# More Mechanisms for Generating Power-Law Size Distributions II

## Principles of Complex Systems
## CSYS/MATH 300, Fall, 2011

### Prof. Peter Dodds

Department of Mathematics & Statistics | Center for Complex Systems |
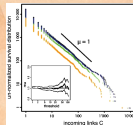Vermont Advanced Computing Center | University of Vermont

# Outline

## Growth Mechanisms
### Random Copying
### Words, Cities, and the Web

## Optimization
### Minimal Cost
### Mandelbrot vs. Simon
### Assumptions
### Model
### Analysis
### Extra
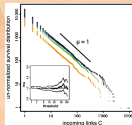### And the winner is...?

## References

# Aggregation

- Random walks represent additive aggregation
- Mechanism: Random addition and subtraction
- Compare across realizations, no competition.
- Next: Random Additive/Copying Processes involving Competition.
- Widespread: Words, Cities, the Web, Wealth, Productivity (Lotka), Popularity (Books, People, ...)
- Competing mechanisms (trickiness)

# Work of Yore

- ▶ 1924: G. Udny Yule[23]:
  # Species per Genus
- ▶ 1926: Lotka[10]:
  # Scientific papers per author (Lotka's law)
- ▶ 1953: Mandelbrot[12]:
  Optimality argument for Zipf's law; focus on
  language.
- ▶ 1955: Herbert Simon[19, 25]:
  Zipf's law for word frequency, city size, income,
  publications, and species per genus.
- ▶ 1965/1976: Derek de Solla Price[17, 18]:
  Network of Scientific Citations.
- ▶ 1999: Barabasi and Albert[1]:
  The World Wide Web, networks-at-large.



The
UNIVERSITY
of VERMONT

# Examples

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web

Optimization
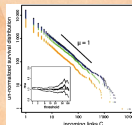Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
And the winner is...?
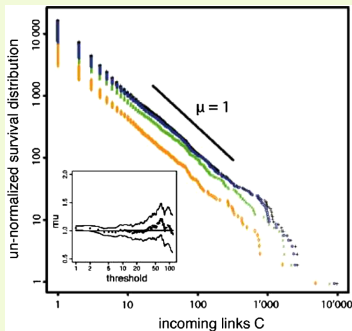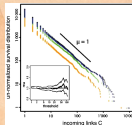
References

## Recent evidence for Zipf's law...



FIG. 1 (color online). (Color Online) Log-log plot of the number of packages in four Debian Linux Distributions with more than $C$ in-directed links. The four Debian Linux Distributions are Woody (19.07.2002) (orange diamonds), Sarge (06.06.2005) (green crosses), Etch (15.08.2007) (blue circles), Lenny (15.12.2007) (black+'s). The inset shows the maximum likelihood estimate (MLE) of the exponent $\mu$ together with two boundaries defining its 95% confidence interval (approximately given by $1 \pm 2/\sqrt{n}$, where $n$ is the number of data points using in the MLE), as a function of the lower threshold. The MLE has been modified from the standard Hill estimator to take into account the discreteness of $C$.

Maillart et al., PRL, 2008:
"Empirical Tests of Zipf's Law Mechanism in Open Source Linux Distribution"[11]

UNIVERSITY
of VERMONT

# Essential Extract of a Growth Model

Random Competitive Replication (RCR):

1. Start with 1 element of a particular flavor at $t = 1$
2. At time $t = 2, 3, 4, \ldots$, add a new element in one of two ways:
   - With probability $\rho$, create a new element with a new flavor
     ➤ Mutation/Innovation

   - With probability $1 - \rho$, randomly choose from all existing elements, and make a copy.
     ➤ Replication/Imitation

   - Elements of the same flavor form a group



UNIVERSITY
of VERMONT

# Random Competitive Replication

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web

Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
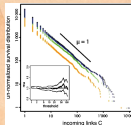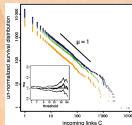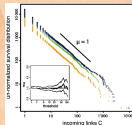And the winner is...?

References

### Example: Words in a text

▶ Consider words as they appear sequentially.

▶ With probability $\rho$, the next word has not previously appeared
  ➤ Mutation/Innovation

▶ With probability $1 - \rho$, randomly choose one word from all words that have come before, and reuse this word
  ➤ Replication/Imitation



UNIVERSITY
of VERMONT

# Random Competitive Replication

- ▶ Competition for replication between elements is random
- ▶ Competition for growth between groups is not random
- ▶ Selection on groups is biased by size
- ▶ Rich-gets-richer story
- ▶ Random selection is easy
- ▶ No great knowledge of system needed

# Random Competitive Replication

- Steady growth of system: +1 element per unit time.
- Steady growth of distinct flavors at rate $\rho$
- We can incorporate
  1. Element elimination
  2. Elements moving between groups
  3. Variable innovation rate $\rho$
  4. Different selection based on group size
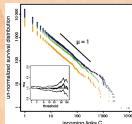     (But mechanism for selection is not as simple...)



