

# MAKING EMPIRICAL POWER-LAW DISTRIBUTIONS USEFUL

OR: HOW TO FIT CURVES & DO STATISTICAL VALIDATION WITH HIGH-  
VARIANCE DATA, AND AN APPLICATION OF THESE TECHNIQUES TO SOME  
EMPIRICAL “POWER LAWS” (WITH PREDICTABLE RESULTS)

Aaron Clauset  
Santa Fe Institute

Algorithms, Inference, and Statistical Physics Workshop  
4 May 2007

Joint with Cosma Shalizi and M. E. J. Newman

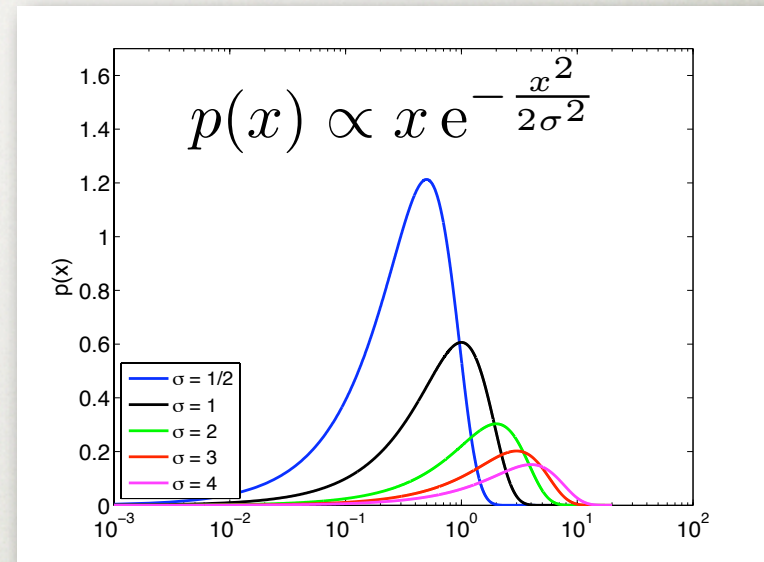


# STATISTICAL DISTRIBUTIONS

---

## Low-variance (“thin-tailed”) distributions

- Maxwell-Boltzmann
- Gaussian (Normal)
- Rayleigh
- etc.



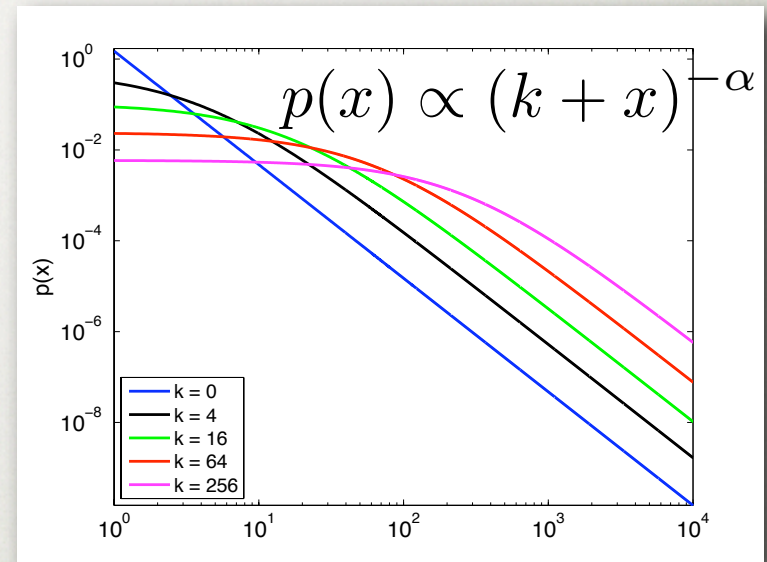
Mean is representative of almost all samples

# STATISTICAL DISTRIBUTIONS

---

## High-variance (“heavy-tailed”) distributions

- Zipf
- Yule-Simon
- “power law”
- etc.



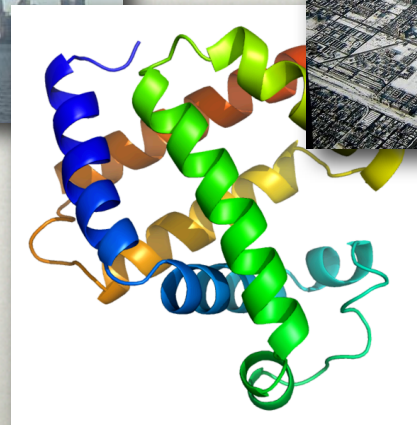
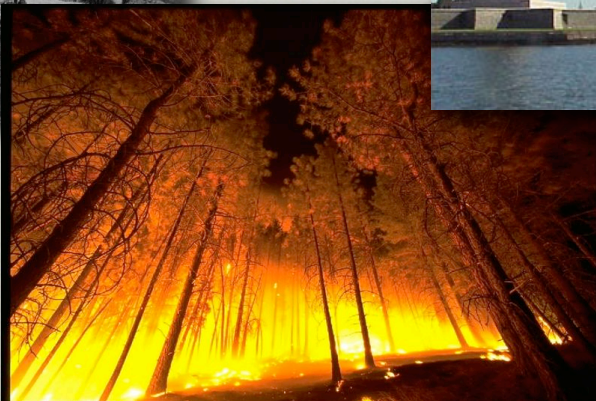
Mean is *not* representative of most samples

# MOTIVATION

---

PL distributions are mathematically **interesting**.  
(imply heterogeneity & scalability)

PL distributions **appear ubiquitous in data**.



# SOME PROBLEMS

---

- Common methods misestimate the scaling behavior  $\alpha$

*“Doesn’t scale the way you think it does.”*

- Behavior rarely validated (i.e., compared to alternative, non-PL models)

