



Learning from Data

Aaron Clauset
@aaronclauset

Assistant Professor of Computer Science
University of Colorado Boulder
External Faculty, Santa Fe Institute

distribution $p(x)$

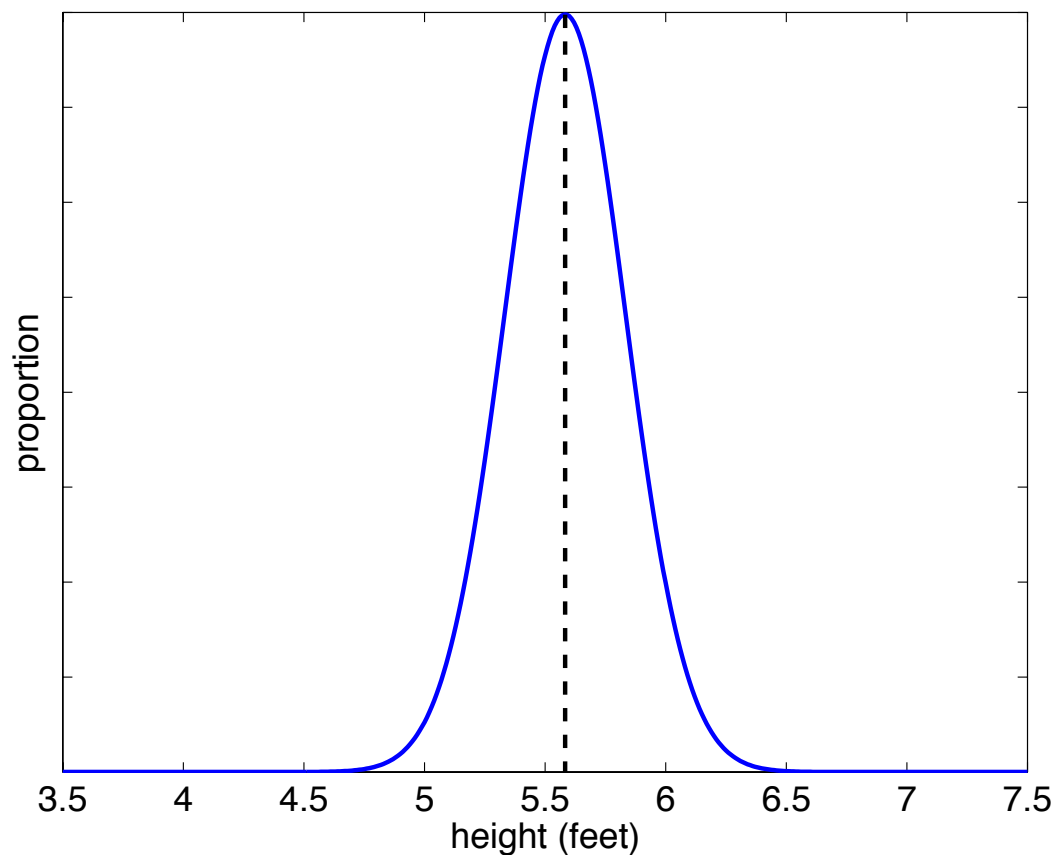
the fraction of times we observe
an event of size x