UNIVERSITY
of VERMONT

# Random Competitive Replication

Definitions:

- ▶ $k_i$ = size of a group $i$
- ▶ $N_k(t)$ = # groups containing $k$ elements at time $t$.

Basic question: How does $N_k(t)$ evolve with time?

First: $\displaystyle\sum_k kN_k(t) = t = $ number of elements at time $t$



THE
UNIVERSITY
of VERMONT

# Random Competitive Replication
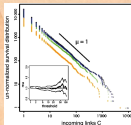
$P_k(t)$ = Probability of choosing an element that belongs to a group of size $k$:

- $N_k(t)$ size $k$ groups
- $\Rightarrow k N_k(t)$ elements in size $k$ groups
- $t$ elements overall

$$P_k(t) = \frac{k N_k(t)}{t}$$

# Random Competitive Replication

$N_k(t)$, the number of groups with $k$ elements, changes at time $t$ if

1. An element belonging to a group with $k$ elements is replicated
   $N_k(t+1) = N_k(t) - 1$
   Happens with probability $(1 - \rho)kN_k(t)/t$

2. An element belonging to a group with $k - 1$ elements is replicated
   $N_k(t+1) = N_k(t) + 1$
   Happens with probability $(1 - \rho)(k-1)N_{k-1}(t)/t$



UNIVERSITY
of VERMONT

# Random Competitive Replication

Special case for $N_1(t)$:
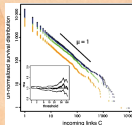
1. The new element is a new flavor:
   $N_1(t+1) = N_1(t) + 1$
   Happens with probability $\rho$

2. A unique element is replicated.
   $N_1(t+1) = N_1(t) - 1$
   Happens with probability $(1-\rho)N_1/t$

# Random Competitive Replication

Put everything together:

For $k > 1$:

$$\langle N_k(t+1) - N_k(t) \rangle = (1-\rho)\left((k-1)\frac{N_{k-1}(t)}{t} - k\frac{N_k(t)}{t}\right)$$

For $k = 1$:

$$\langle N_1(t+1) - N_1(t) \rangle = \rho - (1-\rho)1 \cdot \frac{N_1(t)}{t}$$

# Random Competitive Replication

Assume distribution stabilizes: $N_k(t) = n_k t$

(Reasonable for $t$ large)

- Drop expectations
- Numbers of elements now fractional
- Okay over large time scales
- $n_k/\rho$ = the fraction of groups that have size $k$.

# Random Competitive Replication

Stochastic difference equation:

$$\langle N_k(t+1) - N_k(t) \rangle = (1-\rho)\left((k-1)\frac{N_{k-1}(t)}{t} - k\frac{N_k(t)}{t}\right)$$

becomes

$$n_k(t+1) - n_k t = (1-\rho)\left((k-1)\frac{n_{k-1}t}{t} - k\frac{n_k t}{t}\right)$$

$$n_k(\cancel{t}+1-\cancel{t}) = (1-\rho)\left((k-1)\frac{n_{k-1}\cancel{t}}{\cancel{t}} - k\frac{n_k \cancel{t}}{\cancel{t}}\right)$$

$$\Rightarrow n_k = (1-\rho)\left((k-1)n_{k-1} - kn_k\right)$$

$$\Rightarrow n_k\left(1 + (1-\rho)k\right) = (1-\rho)(k-1)n_{k-1}$$

The
UNIVERSITY
of VERMONT

# Random Competitive Replication

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web

Optimization
Minimal Cost
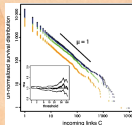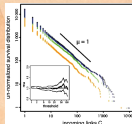Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
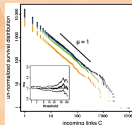And the winner is...?