*“Maybe it’s **not** a power law.”*

# SOME GOALS

---

**How can empirical PL distributions be useful?**

- Useful for model selection

*“What kind of theory should I look for?”*

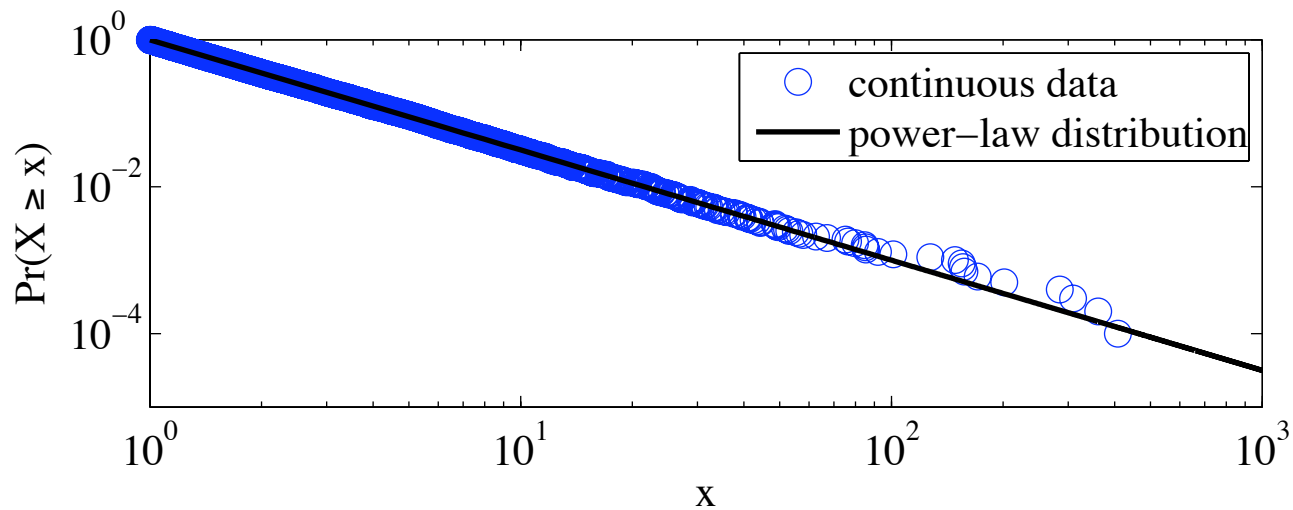
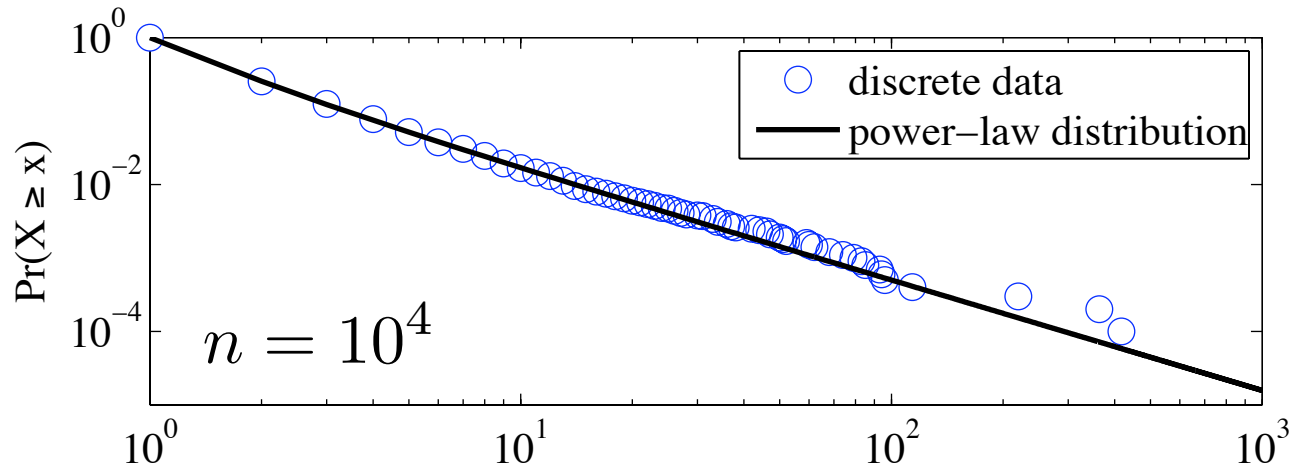
- Useful for (soft) validation

*“My theory explains my data.”*

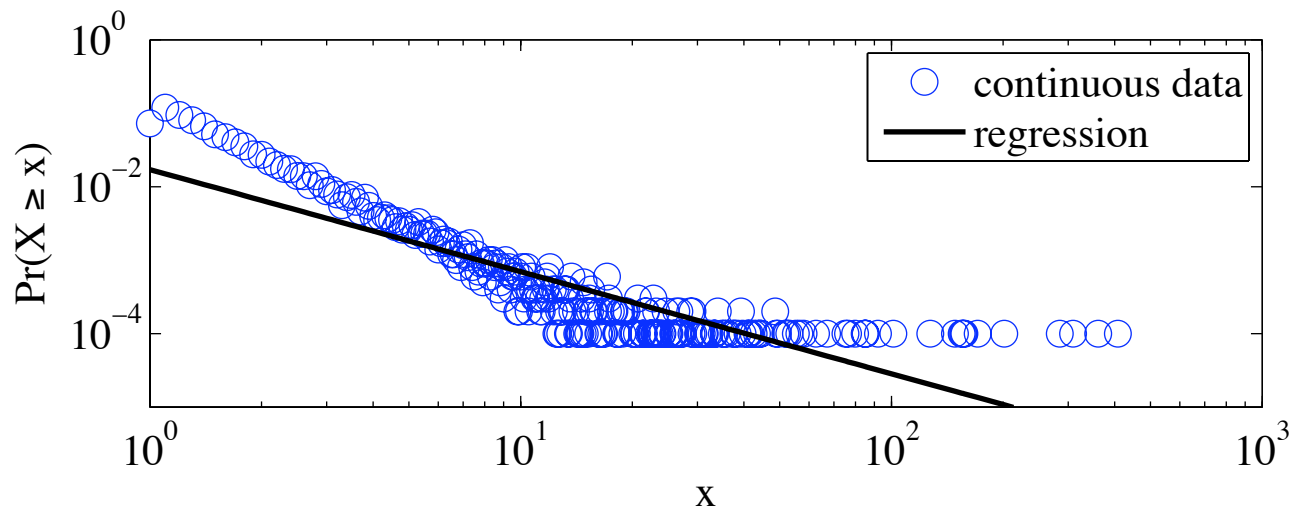
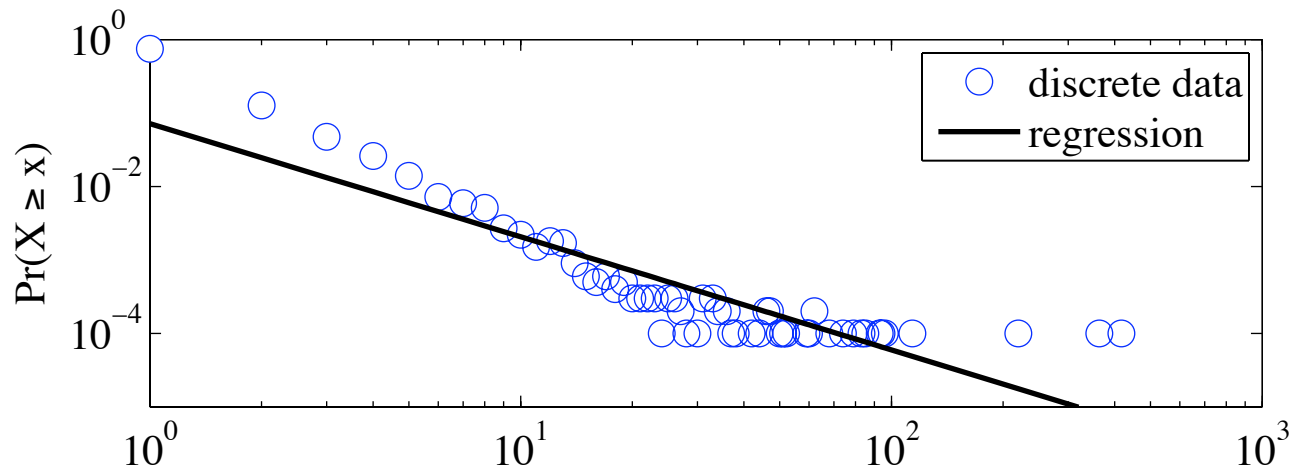
- Deviations point to interesting phenomena

