Background	Theory	Practice	Foreground
		1 61 1	

Mathematics in Medicine Journal Club:

"Power law distributions in empirical data" Aaron Clauset, Cosma Rohilla Shalizi, and M. E. J. Newman

Cory Brunson Center for Quantitative Medicine

March 7, 2017

▲ロト ▲ 理 ト ▲ 王 ト ▲ 王 - の Q (~

Background	Theory	Practice	Foreground

BACKGROUND: SCALE-FREE BIOLOGICAL NETWORKS

Background	Theory	Practice	Foreground

BACKGROUND: SCALE-FREE BIOLOGICAL NETWORKS

- Barabási, A-L and Oltvai, Z N (2004). "Network biology: understanding the cell's functional organization". Nature Reviews Genetics 5(2), 101–113.
- Albert, R (2005). "Scale-free networks in cell biology". Journal of Cell Science 118(21), 4947–4957.
- Zhang, B and Horvath, S (2005). "A General Framework for Weighted Gene Co-Expression Network Analysis". Statistical Applications in Genetics and Molecular Biology 4(1), Article 17.

▲ロト ▲ 理 ト ▲ 王 ト ▲ 王 - の Q (~

Background	Theory	Practice	Foreground
Terminology			

continuous power-law probability distributions

 $P(X > x) \propto x^{-\alpha}, \quad x \ge x_{\min} \quad \exists x_{\min}, \alpha > 0$

Background	Theory	Practice	Foreground
Terminology			

continuous power-law probability distributions

 $P(X > x) \propto x^{-\alpha}, \quad x \ge x_{\min} \quad \exists x_{\min}, \alpha > 0$

► ... are **scaling** in the sense that*

$$P(X > sx \mid X > x) = s^{-\alpha} \quad \forall \ x \ge x_{\min}$$

- ロト - (日) - (1)

Background	Theory	Practice	Foreground
Terminology			

discrete power-law probability distributions

 $p_x \propto x^{-\alpha-1}, \quad x > x_{\min} \quad \exists x_{\min}, \alpha > 0$

Background	Theory	Practice	Foreground
Terminology			

discrete power-law probability distributions

 $p_x \propto x^{-\alpha-1}, \quad x > x_{\min} \quad \exists x_{\min}, \alpha > 0$

► ... produce scaling sequences*

$$x_1 \ge x_2 \ge \cdots \ge x_n, \quad k \propto x_k^{-\alpha}$$

< □ > < @ > < E > < E > E のQ@

Background	Theory	Practice	Foreground
DIAGNOS	TICS		
Size-ranl	< plots:		
	$k = C x_k^{-\alpha}$	$\Rightarrow \log k = \log(C) - \alpha \log x_k$	
library(m <- dis	poweRlaw) pl\$new()	(X);1 0 1000 0 10000	



< □ > < @ > < E > < E >

€ 900

Background	Theory	Practice	Foreground
DIAGNOS	TICS		
Size-ranl	k plots:		
	$k = C x_k^{-\alpha} \Rightarrow $	$\log k = \log(C) - \alpha \log x_k$	
library(noweRlaw)	1000	



< □ > < @ > < E > < E >

€

590

Background	Theory	Practice	Foreground
DIAGNOSTIC	5		
Size-frequenc	cy plots:		
$p_x = Cx$	$^{-\alpha-1} \Rightarrow \log$	$p_x = \log(C) - (\alpha + 1)\log(x)$)
<pre>library(powe m <- displ\$r m\$setXmin(7) m\$setPars(2, x <- dist_ra plot(x = son y = tak log = '</pre>	eRlaw) new() .3) and(m, 1e4) rt(unique(x) ole(sort(x)) 'xy")) , (10 50 200 1000 12 12 48 154 576 13 18 154 576 10 50 200 1000	<u>oasoo</u> 5000
2	<u> </u>	sort(unique(x))	