References

We have a simple recursion:

$$\frac{n_k}{n_{k-1}} = \frac{(k-1)(1-\rho)}{1+(1-\rho)k}$$

▶ Interested in $k$ large (the tail of the distribution)
▶ Can be solved exactly.
  Insert question from assignment 4 (⊞)
▶ To get at tail: Expand as a series of powers of $1/k$
  Insert question from assignment 4 (⊞)

# Random Competitive Replication

► We (okay, you) find

$$\frac{n_k}{n_{k-1}} \simeq \left(1 - \frac{1}{k}\right)^{\frac{(2-\rho)}{(1-\rho)}}$$

►

$$\frac{n_k}{n_{k-1}} \simeq \left(\frac{k-1}{k}\right)^{\frac{(2-\rho)}{(1-\rho)}}$$

►

$$n_k \propto k^{-\frac{(2-\rho)}{(1-\rho)}} = k^{-\gamma}$$

$$\boxed{\gamma = \frac{(2-\rho)}{(1-\rho)} = 1 + \frac{1}{(1-\rho)}}$$

# Random Competitive Replication

$$\gamma = \frac{(2-\rho)}{(1-\rho)} = 1 + \frac{1}{(1-\rho)}$$

▶ Micro to macros story with $\gamma$ and $\rho$ measurable.

▶ Observe $2 < \gamma < \infty$ as $\rho$ varies.

▶ For $\rho \simeq 0$ (low innovation rate):

$$\gamma \simeq 2$$

▶ Recalls Zipf's law: $s_r \sim r^{-\alpha}$
  ($s_r$ = size of the $r$th largest element)

▶ We found $\alpha = 1/(\gamma - 1)$

▶ $\gamma = 2$ corresponds to $\alpha = 1$

# Random Competitive Replication

- We (roughly) see Zipfian exponent [25] of $\alpha = 1$ for many real systems: city sizes, word distributions, ...
- Corresponds to $\rho \to 0$ (Krugman doesn't like it) [9]
- But still other mechanisms are possible...
- Must look at the details to see if mechanism makes sense... more later.

# Random Competitive Replication

We had one other equation:

►
$$\langle N_1(t+1) - N_1(t) \rangle = \rho - (1-\rho)1 \cdot \frac{N_1(t)}{t}$$

► As before, set $N_1(t) = n_1 t$ and drop expectations

►
$$n_1(t+1) - n_1 t = \rho - (1-\rho)1 \cdot \frac{n_1 t}{t}$$

►
$$n_1 = \rho - (1-\rho)n_1$$

► Rearrange:
$$n_1 + (1-\rho)n_1 = \rho$$

►
$$\boxed{n_1 = \frac{\rho}{2-\rho}}$$

UNIVERSITY of VERMONT

# Random Competitive Replication

$$\text{So...} \qquad N_1(t) = n_1 t = \frac{\rho t}{2 - \rho}$$

▶ Recall number of distinct elements = $\rho t$.
▶ Fraction of distinct elements that are unique (belong to groups of size 1):

$$\frac{N_1(t)}{\rho t} = \frac{1}{2 - \rho}$$

(also = fraction of groups of size 1)
▶ For $\rho$ small, fraction of unique elements $\sim 1/2$
▶ Roughly observed for real distributions
▶ $\rho$ increases, fraction increases
▶ Can show fraction of groups with two elements $\sim 1/6$
▶ Model does well at both ends of the distribution



UNIVERSITY
of VERMONT

From Simon[19]:

Estimate $\rho_{est}$ = # unique words/# all words

For Joyce's Ulysses: $\rho_{est} \simeq 0.115$

| $N_1$ (real) | $N_1$ (est) | $N_2$ (real) | $N_2$ (est) |
|---|---|---|---|
| 16,432 | 15,850 | 4,776 | 4,870 |



The
UNIVERSITY
of VERMONT

# Evolution of catch phrases

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web
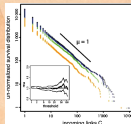Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
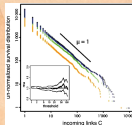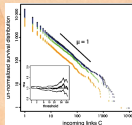And the winner is...?

References

- Yule's paper (1924)[23]:
  "A mathematical theory of evolution, based on the conclusions of Dr J. C. Willis, F.R.S."
- Simon's paper (1955)[19]:
  "On a class of skew distribution functions" (snore)

## From Simon's introduction:

It is the purpose of this paper to analyse a class of distribution functions that appear in a wide range of empirical data—particularly data describing sociological, biological and economic phenomena.

Its appearance is so frequent, and the phenomena so diverse, that one is led to conjecture that if these phenomena have any property in common it can only be a similarity in the structure of the underlying probability mechanisms.



THE
UNIVERSITY
of VERMONT

# Evolution of catch phrases

More on Herbert Simon (1916–2001):

► Political scientist

► Involved in Cognitive Psychology, Computer Science, Public Administration, Economics, Management, Sociology

► Coined 'bounded rationality' and 'satisficing'

► Nearly 1000 publications

► An early leader in Artificial Intelligence, Information Processing, Decision-Making, Problem-Solving, Attention Economics, Organization Theory, Complex Systems, And Computer Simulation Of Scientific Discovery.