normal distribution

the bell curve

height

average height
of an American
is 5.583 feet

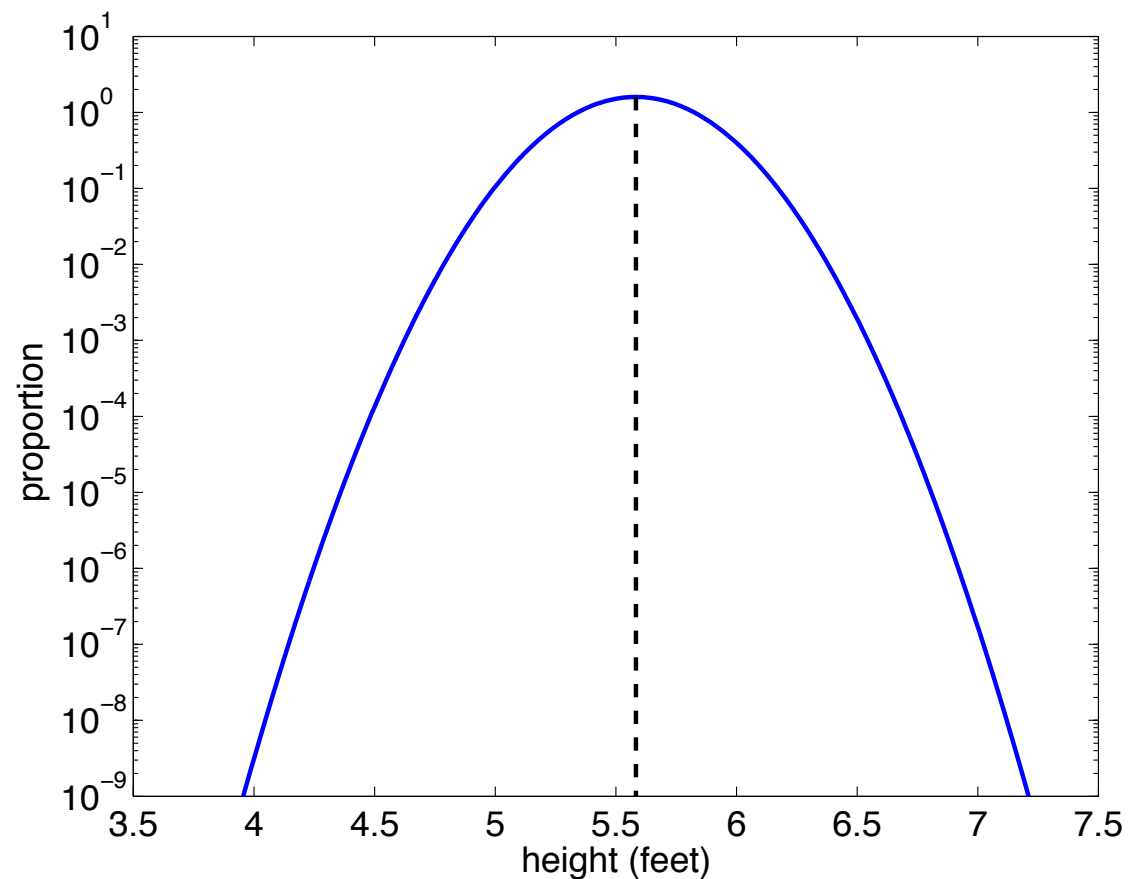


normal distribution

the bell curve

height

average height
of an American
is 5.583 feet

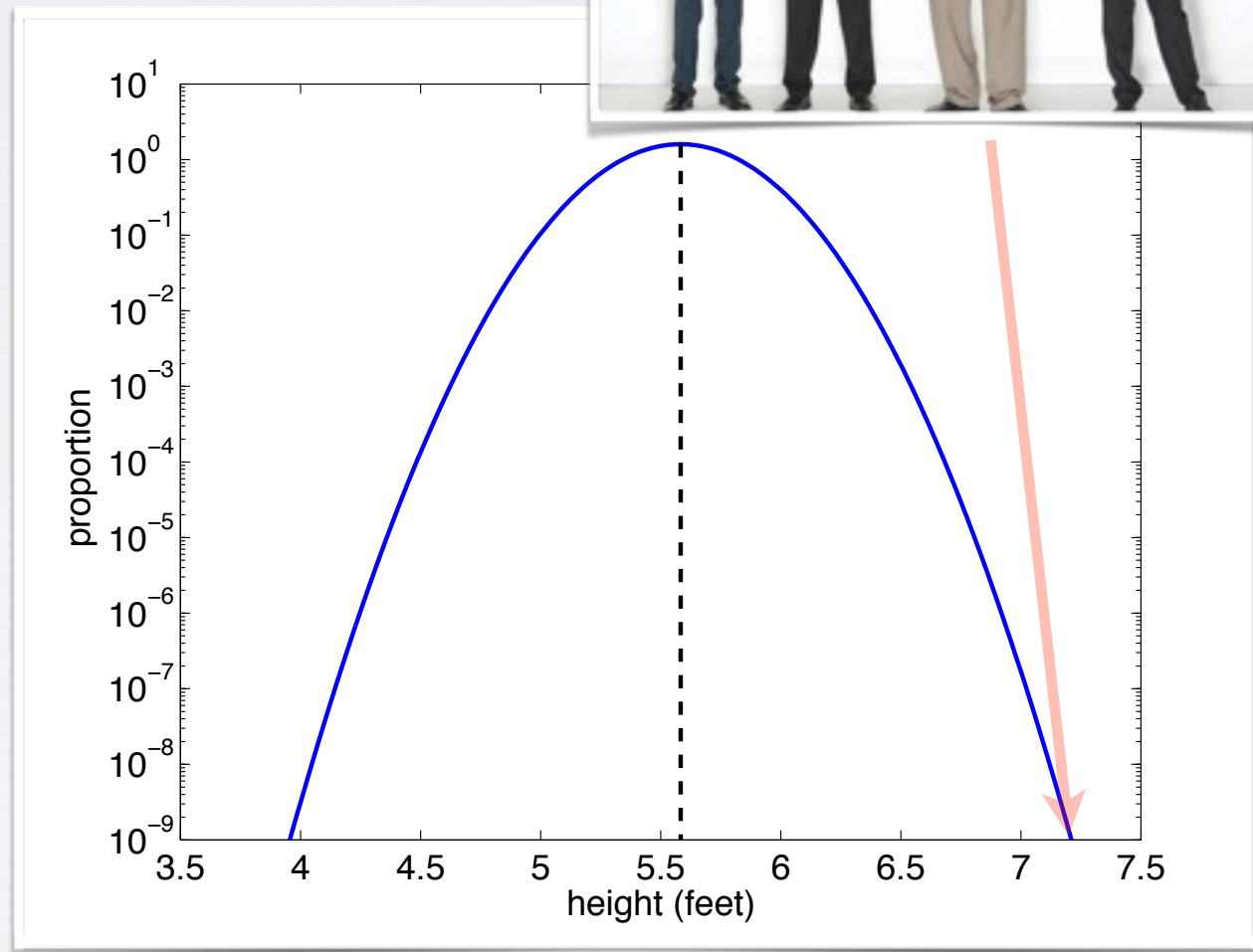


normal distribution

the bell curve

height

average height
of an American
is 5.583 feet



normal distribution

the bell curve

average is representative

lengths (human height, etc.)

weights (human or otherwise)

speeds (highway, running, etc.)