Background	Theory	Practice	Foreground
DIAGNOSTICS	5		
Size-frequenc	y plots:		
$p_x = Cx$	$^{-\alpha-1} \Rightarrow \log$	$p_x = \log(C) - (\alpha + 1)\log(x)$)
<pre>library(powe m <- displ\$r m\$setXmin(7) m\$setPars(2, x <- dist_ra plot(x = son y = tak log = '</pre>	eRlaw) new() .3) and(m, 1e4) rt(unique(x) ole(sort(x)) 'xy")) , (10 50 200 1000 sort(unique(x))	

Background	Theory	Practice	Foreground

Networks with scaling degree sequences are often called "scale-free", assumed to have specific origins:

- evolution via cumulative advantage (preferential attachment)
- emergence at phase transitions / criticality

Background	Theory	Practice	Foreground

Networks with scaling degree sequences are often called "scale-free", assumed to have specific origins:

• evolution via cumulative advantage (preferential attachment)

▲ロト ▲ 理 ト ▲ 王 ト ▲ 王 - の Q (~

• emergence at phase transitions / criticality

and asserted to have many topological properties:

- highly central high-degree nodes (hubs)
- self-similarity / hierarchicality
- invariance under degree-preserving rewiring
- robustness to failure but vulnerability to attack

Background	Theory	Practice	Foreground

Networks with scaling degree sequences are often called "scale-free", assumed to have specific origins:

- evolution via cumulative advantage (preferential attachment)
- emergence at phase transitions / criticality

and asserted to have many topological properties:

- highly central high-degree nodes (hubs)
- self-similarity / hierarchicality
- invariance under degree-preserving rewiring
- robustness to failure but vulnerability to attack

The apparent ubiquity of scaling degree distributions in empirical networks then has profound implications:

◆ロト ◆昼 ト ◆ 臣 ト ◆ 臣 ト ○ 臣 - のへぐ

- significance to natural order (*surprising*)
- universality to complex systems (*determinitive*)

Background	Theory	Practice	Foreground
Example			

To visually inspect whether approximate scale-free topology is satisfied, one plots $\log_{10}(p(k))$ versus $\log_{10}(k)$. A straight line is indicative of scale-free topology[.]

To measure how well a network satisfies a scale-free topology, we propose to use the square of the correlation between $\log(p(k))$ and $\log(k)$, i.e. the model fitting index R^2 of the linear model that regresses $\log(p(k))$ on $\log(k)$.



Zhang and Horvath, 2005

500

Background	Theory	Practice	Foreground
Context			

NETWORK BIOLOGY: UNDERSTANDING THE CELL'S FUNCTIONAL ORGANIZATION

Albert-László Barabási* & Zoltán N. Oltvai‡

Commentary

Scale-free networks in cell biology

Réka Albert

Department of Physics and Huck Institutes of the Life Sciences, Pennsylvania State University, University Park, PA 16802, USA (e-mail: raibert@phys.psu.edu)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Journal of Cell Science 118, 4947-4957 Published by The Company of Biologists 2005 doi:10.1242/jcs.02714

Background	Theory	Practice	Foreground
Context			

1. How widely do biologists observe scale-free networks?

< □ > < @ > < E > < E > E のQ@

- 2. What evidence supports these observations?
- 3. What implications follow from these observations?

Background	Theory	Practice	Foreground
CLAIMS ANI) EVIDENCE		



Claims made about distribution or generation of data

literature cited in Barabási & Oltvai, 2004 and in Albert, 2005

Background	Theory	Practice	Foreground
_			

CLAIMS AND EVIDENCE



literature cited in Barabási & Oltvai, 2004 and in Albert, 2005

Background	Theory	Practice	Foreground
IMPLICATIONS			
A Random network	B Scale-free network	C Hierarchical network	
Aa	Ba	Ca	
Ab	Bb	Cb	
		1,000 10 ⁻¹ 10 ⁻² 10 ⁻⁴ 10 ⁻⁶ 10 ⁻⁷ 10 ⁻⁸ 10 ⁻⁸ 10 ⁻¹⁰ 10	00 10,000