► Nobel Laureate in Economics

# Evolution of catch phrases

## Derek de Solla Price:

▶ First to study network evolution with these kinds of models.

▶ Citation network of scientific papers

▶ Price's term: Cumulative Advantage

▶ Idea: papers receive new citations with probability proportional to their existing # of citations

▶ Directed network

▶ Two (surmountable) problems:
  1. New papers have no citations
  2. Selection mechanism is more complicated

# Evolution of catch phrases

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web
Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
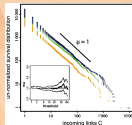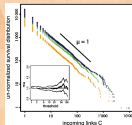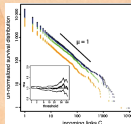And the winner is...?

References

Robert K. Merton: the Matthew Effect (⊞)

▶ Studied careers of scientists and found credit flowed
   disproportionately to the already famous

   From the Gospel of Matthew:
   "For to every one that hath shall be given...
   (Wait! There's more....)
   but from him that hath not, that also which he
   seemeth to have shall be taken away.
   And cast the worthless servant into the outer
   darkness; there men will weep and gnash their teeth."

▶ (Hath = suggested unit of purchasing power.)

▶ Matilda effect: (⊞) women's scientific achievements
   are often overlooked
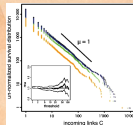


The
UNIVERSITY
of VERMONT

# Evolution of catch phrases

Merton was a catchphrase machine:

1. Self-fulfilling prophecy
2. Role model
3. Unintended (or unanticipated) consequences
4. Focused interview $\rightarrow$ focus group

And just to be clear...

Merton's son, Robert C. Merton, won the Nobel Prize for Economics in 1997.

# Evolution of catch phrases

- ▶ Barabasi and Albert [1]—thinking about the Web
- ▶ Independent reinvention of a version of Simon and Price's theory for networks
- ▶ Another term: "Preferential Attachment"
- ▶ Considered undirected networks (not realistic but avoids 0 citation problem)
- ▶ Still have selection problem based on size (non-random)
- ▶ Solution: Randomly connect to a node (easy) …
- ▶ … and then randomly connect to the node's friends (also easy)
- ▶ Scale-free networks = food on the table for physicists



UNIVERSITY of VERMONT

# Benoît Mandelbrot (⊞)

Nassim Taleb's tribute:

> Benoit Mandelbrot, 1924-2010
>
> *A Greek among Romans*

- ▶ Mandelbrot = father of fractals
- ▶ Mandelbrot = almond bread
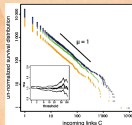- ▶ Bonus Mandelbrot set action: here (⊞).

# Another approach

## Benoît Mandelbrot

- ► Derived Zipf's law through optimization [12]
- ► Idea: Language is efficient
- ► Communicate as much information as possible for as little cost
- ► Need measures of information ($H$) and average cost ($C$)...
- ► Language evolves to maximize $H/C$, the amount of information per average cost.
- ► Equivalently: minimize $C/H$.
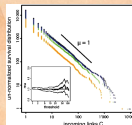- ► Recurring theme: what role does optimization play in complex systems?



UNIVERSITY
of VERMONT

# Not everyone is happy...

Mandelbrot vs. Simon:

▶ Mandelbrot (1953): "An Informational Theory of the Statistical Structure of Languages"[12]

▶ Simon (1955): "On a class of skew distribution functions"[19]

▶ Mandelbrot (1959): "A note on a class of skew distribution function: analysis and critique of a paper by H.A. Simon"[13]

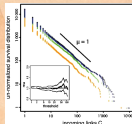▶ Simon (1960): "Some further notes on a class of skew distribution functions"[20]



THE UNIVERSITY OF VERMONT

# Not everyone is happy... (cont.)

Mandelbrot vs. Simon:

► Mandelbrot (1961): "Final note on a class of skew distribution functions: analysis and critique of a model due to H.A. Simon" [15]

► Simon (1961): "Reply to 'final note' by Benoit Mandelbrot" [22]

► Mandelbrot (1961): "Post scriptum to 'final note'" [15]

► Simon (1961): "Reply to Dr. Mandelbrot's post scriptum" [21]

# Not everyone is happy... (cont.)

### Mandelbrot:

"We shall restate in detail our 1959 objections to Simon's 1955 model for the Pareto-Yule-Zipf distribution. Our objections are valid quite irrespectively of the sign of p-1, so that most of Simon's (1960) reply was irrelevant." [14]

### Simon:

"Dr. Mandelbrot has proposed a new set of objections to my 1955 models of the Yule distribution. Like his earlier objections, these are invalid." [22]

# Zipfarama via Optimization

More Power-Law
Mechanisms II

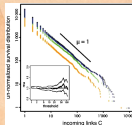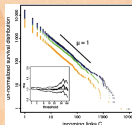Growth
Mechanisms
Random Copying
Words, Cities, and the Web

Optimization
Minimal Cost
Mandelbrot vs. Simon
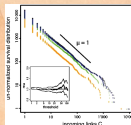Assumptions
Model
Analysis
Extra
And the winner is...?