power law distributions

“heavy-tailed” patterns

average is not representative

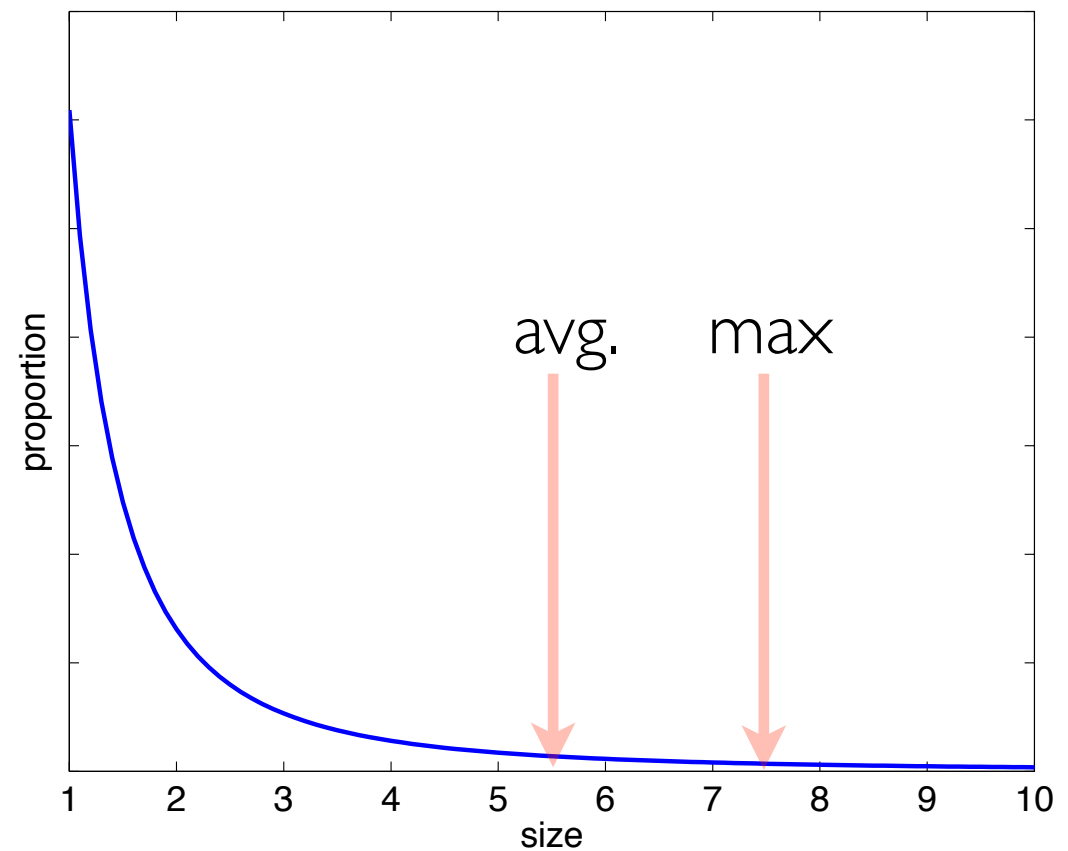
power law distributions

“heavy-tailed” patterns

power law

with average
5.583

most values
very small!