*“My theory doesn’t explain this part.”*

# ESTIMATING $\alpha$



# REGRESSION ON PDF





# ML ESTIMATORS

---

## Continuous power-law distribution

$$\hat{\alpha} = 1 + n \left/ \left[ \sum_{i=1}^n \ln \frac{x_i}{x_{\min}} \right] \right. \quad \sigma_{\hat{\alpha}} \simeq \frac{\hat{\alpha} - 1}{\sqrt{n}}$$

## Discrete power-law distribution

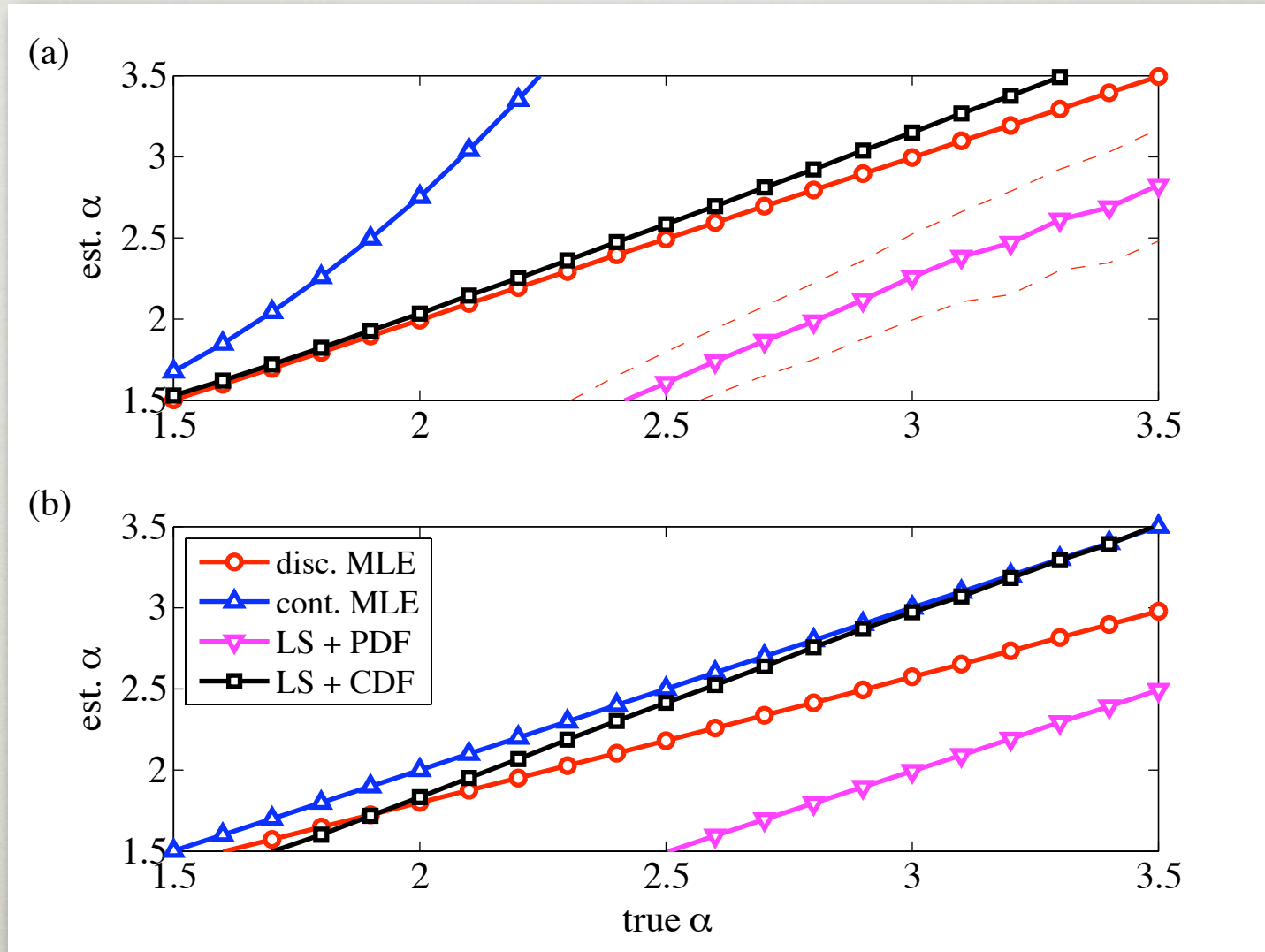
$$\frac{\zeta'(\hat{\alpha}, x_{\min})}{\zeta(\hat{\alpha}, x_{\min})} = -\frac{1}{n} \sum_{i=1}^n \ln x_i$$