References

## Mandelbrot's Assumptions:

- Language contains $n$ words: $w_1, w_2, \ldots, w_n$.
- $i$th word appears with probability $p_i$
- Words appear randomly according to this distribution (obviously not true...)
- Words = composition of letters is important
- Alphabet contains $m$ letters
- Words are ordered by length (shortest first)

# Zipfarama via Optimization

### Word Cost
- ▶ Length of word (plus a space)
- ▶ Word length was irrelevant for Simon's method

### Objection
- ▶ Real words don't use all letter sequences

### Objections to Objection
- ▶ Maybe real words roughly follow this pattern (?)
- ▶ Words can be encoded this way
- ▶ Na na na-na naaaaa...



UNIVERSITY
of VERMONT

# Zipfarama via Optimization

### Binary alphabet plus a space symbol

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| word | 1 | 10 | 11 | 100 | 101 | 110 | 111 | 1000 |
| length | 1 | 2 | 2 | 3 | 3 | 3 | 3 | 4 |
| $1 + \ln_2 i$ | 1 | 2 | 2.58 | 3 | 3.32 | 3.58 | 3.81 | 4 |

- ▶ Word length of $2^k$th word: $= k + 1 = 1 + \log_2 2^k$
- ▶ Word length of $i$th word $\simeq 1 + \log_2 i$
- ▶ For an alphabet with $m$ letters,
  word length of $i$th word $\simeq 1 + \log_m i$.

# Zipfarama via Optimization

## Total Cost $C$

- Cost of the $i$th word: $C_i \simeq 1 + \log_m i$
- Cost of the $i$th word plus space: $C_i \simeq 1 + \log_m(i+1)$
- Subtract fixed cost: $C_i' = C_i - 1 \simeq \log_m(i+1)$
- Simplify base of logarithm:

$$C_i' \simeq \log_m(i+1) = \frac{\log_e(i+1)}{\log_e m} \propto \ln(i+1)$$

- Total Cost:

$$C \sim \sum_{i=1}^{n} p_i C_i' \propto \sum_{i=1}^{n} p_i \ln(i+1)$$

The UNIVERSITY of VERMONT

# Zipfarama via Optimization

## Information Measure

▶ Use Shannon's Entropy (or Uncertainty):

$$H = -\sum_{i=1}^{n} p_i \log_2 p_i$$

▶ (allegedly) von Neumann suggested 'entropy'...

▶ Proportional to average number of bits needed to encode each 'word' based on frequency of occurrence

▶ $-\log_2 p_i = \log_2 1/p_i$ = minimum number of bits needed to distinguish event $i$ from all others

▶ If $p_i = 1/2$, need only 1 bit ($log_2 1/p_i = 1$)

▶ If $p_i = 1/64$, need 6 bits ($log_2 1/p_i = 6$)



UNIVERSITY of VERMONT

# Zipfarama via Optimization

## Information Measure

▶ Use a slightly simpler form:

$$H = - \sum_{i=1}^{n} p_i \log_e p_i / \log_e 2 = -g \sum_{i=1}^{n} p_i \ln p_i$$

where $g = 1/\ln 2$



UNIVERSITY
of VERMONT

# Zipfarama via Optimization

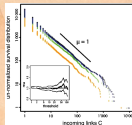More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web
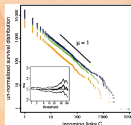
Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
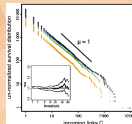Model
Analysis
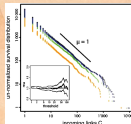Extra
And the winner is...?