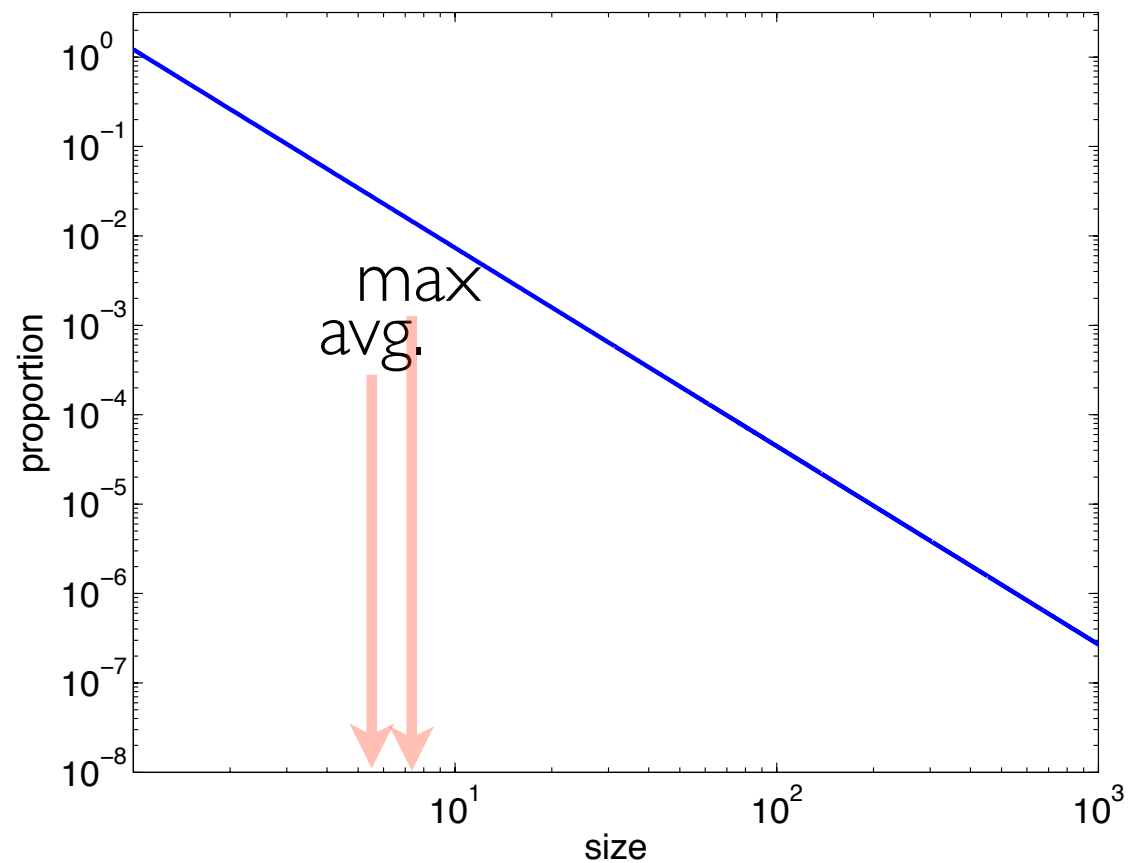
power law distributions

“heavy-tailed” patterns

power law

with average
5.583

most values
very small



power law distributions

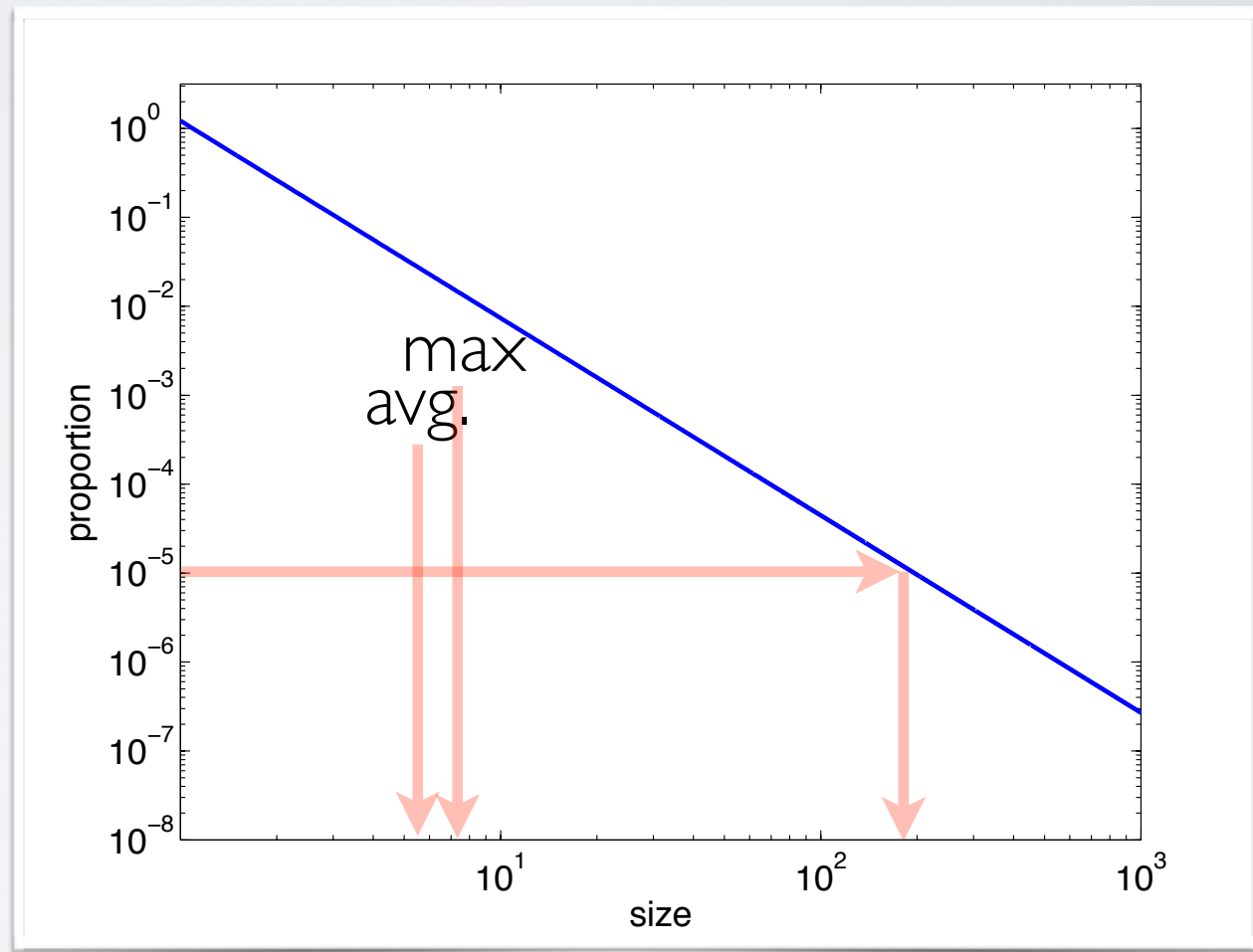
“heavy-tailed” patterns

power law

with average
5.583

most values
very small

but a few values
VERY big



power law distributions

“heavy-tailed” patterns

power law distributions

average is not representative

earthquake energy

solar flare energy

flood water volume

forest fire size

landslide size

lunar crater size

terrorism events

international wars

book sales

electrical blackouts