Barabási & Oltvai, 2004 🔊 < 🔿

Background	Theory	Practice	Foreground

- "Presence of scale-free behavior indicates a high degree of self-organization in the system and is known to be a characteristic of natural systems....The scale-free character of coexpressed gene networks means that these networks are extremely inhomogeneous and contain few genes that are very highly connected and a large number of genes with low connectivity." (Agrawal, 2002)
- "Since hubs are rare relative to other nodes, random elimination of nodes has minimal effect on network topology because statistically, a randomly eliminated node is likely to have low connectivity. ... We propose that organisms' general ability to compensate for individual mutations is largely a result of the scale-free properties of the gene expression network." (Featherstone & Broadie, 2002)
- "This result suggests the existence of a selective force in the overall design of genetic pathways to maintain a highly connected class of genes." (Stuart et al, 2003)

Background	Theory	Practice	Foreground

CONFUSION

- "Surprisingly, genes such as calcium/calmodulin-dependant kinase (CAM- KII) and the signaling GTPase Ras1, which are central to intracellular signaling and therefore would be most expected to 'tie' together various aspects of cell biology, are not hubs....It is also somewhat surprising that transcription factors do not head the list of hubs, since deletion of a transcription factor would be expected to alter the expression of all its downstream genes. (Featherstone & Broadie, 2002)
- "The distribution of interactions per protein decays faster than the power law predicted by a "rich-get-richer" model of scale-free networks[.] This rapid decay suggests that highly connected proteins may be suppressed in biological networks and supports a previous observation that connections between highly connected proteins are also suppressed." (Giot et al, 2003)
- "[W]ith our current metabolic information, the [average path length] of the *E. coli* network remains ≈ 8, much larger than that of a random graph. The metabolic world of *E. coli* is therefore not small with respect to biosynthesis/degradation pathways on the traditional metabolic map." (Masanori, 2004)

Background	Theory	Practice	Foreground

THEORY: POWER LAWS AND SCALE-FREE GRAPHS

Background	Theory	Practice	Foreground
THEORY	DOLUTER LANC	AND CONTE EDEE	OD A DITO

I'HEORY: POWER LAWS AND SCALE-FREE GRAPHS

- Mitzenmacher, M (2004). "A Brief History of Generative Models for Power Law and Lognormal Distributions". *Internet Mathematics* 1(2), 226–251.
- Willinger, W, Alderson, D, Doyle, J C, and Li, L (2004). "More "normal" than normal: scaling distributions and complex systems". Proceedings of the 2004 Winter Simulation Conference, 130–141.
- ► Fox Keller, E (2005). "Revisiting "scale-free" networks". *BioEssays* 27(10), 1060–1068.
- ► Li, L, Alderson, D, Doyle, J C, Willinger, W (2005). "Towards a Theory of Scale-Free Graphs: Definition, Properties, and Implications". *Internet Mathematics* 2(4), 431–523.
- Clauset, A, Shalizi, C R, Newman, M E J (2009). "Power-Law Distributions in Empirical Data" SIAM Review 51(4), 661–703.

Background	Theory	Practice	Foreground

Networks with scaling degree sequences are often called "scale-free", **assumed** to have specific origins:

- evolution via cumulative advantage (preferential attachment)
- emergence at phase transitions / criticality

and asserted to have many topological properties:

- highly central high-degree nodes (hubs)
- self-similarity / hierarchicality
- ► invariance under degree-preserving rewiring
- robustness to failure but vulnerability to attack

The **apparent** ubiquity of scaling degree distributions in empirical networks then has profound implications:

- significance to natural order (*surprising*)
- universality to complex systems (*determinitive*)

Background	Theory	Practice	Foreground
			(

GENERATIVE MODELS

Power-law behavior arises from many processes.

- Cumulative advantage
- Optimization \rightarrow
- Multiplicative processes
- Durations between events
- Double Pareto distributions



Figure 1. Rank-frequency distributions of Helen B. (with paranoid schizophrenia) in samples of (1) 50,000 words; (11) 30,000 words; (111) 20,000 words; (IV) 10,000 words; (V and VI) 5,000 words, and (VII and VIII) 2,000 words. (From Arch. Neurol. Psychiat. 49 (1943) 831)

Mitzenmacher, 2004 Sale

9		 0.0
INVARIANCE PRC	PERTIES	

Practice

Power-law behavior survives many transformations.