References

▶ Minimize

$$F(p_1, p_2, \ldots, p_n) = C/H$$

subject to constraint

$$\sum_{i=1}^{n} p_i = 1$$

▶ Tension:
(1) Shorter words are cheaper
(2) Longer words are more informative (rarer)

# Zipfarama via Optimization

Time for Lagrange Multipliers:

▶ Minimize

$$\Psi(p_1, p_2, \ldots, p_n) =$$

$$F(p_1, p_2, \ldots, p_n) + \lambda G(p_1, p_2, \ldots, p_n)$$

where

$$F(p_1, p_2, \ldots, p_n) = \frac{C}{H} = \frac{\sum_{i=1}^{n} p_i \ln(i+1)}{-g \sum_{i=1}^{n} p_i \ln p_i}$$

and the constraint function is

$$G(p_1, p_2, \ldots, p_n) = \sum_{i=1}^{n} p_i - 1 = 0$$

Insert question from assignment 5 (⊞)



UNIVERSITY of VERMONT

# Zipfarama via Optimization

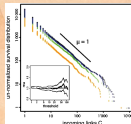Some mild suffering leads to:

► 
$$p_j = e^{-1-\lambda H^2/gC}(j+1)^{-H/gC} \propto (j+1)^{-H/gC}$$

► A power law appears [applause]: $\boxed{\alpha = H/gC}$

► Next: sneakily deduce $\lambda$ in terms of $g$, $C$, and $H$.

► Find
$$p_j = (j+1)^{-H/gC}$$

# Zipfarama via Optimization

### Finding the exponent

▶ Now use the normalization constraint:

$$1 = \sum_{j=1}^{n} p_j = \sum_{j=1}^{n} (j+1)^{-H/gC} = \sum_{j=1}^{n} (j+1)^{-\alpha}$$

▶ As $n \to \infty$, we end up with $\zeta(H/gC) = 2$
where $\zeta$ is the Riemann Zeta Function

▶ Gives $\alpha \simeq 1.73$ ($> 1$, too high)

▶ If cost function changes ($j + 1 \to j + a$) then
exponent is tunable

▶ Increase $a$, decrease $\alpha$

# Zipfarama via Optimization

All told:

- ▶ Reasonable approach: Optimization is at work in evolutionary processes
- ▶ But optimization can involve many incommensurate elements: monetary cost, robustness, happiness,...
- ▶ Mandelbrot's argument is not super convincing
- ▶ Exponent depends too much on a loose definition of cost



The
UNIVERSITY
of VERMONT

# More

## Reconciling Mandelbrot and Simon

► Mixture of local optimization and randomness

► Numerous efforts...

1. Carlson and Doyle, 1999:
   Highly Optimized Tolerance
   (HOT)—Evolved/Engineered Robustness [4, 5]

2. Ferrer i Cancho and Solé, 2002:
   Zipf's Principle of Least Effort [8]

3. D'Souza et al., 2007:
   Scale-free networks [6]



THE
UNIVERSITY
of VERMONT

# More

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web
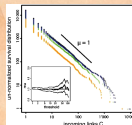
Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
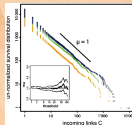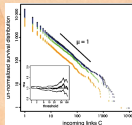And the winner is...?

References

Other mechanisms:

▶ Much argument about whether or not monkeys typing could produce Zipf's law... (Miller, 1957) [16]

▶ Miller gets to slap Zipf a little in an introduction to a 1965 reprint of Zipf's "Psycho-biology of Language" [24]

▶ Still fighting: "Random Texts Do Not Exhibit the Real Zipf's Law-Like Rank Distribution" [7] by Ferrer-i-Cancho and Elvevåg, 2010.



The
UNIVERSITY
of VERMONT

# Others are also not happy

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web

Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
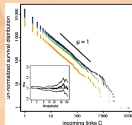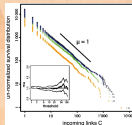And the winner is...?

References

## Krugman and Simon

▶ "The Self-Organizing Economy" (Paul Krugman, 1995) [9]

▶ Krugman touts Zipf's law for cities, Simon's model

▶ "Déjà vu, Mr. Krugman" (Berry, 1999)

▶ Substantial work done by Urban Geographers

# Who needs a hug?

## From Berry [2]

▶ Déjà vu, Mr. Krugman. Been there, done that. The Simon-Ijiri model was introduced to geographers in 1958 as an explanation of city size distributions, the first of many such contributions dealing with the steady states of random growth processes, ...

▶ But then, I suppose, even if Krugman had known about these studies, they would have been discounted because they were not written by professional economists or published in one of the top five journals in economics!

The
UNIVERSITY
of VERMONT

# Who needs a hug?

From Berry [2]

- ... [Krugman] needs to exercise some humility, for his world view is circumscribed by folkways that militate against recognition and acknowledgment of scholarship beyond his disciplinary frontier.
- Urban geographers, thank heavens, are not so afflicted.