financial wealth

city population

financial returns

surname frequency

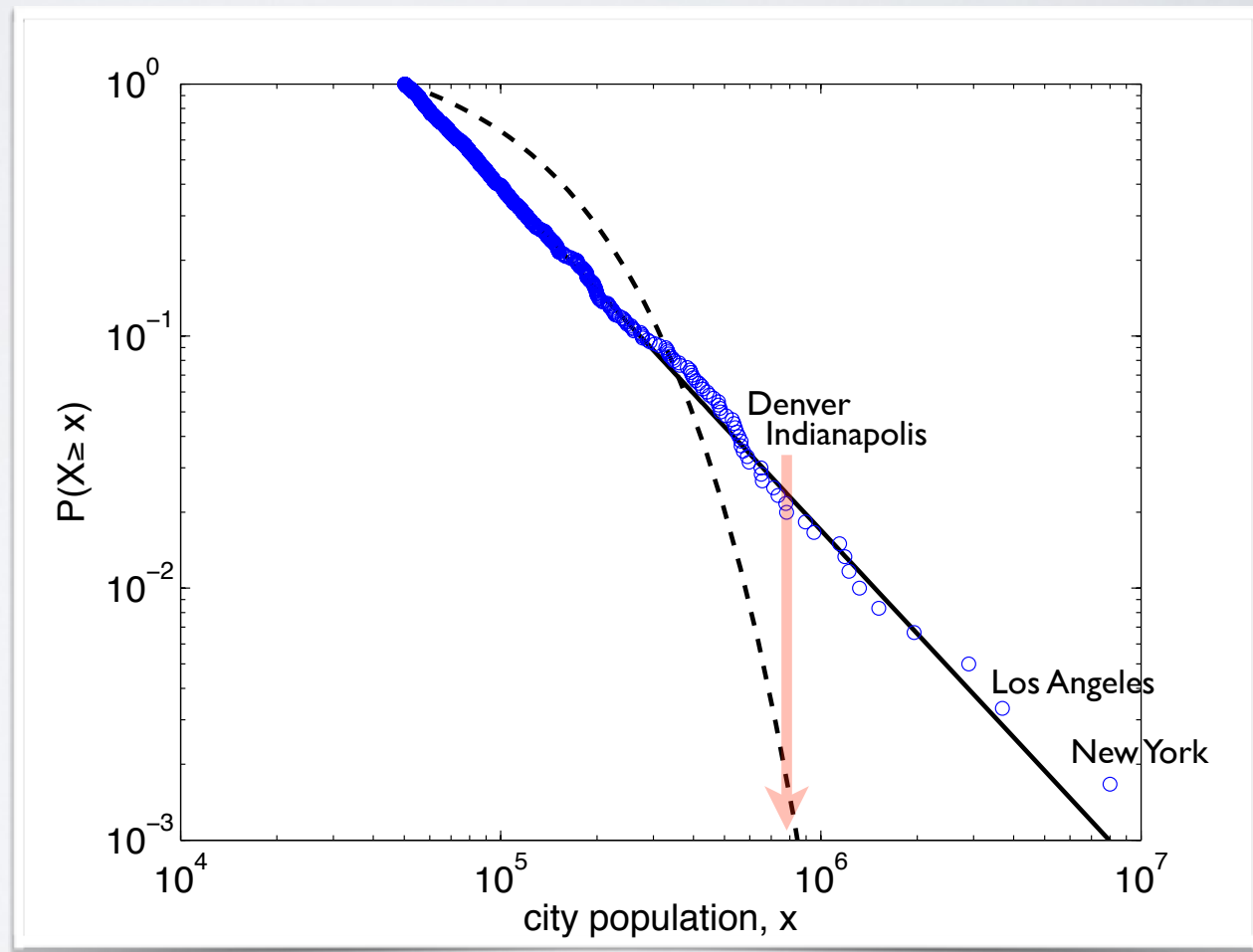
power laws vs. normals

counter-intuitive

power laws vs. normals

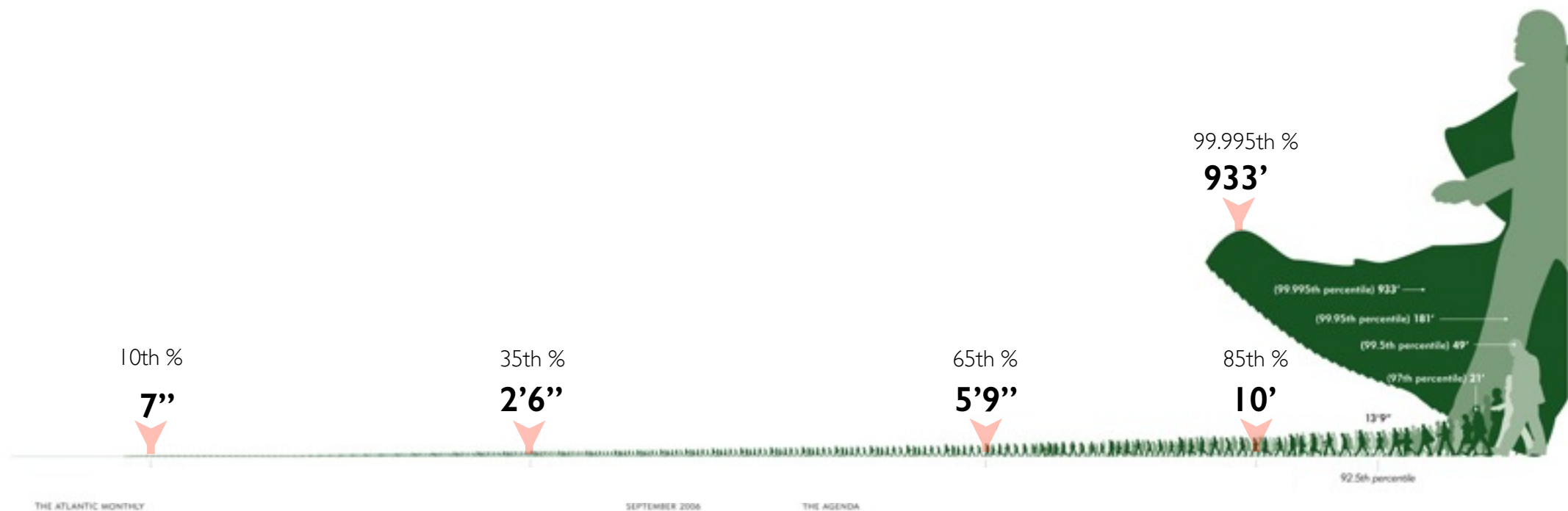
counter-intuitive

500 largest US cities



power laws vs. normals
counter-intuitive

financial wealth as human heights



"The Height of Inequality" (2006) in The Atlantic

terrorism

terrorism

MIPT **TERRORISM**
KNOWLEDGE BASE SM

 **incident profile**

ABU HAFS AL-MASRI BRIGADE AND SECRET ORGANIZATION OF AL-QAEDA IN EUROPE ATTACKED TRANSPORTATION TARGET (JULY 7, 2005, UNITED KINGDOM)