If X_1, \ldots, X_n

Background

• follow scaling distributions $X_k \sim P_k$

Theory

• with common scaling parameter $1 < \alpha < 2$

then so do

$$S_n = \sum_k X_k$$
 (aggregation; Central Limit Theorem)
 $M_n = \max(X_n)$ (maximizing choices)
 $W_n \sim \sum_k w_k P_k$ (weighted mixtures)

Willinger et al, 2004 $_{\mathcal{OQC}}$

Foreground

TOPOLOGICAL DIVERSITY

Power-law behavior admits many network topologies.



TOPOLOGICAL DIVERSITY

Power-law behavior admits many network topologies.



Background	Theory	Practice	Foreground
STATISTICA	L INFERENCE		

Power-law behavior cannot be reliably eyeballed.



□ Lauset et al, 2009 SQC

Background	Theory	Practice	Foreground
STATISTICAL	INFERENCE		

Power-law behavior cannot be reliably eyeballed.



Clauset et al, 2009 Sac

Background		Theory	Practice		Foreground
Impli	CATIONS				
	phase transition	power-law relationship	linearity on log-log plot		
	cumulative advantage	"scale-			
		freeness"	ŢŢ	<pre>}</pre>	(degree sequence of) a network
	central hubs self	-similarity inva	wiring ariance fragile]]	



Background	Theory	Practice	Foreground

PRACTICE: HYPOTHESIZING POWER LAWS

Background	Theory	Practice	Foreground

PRACTICE: HYPOTHESIZING POWER LAWS

- ▶ Vuong, Q H (1989). "Likelihood Ratio Tests for Model Selection and Non-Nested Hypotheses". *Econometrica* **57**(2), 307–333.
- Clauset, A, Young, M, (2007). "On the frequency of severe terrorist events" *Journal of Conflict Resolution* 51(1), 58–87.
- Clauset, A, Shalizi, C R, Newman, M E J (2009). "Power-Law Distributions in Empirical Data" SIAM Review 51(4), 661–703. http://tuvalu.santafe.edu/~aaronc/powerlaws/
- Gillespie, C S (2015). "Fitting Heavy Tailed Distributions: The poweRlaw Package". Journal of Statistical Software 64(2), 1–16. http://www.jstatsoft.org/v64/i02/

Background	Theory	Practice	Foreground
DECIDE			
KECIPE			

How to analyze (discrete) empirical data hypothesized as

$$y = Cx^{-\alpha}, \quad x \ge x_{\min}$$

- 1. Estimate lower bound x_{\min} and scaling parameter α
- 2. Assess goodness-of-fit (sampling + discernibility test)
- 3. Compare power-law to alternative hypotheses (LRT)

Background	Theory	Practice	Foreground
ESTIMATE TH	HE SCALING PA	RAMETER	

For a fixed x_{\min} , either

numerically maximize the likelihood function

$$-n\log\sum_{n=0}^{\infty}(n+x_{\min})^{-\alpha}-\alpha\sum_{i=1}^{n}\log x_{i}$$

or

use the approximation

$$\hat{\alpha} \simeq 1 + n \left[\sum_{i=1}^{n} \log \frac{x_i}{x_{\min} - \frac{1}{2}} \right]^{-1}$$

(best when $x_{\min} \gtrsim 6$)

Background	Theory	Practice	Foreground
ESTIMATE TE	IE LOWER BOU	JND	

1. For each x_{\min} , define

$$S(x) = \text{CDF}$$
 of data satisfying $x \ge x_{\min}$
 $P(x) = \text{CDF}$ of best-fit power-law model over $x \ge x_{\min}^*$

- 2. Pick a distribution distance measure *D*
- 3. Select x_{\min} that minimizes D(S(x), P(x))

Background	Theory	Practice	Foreground
-			
ESTIMATE T	HE LOWER BOL	IND	

1. For each x_{\min} , define

$$S(x) = \text{CDF}$$
 of data satisfying $x \ge x_{\min}$
 $P(x) = \text{CDF}$ of best-fit power-law model over $x \ge x_{\min}^*$

