models and use them to control systems.

As an example, consider a power law that has been the subject of significant research: the distribution of in-links and out-links on the directed graph corresponding to the pages of the World Wide Web (see, e.g., [Barabási et al. 99, Broder et al. 00, Kleinberg et al. 99]). Early work on the web graph discovered that these distributions followed a power law, and the significance for applications was clear. In-degree was used early on by search engines such as AltaVista to rank web pages, and the power law distribution ensured that this approach was reasonably useful, as only a small number of pages were likely to have a significantly larger number of in-links.

A significant amount of historical and recent theoretical work has led to models for why power law behavior occurs, as detailed in various surveys [Mitzenmacher 04, Newman 05]. For the web graph, the primary model of study has been variations of preferential attachment: the more links a web page has, the more links it is likely to obtain in the future. But even for this specific problem, there has been little or no systemic or theoretical work covering validation or control. Validation is important, because there are numerous models that yield power laws, and in many cases more than one of these models can reasonably be used to explain an observed power law behavior. For example, models other than preferential attachment have been used to describe link behavior, and some evidence has been given that preferential attachment is not a suitable explanation (see, e.g., [Huberman and Adamic 99]). On the issue of control, it now seems clear that the presence of search engines (especially Google) has affected how people link to web pages. There is less need for people to explicitly put in links in their web pages when the links are easily found on Google; however, if people reduce how often they link to pages, this might affect the performance of Google's algorithms. Could this impact and its effects have been modeled appropriately ahead of time? (A model for this situation was examined in [Chakrabarti et al. 05].)

This path of observation, interpretation, and modeling has been repeated in many domains. I would argue, however, that without validating a model it is

not clear that one understands the underlying behavior and therefore how the behavior might change over time. It is not enough to plot data and demonstrate a power law, allowing one to say things about current behavior; one wants to ensure that one can accurately predict *future* behavior appropriately, and that requires understanding the correct underlying model.

For example, barring major changes, will the degree distribution of web pages still look the same two years from now? Recent work suggests that the density of edges in various power law graphs, such as citation graphs and the Autonomous Systems graph, may be increasing over time [Leskovec et al. 05]. If the same holds in the graph of web pages, the answer may be no. At this point we do not know enough about the web graph to answer the question. Further, from the practical point of view, one would like to take the model and use it to control the system. Suppose that Google could somehow encourage people to create more high-quality links, and it wanted to determine if it could do so in a way that was worth the cost. This would require not only validating the underlying model but also understanding how the model would respond to control mechanisms, such as advertising, contests, or direct payments, that might increase links.

2. New Directions

I can happily say that I believe that the community is currently moving toward validation. A key step in this direction, perhaps between modeling and validation, is work on *invalidation*. Invalidation can take multiple forms, including showing that a model is insufficient to explain a power law by providing a reasonable or more compelling alternative, or considering additional properties beyond the power law behavior suggested by a proposed model and showing that these properties do not fit the data. For example, an insightful recent paper by Lakhina, Byers, Crovella, and Xie [Lakhina et al. 03] questions the methodology of using *traceroute*-based maps (as in, for example, the highly-cited [Faloutsos et al. 99]) of the Internet to conclude that the degree distribution of the Internet-connectivity graphs follows a power law. With *traceroute*-based maps, one attempts to obtain a snapshot of the edges in the network by tracing the route of packets from one or more specialized sources and various destinations in the network. Unfortunately, as shown in [Lakhina et al. 03], this leads to a very biased sample of the edges in the network; they show that under this methodology the edges found even from random regular graphs would yield a subgraph that also exhibits power law behavior. This is because there is an inherent correlation between the observed degree of a vertex in the resulting subgraph and the proximity of the source of the *traceroute* queries. This paper does not conclude

that the Internet does not have such a degree distribution, but it does provide reason to believe that the methodology needs to be reexamined. Similarly, the paper by Chen et al. [Chen et al. 02] considers several different characteristics of the router graph and finds that models based on preferential attachment do not have the appropriate properties that this graph seems to have. It then proposes a new model, which appears to have these additional properties.