# COMPARISON

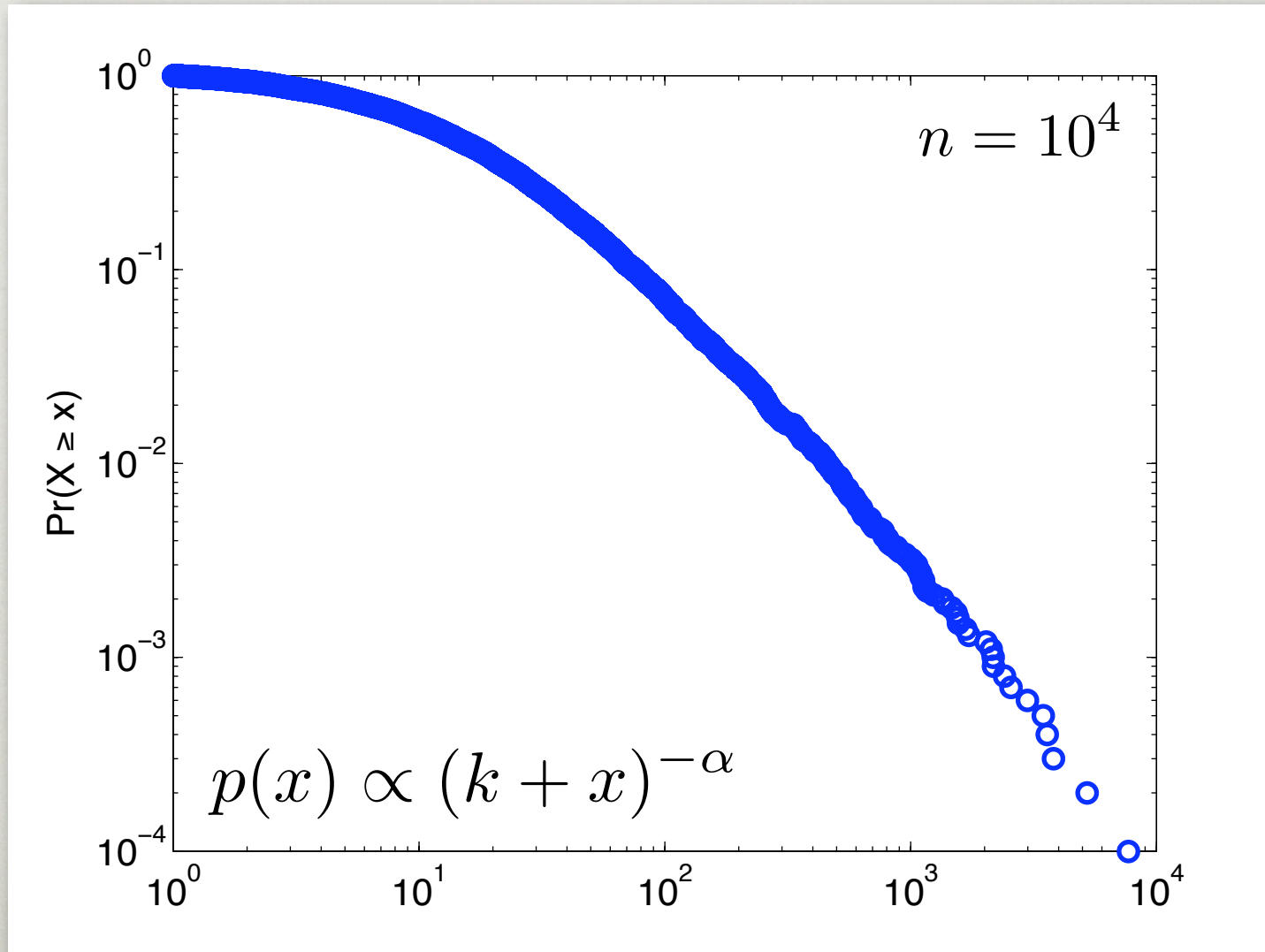
---

method	notes	est. $\alpha$ (discrete)	est. $\alpha$ (continuous)
LS + PDF	const. width	1.5(1)	1.39(5)
LS + CDF	const. width	2.37(2)	2.480(4)
LS + PDF	log. width	1.5(1)	1.19(2)
LS + CDF	rank-freq	2.570(6)	2.4869(3)
cont. MLE	—	4.46(3)	<b>2.50(2)</b>
disc. MLE	—	<b>2.49(2)</b>	2.19(1)