# So who's right?

**Empirical Tests of Zipf's Law Mechanism in Open Source Linux Distribution**

T. Maillart,[1] D. Sornette,[1] S. Spaeth,[2] and G. von Krogh[2]

[1]Chair of Entrepreneurial Risks, Department of Management, Technology and Economics, ETH Zurich, CH-8001 Zurich, Switzerland
[2]Chair of Strategic Management and Innovation, Department of Management, Technology and Economics,
ETH Zurich, CH-8001 Zurich, Switzerland

Zipf's power law is a ubiquitous empirical regularity found in many systems, thought to result from proportional growth. Here, we establish empirically the usually assumed ingredients of stochastic growth models that have been previously conjectured to be at the origin of Zipf's law. We use exceptionally detailed data on the evolution of open source software projects in Linux distributions, which offer a remarkable example of a growing complex self-organizing adaptive system, exhibiting Zipf's law over four full decades.

# So who's right?

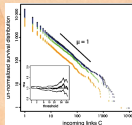More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web

Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
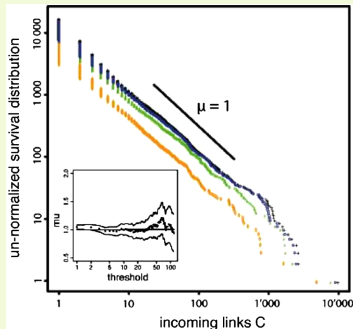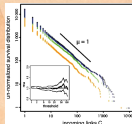Model
Analysis
Extra
And the winner is...?

References

FIG. 1 (color online). (Color Online) Log-log plot of the number of packages in four Debian Linux Distributions with more than $C$ in-directed links. The four Debian Linux Distributions are Woody (19.07.2002) (orange diamonds), Sarge (06.06.2005) (green crosses), Etch (15.08.2007) (blue circles), Lenny (15.12.2007) (black + 's). The inset shows the maximum likelihood estimate (MLE) of the exponent $\mu$ together with two boundaries defining its 95% confidence interval (approximately given by $1 \pm 2/\sqrt{n}$, where $n$ is the number of data points using in the MLE), as a function of the lower threshold. The MLE has been modified from the standard Hill estimator to take into account the discreteness of $C$.

Maillart et al., PRL, 2008:
"Empirical Tests of Zipf's Law Mechanism in Open Source Linux Distribution"[11]
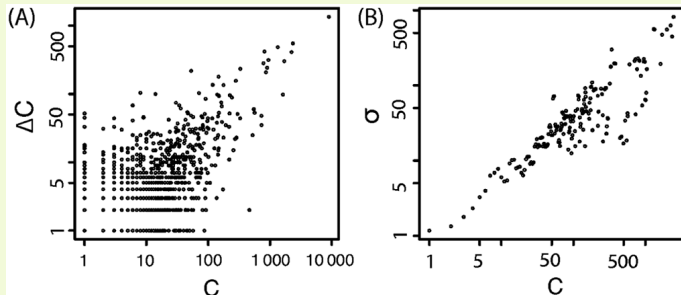
# So who's right?

FIG. 2. Left panel: Plots of $\Delta C$ versus $C$ from the Etch release (15.08.2007) to the latest Lenny version (05.05.2008) in double logarithmic scale. Only positive values are displayed. The linear regression $\Delta C = R \times C + C_0$ is significant at the 95% confidence level, with a small value $C_0 = 0.3$ at the origin and $R = 0.09$. Right panel: same as left panel for the standard deviation of $\Delta C$.

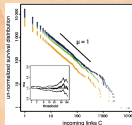► Rough, approximately linear relationship between $C$

UNIVERSITY
of VERMONT

# So who's right?

Bornholdt and Ebel (PRE), 2001:
"World Wide Web scaling exponent from Simon's 1955
model" [3].

- ► Show Simon's model fares well.
- ► Recall $\rho$ = probability new flavor appears.
- ► Alta Vista (⊞) crawls in approximately 6 month period
  in 1999 give $\rho \simeq 0.10$
- ► Leads to $\gamma = 1 + \frac{1}{1-\rho} \simeq 2.1$ for in-link distribution.
- ► Cite direct measurement of $\gamma$ at the time: $2.1 \pm 0.1$
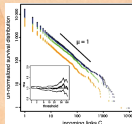  and 2.09 in two studies.

# So who's right?

Nutshell:

- Simonish random 'rich-get-richer' models agree in detail with empirical observations.
- Power-lawfulness: Mandelbrot's optimality is still apparent.
- Optimality arises for free in Random Competitive Replication models.