The article in this issue by Li, Alderson, Doyle, and Willinger [Li et al. 05] is in part a paper on invalidation, although it is also a great deal more. A major part of the paper focuses on clarifying that there are many different types of graphs that yield power law distributions for their degree sequences, and these graphs may appear extremely different with respect to other properties. They argue that many people have applied certain types of graphs, such as graphs created by preferential attachment, to describe real objects, such as the Internet connectivity graph, based on the fact that both have a power law degree distribution. They show by an examination of the data that this characterization does not match the real world. Beyond such invalidation arguments, they also provide a large framework and foundations that will certainly play a role in future invalidation and validation arguments. In particular, they attempt to give a simple but useful way of distinguishing various types of graphs with heavy-tailed in-degree and/or out-degree distributions and to build a proper vocabulary for making appropriate distinctions. I believe that this paper should be and will be read by everyone working in the area. It will provide a healthy jumping-off point for a richer, more refined view of the various power law models and their relation to real structures.

This ad hoc approach based on invalidation, whereby we repeatedly find alternatives or flaws in existing models and adapt accordingly, has already been useful in developing models with appropriate properties for specific experiments and in bringing important variations to light. It is, however, incomplete and unsatisfying both theoretically and practically. It is never clear that one has captured all relevant statistics appropriately. Invalidation leads us naturally to consider whether validation is possible from the outset.

What other approaches to validation are there? A compelling possibility is to make use of *time series analysis*, with the goal that one observes the system over time to judge the underlying assumptions of the proposed model. That is, instead of using the model to generate samples and see if they appear to have the same features from some limited feature list as the object being studied, one examines the dynamics of the object being studied to see if it fits the model. Naturally, another advantage of this approach is that time series analysis might suggest variations or outright changes in the underlying model that would lead to better modeling. While time series analysis is certainly at least implicit in other works,

I am not aware of any clear theoretical formulation in the setting of analysis of power law graphs. An appropriate theoretical framework is imperatively needed.

In most cases, time series analysis for computer network problems will require some type of *trace-based analysis*, obtained by monitoring the system in some way. Unfortunately, trace-based analysis is often difficult and time-consuming, with many potential problems. The large quantity of data that must be gathered is one obvious drawback. Heterogeneity in the network, including heterogeneity over users and heterogeneity in behavior over time, is another difficulty. Inherent noise may be a challenge. Finally, when dealing with existing networks, relevant past data may not have been captured, making it difficult to determine if the model was appropriate from the beginning.

Sampling provides efficient means of coping with some of the problems of large traces. There are two natural sampling approaches in this setting. The first is to consider global traces at discrete time steps, or *snapshots*. The second approach is to capture more detailed information at a subset of the sites. Which type of sampling is appropriate may depend on the object being studied; in some cases both would be appropriate, and they may reinforce each other.

Coping with these various problems will require a rich collaboration between theorists and system builders, as implementing appropriate sampling, recording, and analysis tools for such analyses will require sophisticated understanding of both the underlying theory and systems.

While time series analysis with trace-based analysis may be one path toward validation, there should certainly be others. The underlying question is whether or not we can break down sufficiently the systems that we are studying so that we can successfully map actions being taken onto steps in the model. It seems to be a goal within reach and certainly a challenge that should be tackled.

There is still the further issue of finding ways to control system behavior. In many natural sciences, simply observing a power law is a significant result, and there is little hope of changing the underlying resulting power law. For example, the size of earthquakes roughly follows a power law distribution, and without significant human intervention in the movements of the earth's crust, there is little hope that there is much we can do to affect it. In computer systems, however, we might expect that in many cases we can change the system behavior, potentially modifying the causes that give rise to the power law in order to better engineer the system. In our world, we can to some extent control the laws of nature, and we need to know how to use that power effectively.