Incident Date: July 7, 2005

Terrorist Organization(s): Abu Hafs al-Masri Brigade , Secret Organization of al-Qaeda in Europe

Target: Transportation

City: London

Country: United Kingdom

Region: Western Europe

Tactic: Bombing

Weapon: Explosives

Fatalities: 27

Injuries: 0



RAND-MIPT data

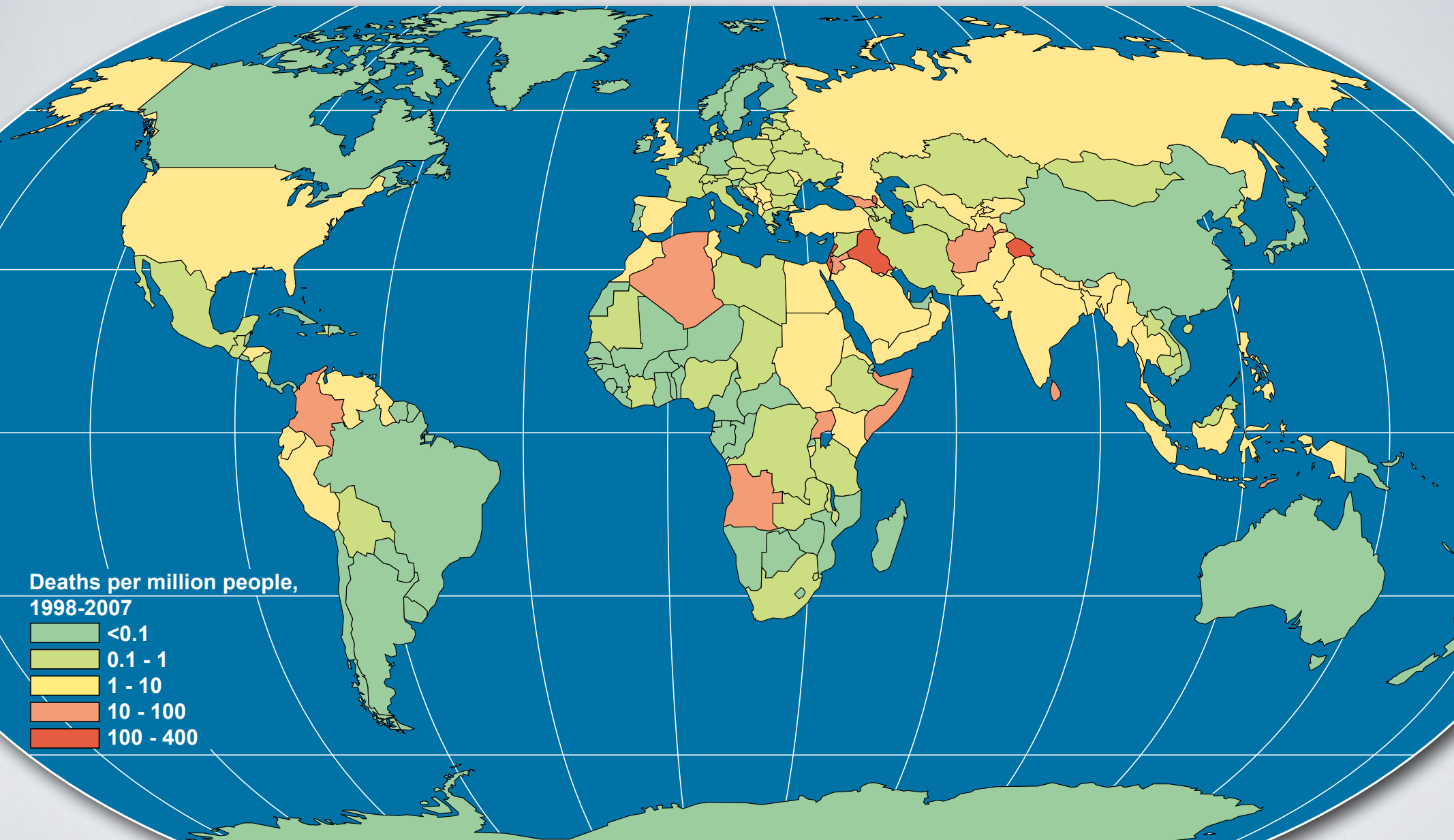
- 40 years (1968-2008)
- domestic + international
- 5000+ cities, 187 countries
- 36,018 events (37% deadly)

terrorism

- where does terrorism occur?
- what is risk of dying from terrorism?

data analysis:

1. *take all events 1998-2007*
2. *count deaths in each country*
3. *divide total by country's population*
4. *yields per capita risk of death*
5. *visualize on a world map*



deaths per million people, USA 2007

terrorism	0	--
lightning	0.15	1
bee sting	0.18	x1.2
airplane crash	0.23*	x1.5
homicide	61.74*	x408.3
car crash	124.43	x829.5