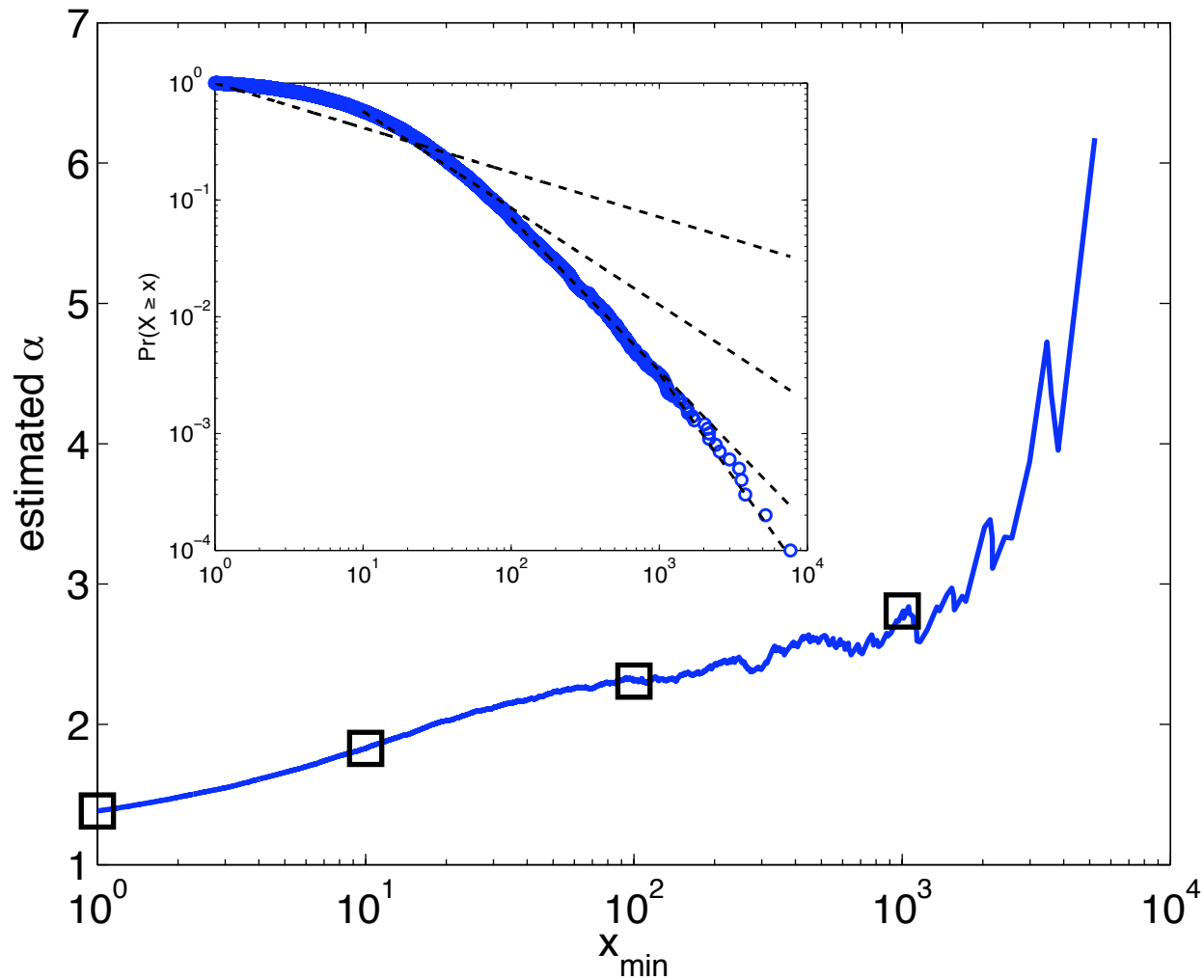
# IN GENERAL



# ESTIMATING $x_{\min}$



# A FEW CHOICES



# AN OBJECTIVE METHOD

---

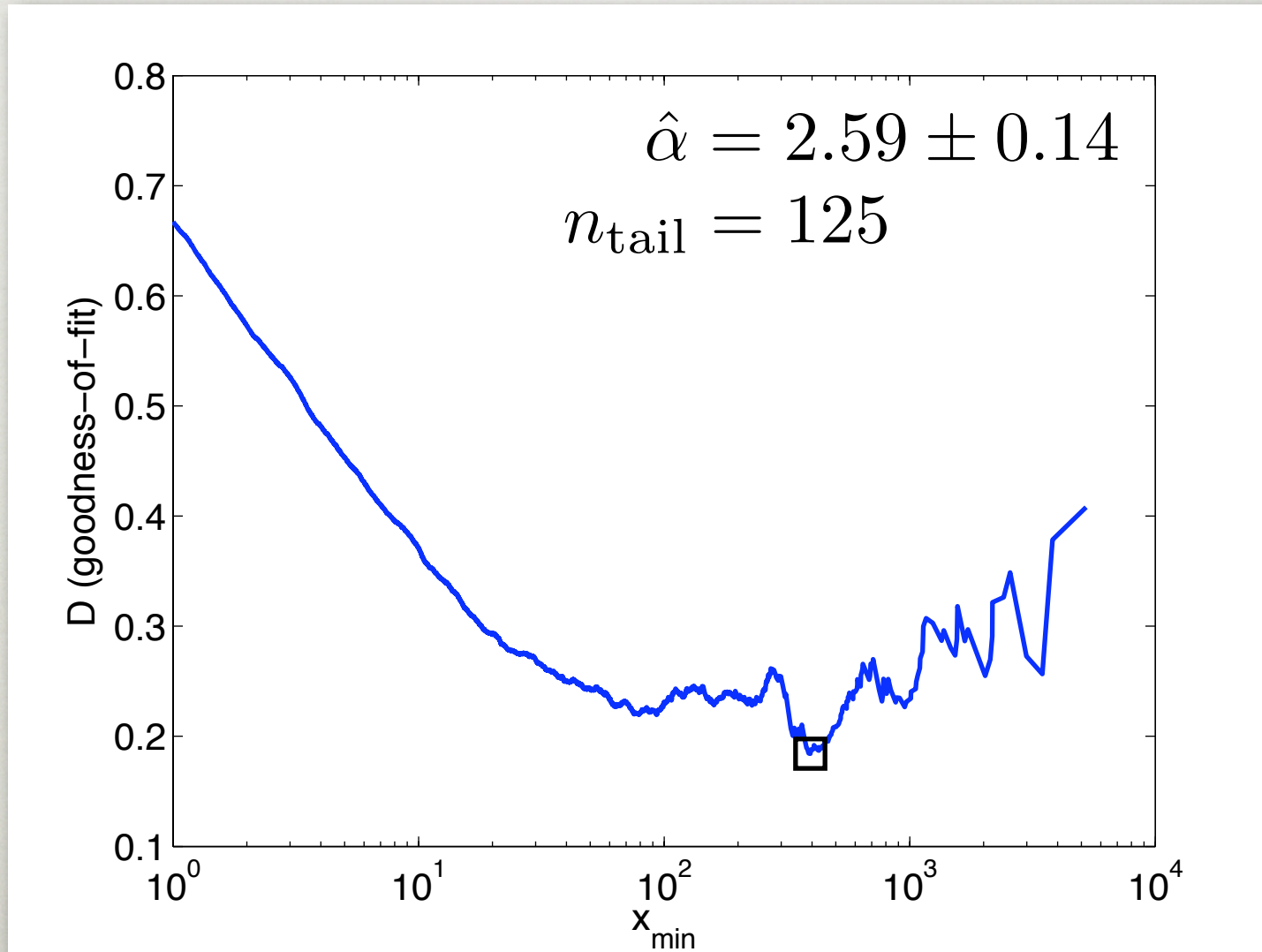
Choose power law that minimizes “distance”  
between model and data:

Kolmogorov-Smirnov goodness-of-fit statistic

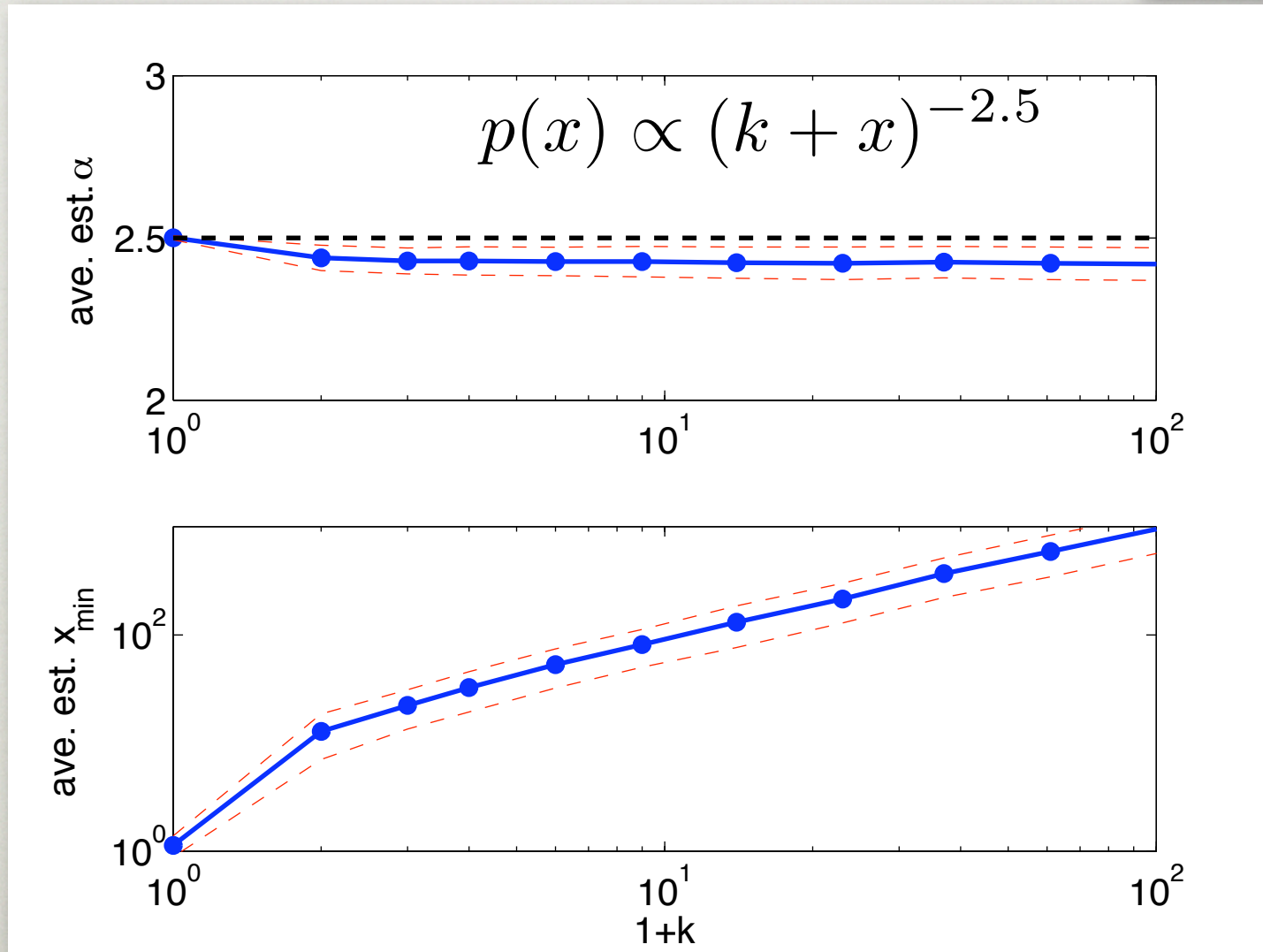
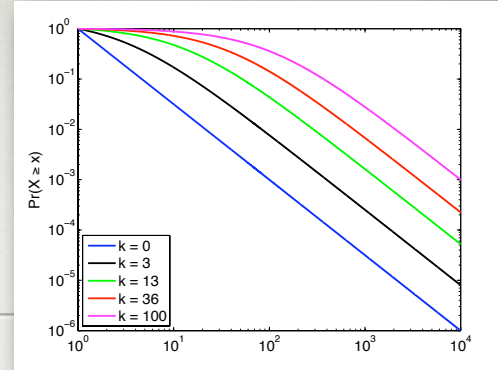
$$D(x_{\min}) = \max_{x \geq x_{\min}} |S_n(x) - P(x)|$$

We choose  $\hat{x}_{\min} = \min_y D(y)$

# GoF RESULTS



# IN GENERAL



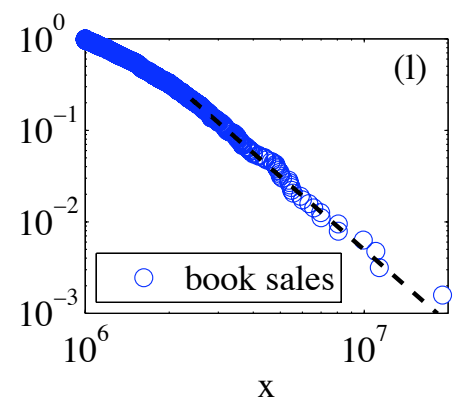
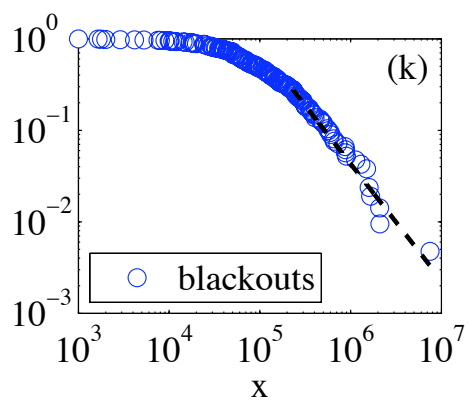
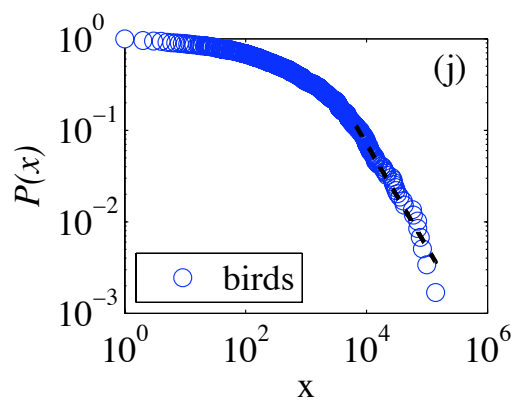
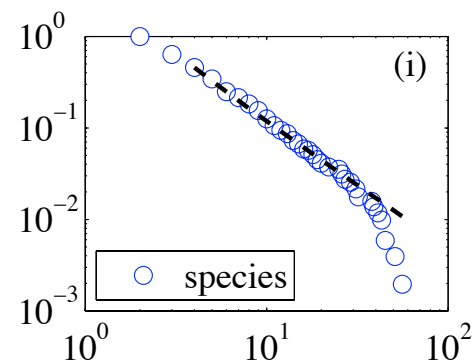
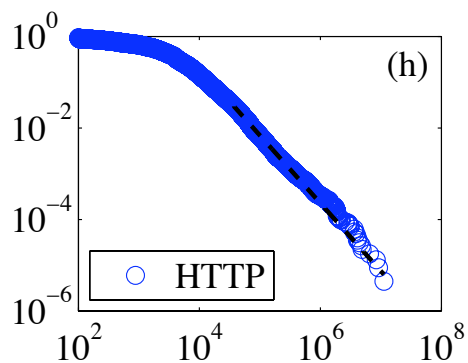
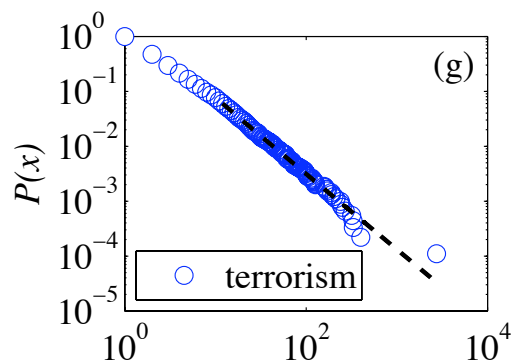
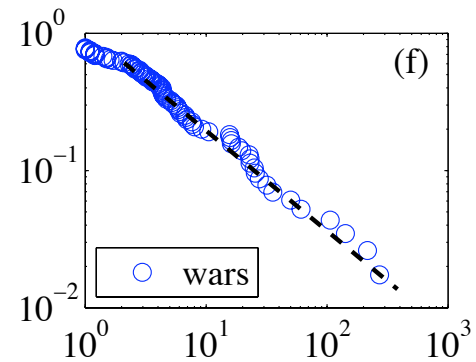
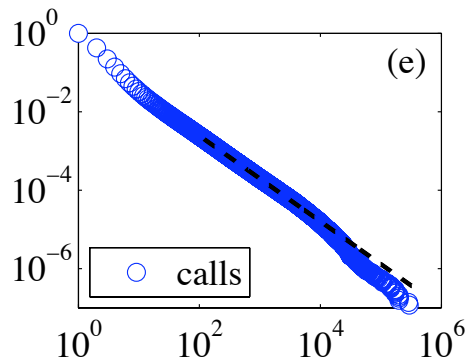
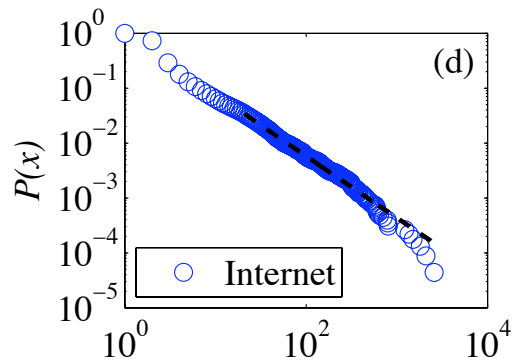
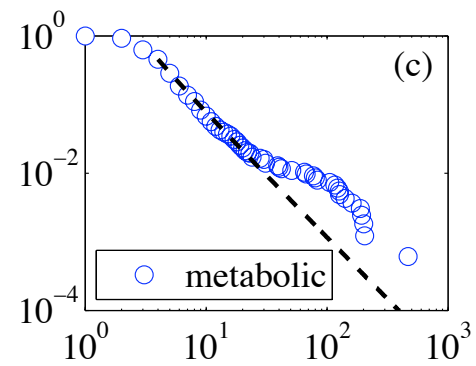
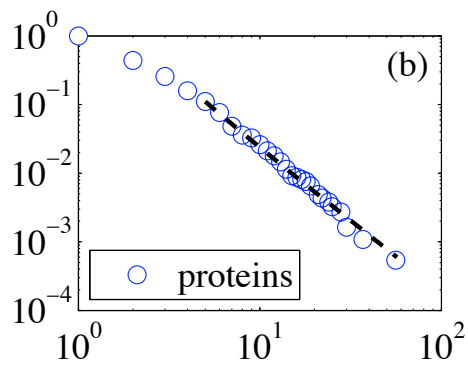
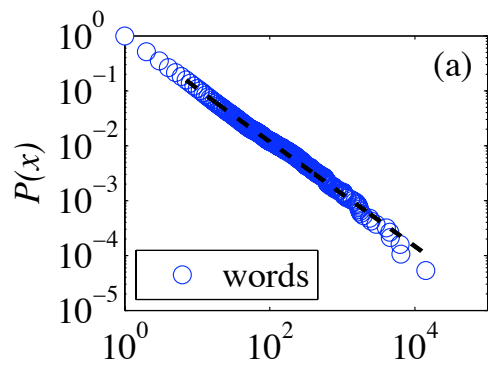