There are several basic means of implementing control. For example, we might find additional constraints to add to the system. Imagine, for example, that we introduced a geometric restriction, allowing links only of a certain length. This is the idea behind the geometric preferential attachment model studied in [Flaxman

et al. 04]; links arise only between nodes sufficiently close in some suitable metric space. Rather than adding direct constraints on the users or the system, one might achieve similar effects through the softer approach of adding incentives or costs to the system. Such constraints arise in the various HOT (Highly Optimized Tolerance, Heuristically Optimized Tradeoffs) proposals [Carlson and Doyle 99, Fabrikant et al. 01], but the effects of trying to control the system through such constraints remain open. The study of how to affect system behavior would also fit nicely into the growing area of *distributed algorithmic mechanism design*, as introduced in [Feigenbaum et al. 02, Feigenbaum and Shenker 02].

The truth is, however, that at this point our understanding of the actual processes that generate power laws in networks is so limited that we do not know what kinds of control mechanisms are suitable. The foundations of this area are still wide open. In addition, by gaining a better understanding of the mechanisms that underly power laws and how they might be modified in the computer network domain, we might also find approaches that work in other domains, including social networks, economics, or even biological systems. In this way, our insights might have a powerful impact outside the field of computer networks, as well.

3. Conclusion

The proposed agenda is that we must move from observation, interpretation, and modeling toward validation and control in power law research. Validation holds the promise of making power law research more scientifically sound; control holds the promise of making it more directly useful. Recent signs suggest that we have already started moving along this path, as people have begun to realize that claims being made for many power law models go beyond what can reasonably be shown, and papers invalidating previous models are coming to the forefront. The article in this issue by Li, Alderson, Doyle, and Willinger is an outstanding example of this, introducing both what has been done in the past and research directions for the future [Li et al. 05].

The idea behind this agenda yields many interesting possible corollaries on which the community must build a consensus. For quite a while, an NP-completeness result has generally not been considered worthy of publication, unless the result was important for a specific field, especially challenging mathematically, or led to new mathematical or computational insights. I would suggest that the same standard should now hold for a large part of the research in power laws. Already, observing a power law in itself no longer seems sufficient for publication in networks, except in rare circumstances; usually some proposed model

is expected. With the plethora of papers on modeling available, perhaps the bar will start to rise for such papers. At the very least, I would argue that modeling papers should include as a matter of course a proposal for how the model could potentially be validated, even if the authors do not themselves tackle the challenge of gathering and processing the data to perform the validation. Or perhaps papers giving new variations of existing models should be required to explicitly invalidate the old models, using real data.

While this proposed agenda certainly raises a number of challenges, it creates many opportunities and research problems. Research in power laws has already brought together computer scientists from theory and systems in remarkable ways, as well as researchers from many other communities. The agenda of validation and control promises to continue and possibly even enhance these collaborations. The agenda also raises the possibility that this research can begin actively affecting real systems, making them more efficient and better designed. The past decade has been very exciting, with many remarkable discoveries in the theory and practice of power laws. I believe that moving the agenda toward validation and control will lead to continuing excitement from this research area.

Acknowledgments. Thanks to David Parkes for several interesting discussions on this subject, and thanks to Scott Aaronson, John Byers, and Fan Chung for comments on earlier drafts. The author is supported in part by NSF grant CCR-0121154.

References

- [Barabási 02] A.-L. Barabási. *Linked: How Everything is Connected to Everything Else and What It Means*. New York: Perseus Publishing, 2002.
- [Barabási et al. 99] A.-L. Barabási, R. Albert, and H. Jeong. “Mean-Field Theory for Scale-Free Random Networks.” *Physica A* 272 (1999), 173–189.
- [Broder et al. 00] A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener. “Graph Structure in the Web: Experiments and Models.” In *Proceedings of the Ninth World Wide Web Conference*. Available from World Wide Web (<http://www9.org/w9cdrom/160/160.html>), 2000.
- [Carlson and Doyle 99] J. M. Carlson and J. Doyle. “Highly Optimized Tolerance: A Mechanism for Power Laws in Designed Systems.” *Physics Review E* 60:2 (1999), 1412–1427.
- [Chakrabarti et al. 05] S. Chakrabarti, A. Frieze, and J. Vera. “The Influence of Search Engines on Preferential Attachment.” In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 293–300. Philadelphia: SIAM, 2005.
- [Chen et al. 02] Q. Chen, H. Chang, R. Govindan, S. Jamin, S. Shenker, W. Willinger. “The Origin of Power Laws in Internet Topologies Revisited.” In *Proceedings of IEEE INFOCOM 2002*, pp. 608–617. Los Alamitos, CA: IEEE Press, 2002.