 \Box Clauset et al, 2007 \neg_{\circ}

- 2. Pick a distribution distance measure *D*
- 3. Select x_{\min} that minimizes D(S(x), P(x))
- * Decide how to handle P(X = x) for $x < x_{\min}$

Bac	kgroui	nd
	0	

Theory

Practice

EVALUATING AN x_{\min} estimate



Background	Theory	Practice	Foreground

CALCULATE THE GOODNESS-OF-FIT

- -2. Pick D
- -1. Estimate a power-law fit P(x) from the data $S(x)^*$
 - 0. Calculate D = D(P(x), S(x))
 - 1. Generate lots of artificial datasets $S'_i(x)$ from P(x)

 \Box Clauset et al, 2009 \mathcal{A}

- 2. Calculate $D_i = D(P(x), S'_i(x))$
- 3. Estimate $p \approx P(D_i > D)$

Background	Theory	Practice	Foreground
T			_



□ Lauset et al, 2009 9 <</p>

Background	Theory		Practice	Foreground
COMPARE AGAIN	ISI ALII	EKNA	IIVE HYPOI.	HESES
Given the best-fit j	ower-law	PDF a	and best-fit alterr	native PDF
	p(x)	and	q(x)	
calculate the likeli	hoods			
	11		11	

$$L_p = \prod_{i=1}^{n} p(x_i)$$
 and $L_q = \prod_{i=1}^{n} q(x_i)$

and consider the hypotheses

$$H_0: E[L_p/L_q] = 0, \quad H_p: E[L_p/L_q] > 0, \quad H_q: E[L_p/L_q] < 0$$

Vuong, 1989; Clauset et al, 2009 on the second seco

Background	Theory		Practice	Foreground
COMPARE AG. Given the best	AINST ALTI -fit power-law	ERNA PDF a	TIVE HYPOT and best-fit alterr	HESES native PDF
	p(x)	and	q(x)	
calculate the li	kelihoods			

$$L_p = \prod_{i=1}^n p(x_i)$$
 and $L_q = \prod_{i=1}^n q(x_i)$

and consider the hypotheses

$$H_0: E[L_p/L_q] = 0, \quad H_p: E[L_p/L_q] > 0, \quad H_q: E[L_p/L_q] < 0$$

Then the log-likelihood ratio

$$\mathcal{R} = \sum_{i=1}^{n} \left[\log p(x_i) - \log q(x_i) \right]$$

is asymptotically normal with $\sigma_{\mathcal{R}}^2 \xrightarrow{n \to \infty} 0$ Vuong, 1989; Clauset et al, 2009

EVALUATION OF RECIPE

Compare 24 "power-law" datasets for power-law, log-normal, and exponential behavior...

- (a) The frequency of occurrence of unique words in the novel Moby Dick by Herman Melville [43].
- (b) The degrees (i.e., numbers of distinct interaction partners) of proteins in the partially known protein-interaction network of the yeast Saccharomyces cerevisiae [28].
- (c) The degrees of metabolites in the metabolic network of the bacterium Escherichia coli [26].
- (d) The degrees of nodes in the partially known network representation of the Internet at the level of autonomous systems for May 2006 [25]. (An autonomous system is a group of IP addresses on the Internet among which routing is handled internally or "autonomously," rather than using the Internet's large-scale border gateway protocol routing mechanism.)
- (e) The number of calls received by customers of AT&T's long distance telephone service in the United States during a single day [1, 5].
- (f) The intensity of wars from 1816–1980 measured as the number of battle deaths per 10 000 of the combined populations of the warring nations [53, 49].

Clauset et al, 2009 Sac

STATISTICAL INFERENCE

Power-law behavior is not ubiquitous.



Statistical support for power-law distributions

- "good": power-law is plausible; others are not
- "moderate": power-law is plausible; others are also
- "with cut-off": power-law with exponential cut-off is plausible
- "none": power-law is not plausible

 \square , \square ,

STATISTICAL INFERENCE

Power-law behavior is not ubiquitous.

Evidence for power-law versus log-normal distributions



Background	Theory	Practice	Foreground

TUTORIAL IN poweRlaw PACKAGE FOR R

CICATS Study Group 20 March (Monday) @ noon