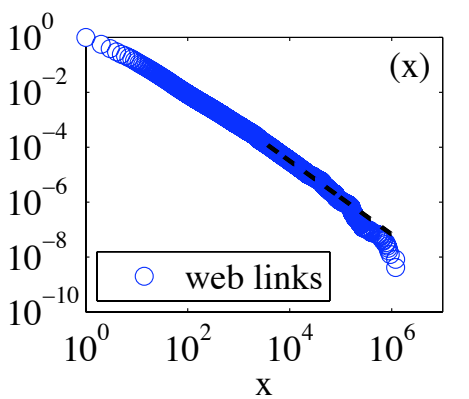
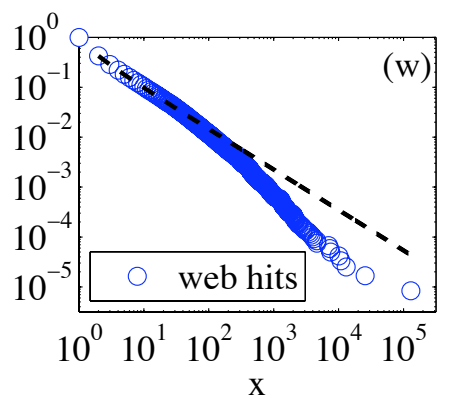
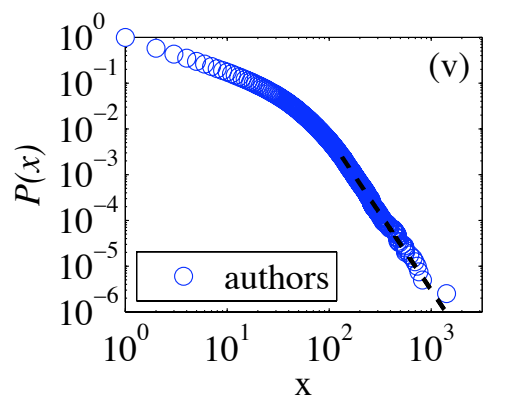
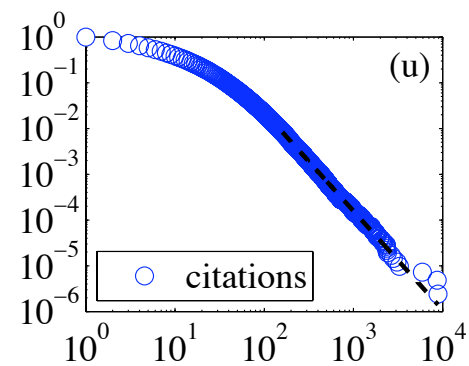
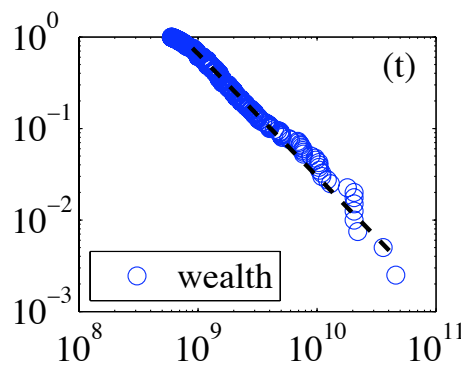
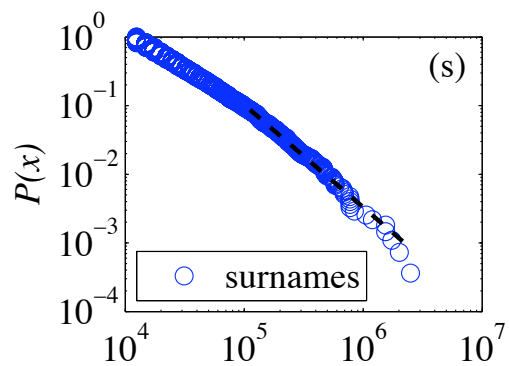
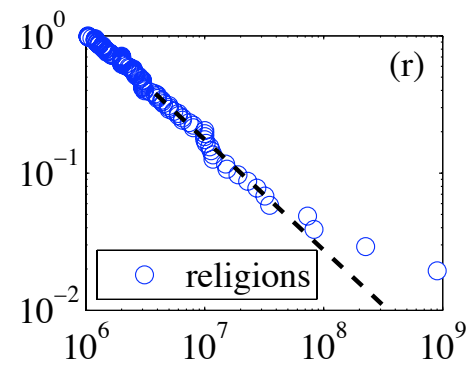
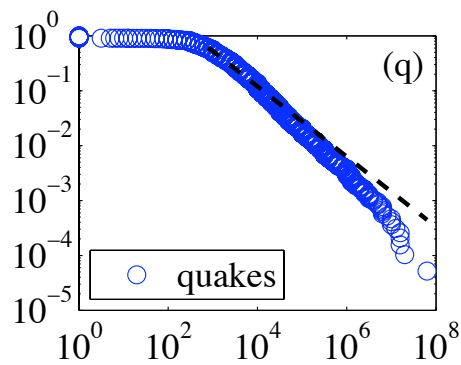
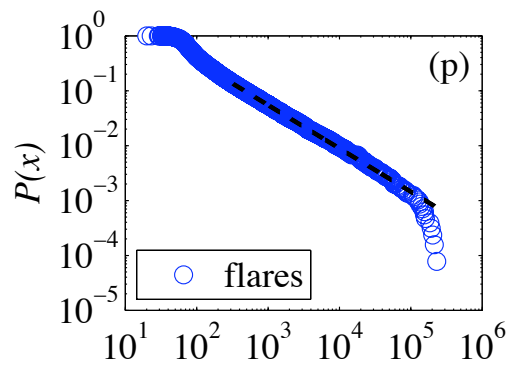
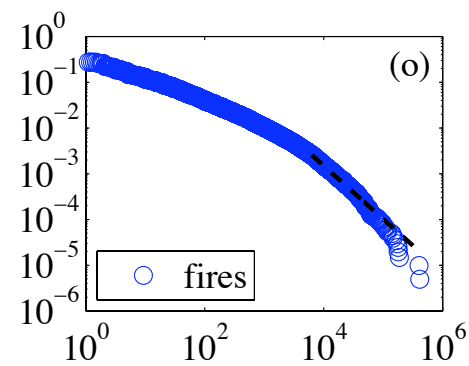
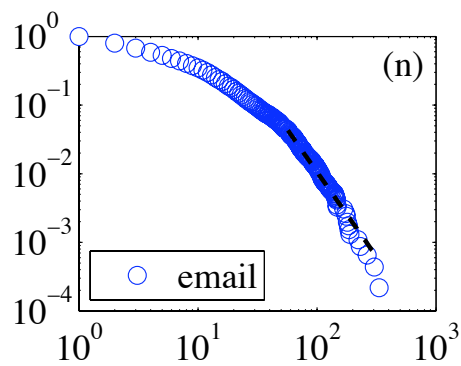
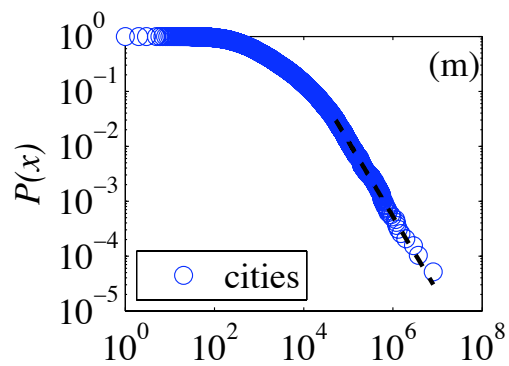
# SOME REAL DATA