THE UNIVERSITY OF VERMONT

# References I

[1] A.-L. Barabási and R. Albert.
Emergence of scaling in random networks.
Science, 286:509–511, 1999. pdf (⊞)

[2] B. J. L. Berry.
Déjà vu, Mr. Krugman.
Urban Geography, 20:1–2, 1999. pdf (⊞)

[3] S. Bornholdt and H. Ebel.
World Wide Web scaling exponent from Simon's
1955 model.
Phys. Rev. E, 64:035104(R), 2001. pdf (⊞)

[4] J. M. Carlson and J. Doyle.
Highly optimized tolerance: A mechanism for power
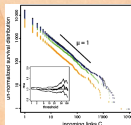laws in design systems.
Phys. Rev. E, 60(2):1412–1427, 1999. pdf (⊞)

# References II

[5] J. M. Carlson and J. Doyle.
Complexity and robustness.
Proc. Natl. Acad. Sci., 99:2538–2545, 2002. pdf (⊞)

[6] R. M. D'Souza, C. Borgs, J. T. Chayes, N. Berger,
and R. D. Kleinberg.
Emergence of tempered preferential attachment
from optimization.
Proc. Natl. Acad. Sci., 104:6112–6117, 2007. pdf (⊞)

[7] R. Ferrer-i Cancho and B. Elvevåg.
Random texts do not exhibit the real Zipf's law-like
rank distribution.
PLoS ONE, 5:e9411, 03 2010.

[8] R. Ferrer i Cancho and R. V. Solé.
Zipf's law and random texts.
Advances in Complex Systems, 5(1):1–6, 2002.

# References III

[9]   P. Krugman.
      The self-organizing economy.
      Blackwell Publishers, Cambridge, Massachusetts,
      1995.

[10]  A. J. Lotka.
      The frequency distribution of scientific productivity.
      Journal of the Washington Academy of Science,
      16:317–323, 1926.

[11]  T. Maillart, D. Sornette, S. Spaeth, and G. von
      Krogh.
      Empirical tests of Zipf's law mechanism in open
      source Linux distribution.
      Phys. Rev. Lett., 101(21):218701, 2008. pdf (⊞)
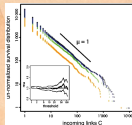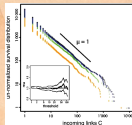
# References IV

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web

Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
And the winner is...?

References

[12] B. B. Mandelbrot.
An informational theory of the statistical structure of
languages.
In W. Jackson, editor, Communication Theory, pages
486–502. Butterworth, Woburn, MA, 1953. pdf (⊞)

[13] B. B. Mandelbrot.
A note on a class of skew distribution function.
analysis and critique of a paper by H. A. Simon.
Information and Control, 2:90–99, 1959.

[14] B. B. Mandelbrot.
Final note on a class of skew distribution functions:
analysis and critique of a model due to H. A. Simon.
Information and Control, 4:198–216, 1961.
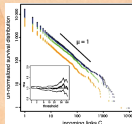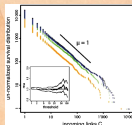
# References V

More Power-Law
Mechanics II

Growth
Mechanics
Random Copying
Words, Cities, and the Web

Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
And the winner is...?

References

[15] B. B. Mandelbrot.
Post scriptum to 'final note'.
Information and Control, 4:300–304, 1961.

[16] G. A. Miller.
Some effects of intermittent silence.
American Journal of Psychology, 70:311–314, 1957.
pdf (⊞)

[17] D. J. d. S. Price.
Networks of scientific papers.
Science, 149:510–515, 1965. pdf (⊞)

[18] D. J. d. S. Price.
A general theory of bibliometric and other cumulative
advantage processes.
J. Amer. Soc. Inform. Sci., 27:292–306, 1976.

# References VI

[19] H. A. Simon.
On a class of skew distribution functions.
Biometrika, 42:425–440, 1955. pdf (⊞)

[20] H. A. Simon.
Some further notes on a class of skew distribution
functions.
Information and Control, 3:80–88, 1960.

[21] H. A. Simon.
Reply to Dr. Mandelbrot's post scriptum.
Information and Control, 4:305–308, 1961.

[22] H. A. Simon.
Reply to 'final note' by Benoît Mandelbrot.
Information and Control, 4:217–223, 1961.
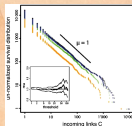
# References VII

More Power-Law
Mechanisms II

Growth
Mechanisms
Random Copying
Words, Cities, and the Web
Optimization
Minimal Cost
Mandelbrot vs. Simon
Assumptions
Model
Analysis
Extra
And the winner is...?

References

[23] G. U. Yule.
A mathematical theory of evolution, based on the
conclusions of Dr J. C. Willis, F.R.S.
Phil. Trans. B, 213:21–, 1924.

[24] G. K. Zipf.
The Psychobiology of Language.
Houghton-Mifflin, New York, NY, 1935.

[25] G. K. Zipf.
Human Behaviour and the Principle of Least-Effort.
Addison-Wesley, Cambridge, MA, 1949.