- 2002, *Málaga, Spain, July 8–13, 2002, Proceedings*, pp. 110–122, Lecture Notes in Computer Science 2380. Berlin: Springer, 2002.
- [Faloutsos et al. 99] M. Faloutsos, P. Faloutsos, and C. Faloutsos. “On Power-Law Relationships of the Internet Topology.” In *Proceedings of the Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication*, pp. 251–261. New York: ACM Press, 2002.
- [Feigenbaum and Shenker 02] J. Feigenbaum and S. Shenker. “Distributed Algorithmic Mechanism Design: Recent Results and Future Directions.” In *Proceedings of the Sixth International Workshop on Discrete Algorithms and Methods for Mobile Computing and Communications*, pp. 1–13. New York: ACM Press, 2002.
- [Feigenbaum et al. 02] J. Feigenbaum, C. Papadimitriou, R. Sami, and S. Shenker. “A BGP-Based Mechanism for Lowest-Cost Routing.” In *Proceedings of the Twenty-First Annual Symposium on Principles of Distributed Computing*, pp. 173–182. New York: ACM Press, 2002.
- [Flaxman et al. 04] A. Flaxman, A. Frieze, and J. Vera. “A Geometric Preferential Attachment Model of Networks.” In *Algorithms and Models for the Web-Graph: Third International Workshop, WAW 2004, Rome, Italy, October 16, 2004, Proceedings*, pp. 44–55, Lecture Notes in Computer Science 3243. Berlin: Springer, 2004.
- [Huberman and Adamic 99] B. A. Huberman and L. A. Adamic. “Evolutionary Dynamics of the World Wide Web.” Technical Report, Xerox Palo Alto Research Center, 1999. Appears as a brief communication in *Nature* 401 (1999), 131.
- [Kleinberg et al. 99] J. Kleinberg, R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. “The Web as a Graph: Measurements, Models, and Methods.” In *Computing and Combinatorics: 5th Annual International Conference, COCOON’99, Tokyo, Japan, July 1999, Proceedings*, pp. 1–17, Lecture Notes in Computer Science 1627. Berlin: Springer, 1999.
- [Lakhina et al. 03] A. Lakhina, J. Byers, M. Crovella, and P. Xie. “Sampling Biases in IP Topology Measurements.” In *Proceedings of IEEE INFOCOM 2003*, pp. 332–341. Los Alamitos, CA: IEEE Press, 2003.
- [Leskovec et al. 05] J. Leskovec, J. Kleinberg, and C. Faloutsos. “Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations.” In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 177–187. New York: ACM Press, 2005.
- [Li et al. 05] L. Li, D. Alderson, C. Doyle, and W. Willinger. “Towards a Theory of Scale-Free Graphs: Definition, Properties, and Implications.” *Internet Mathematics* 2:4 (2005), 431–523.
- [Mitzenmacher 04] M. Mitzenmacher. “A Brief History of Generative Models for Power Law and Lognormal Distributions.” *Internet Mathematics* 1:2 (2004), 226–251.
- [Newman 05] M. E. J. Newman. “Power Laws, Pareto Distributions and Zipf’s Law.” *Contemporary Physics* 46 (2005), 323–351.
- [Watts 03] D. J. Watts. *Six Degrees: The Science of a Connected Age*. New York: W. W. Norton and Company, 2003.

[Watts 03] D. J. Watts. *Six Degrees: The Science of a Connected Age*. New York: W. W. Norton and Company, 2003.

Michael Mitzenmacher, Division of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138 (michaelm@eecs.harvard.edu)

Received November 14, 2005; accepted November 15, 2005.