---

count of word use  
protein interaction degree  
metabolic degree  
Internet degree  
telephone calls received  
intensity of wars  
terrorist attack severity  
HTTP session size (kB)  
species per genus  
bird species sightings  
blackouts  
sales of books

population of cities  
email address books size  
forest fire size (acres)  
solar flare intensity  
quake intensity  
religious followers  
freq. of surnames  
net worth (mil. USD)  
citations to papers  
papers authored  
hits to web sites  
links to web sites





# VALIDATION

---

Two approaches:

- $p$ -value

*“Are deviations from PL model explained by statistical fluctuations?”*

- likelihood ratio (LR)

*“Does this other model look more like the data?”*

# SUPPORT FOR POWER LAW

---

data set	power law $p$	log-normal LR	w. cut-off LR	support for power law
cities	<b>0.761</b>	-0.435	-0.298	reasonable
fire	0.045	-1.78	<b>-5.02</b>	w. cut-off
flares	0.997	-0.803	<b>-4.52</b>	w. cut-off
HTTP	0.000	1.59	0.000	none
quakes	0.000	-7.14	<b>-24.4</b>	w. cut-off
wealth	0.001	-0.0777	-0.198	none
web hits	0.000	0.255	0.000	none

# SUPPORT FOR POWER LAW

---

data set	power law $p$	log-normal LR	w. cut-off LR	support for power law
Internet	<b>0.286</b>	-0.807	-1.97	reasonable
citations	<b>0.204</b>	-0.141	-0.007	reasonable
metabolic	0.000	-1.05	0.000	none
species	0.103	-1.63	<b>-3.80</b>	w. cut-off
terrorism	<b>0.684</b>	-0.278	-0.077	reasonable
words	<b>0.487</b>	0.395	-0.899	good

# CONCLUSIONS

---

- Only maximum likelihood accurately estimates  $\alpha$
- Can now accurately (and objectively) estimate  $x_{\min}$
- Validating power-law distributions possible
- Some “power laws” should be “revisited”
- Deviations from PL can suggest new hypotheses

**Thanks (comments and data):** L. Adamic, A. Boyer, A. Broder, A. Downey, J.D. Farmer, P. Holme, M. Huss, J. Karlin, J. Ladau, M. Mitzenmacher, C. Moore, S. Redner, S. Stoev, M. Wheatland, J. Wiener, W. Willinger and M. Young

**FIN**