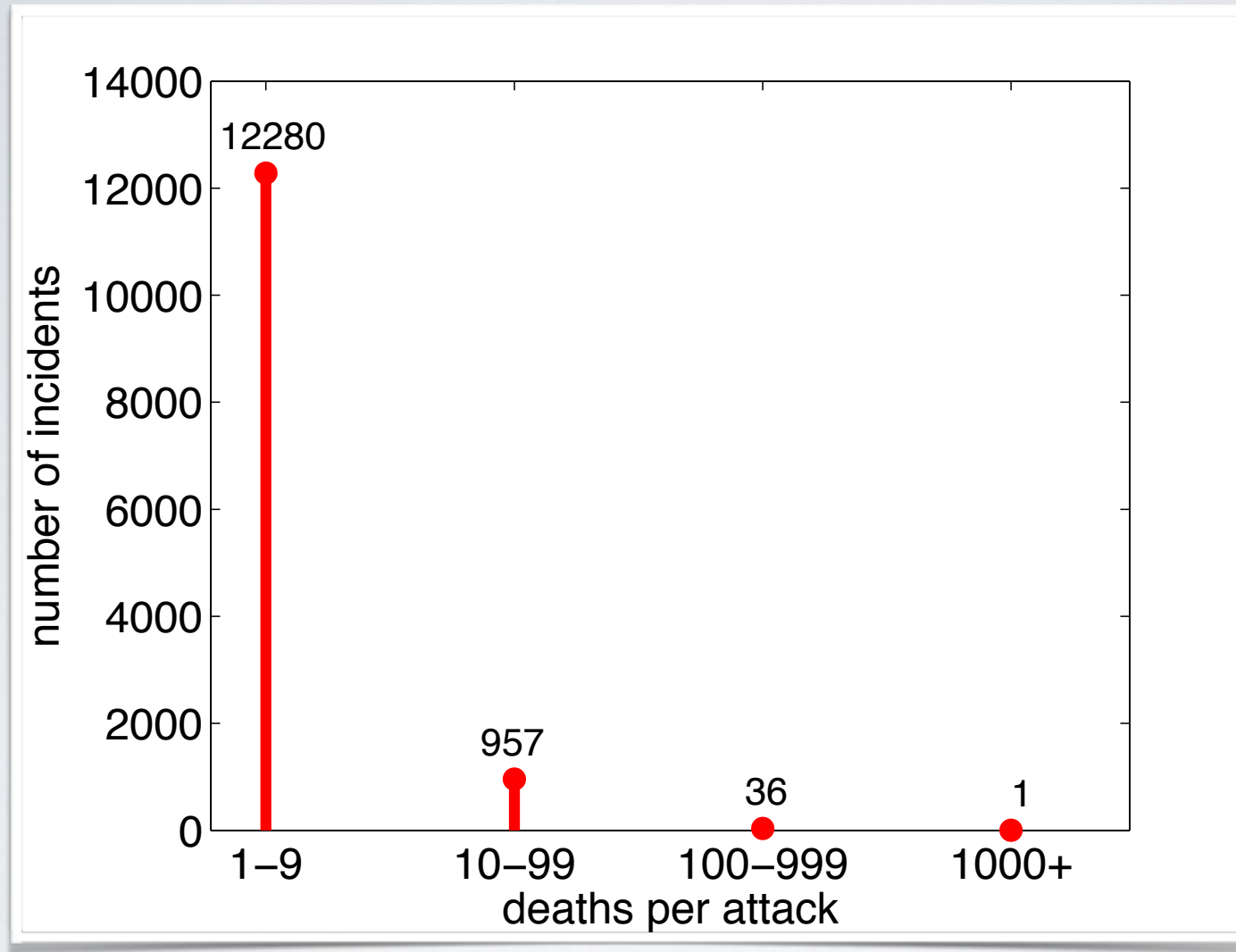
terrorism

- how many people die in a terrorist event?

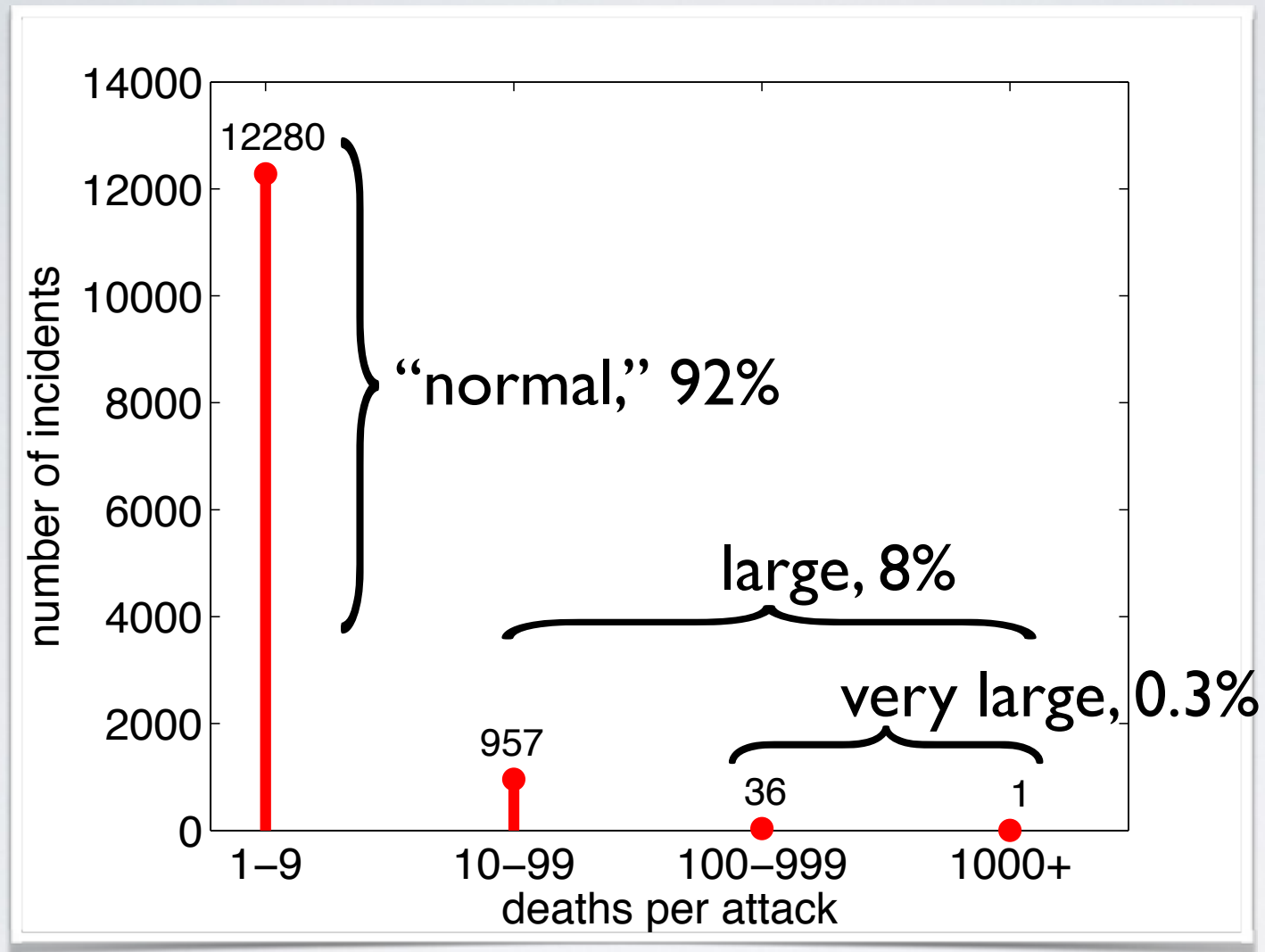
data analysis:

1. *take all events*
2. *count # times 1 death, 2 deaths, etc.*
3. *visualize as distribution*

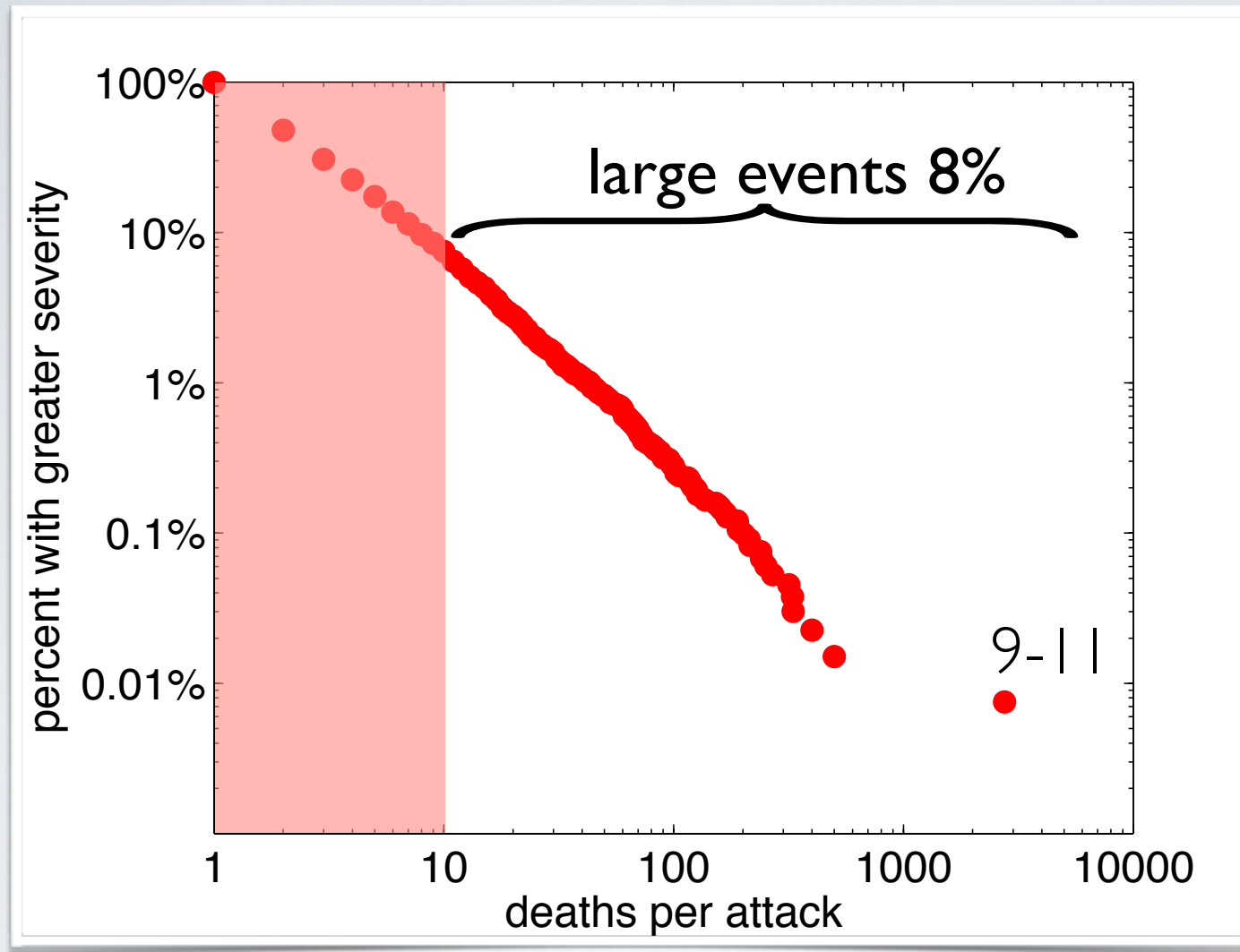
deadly terrorist events, 1968-2008



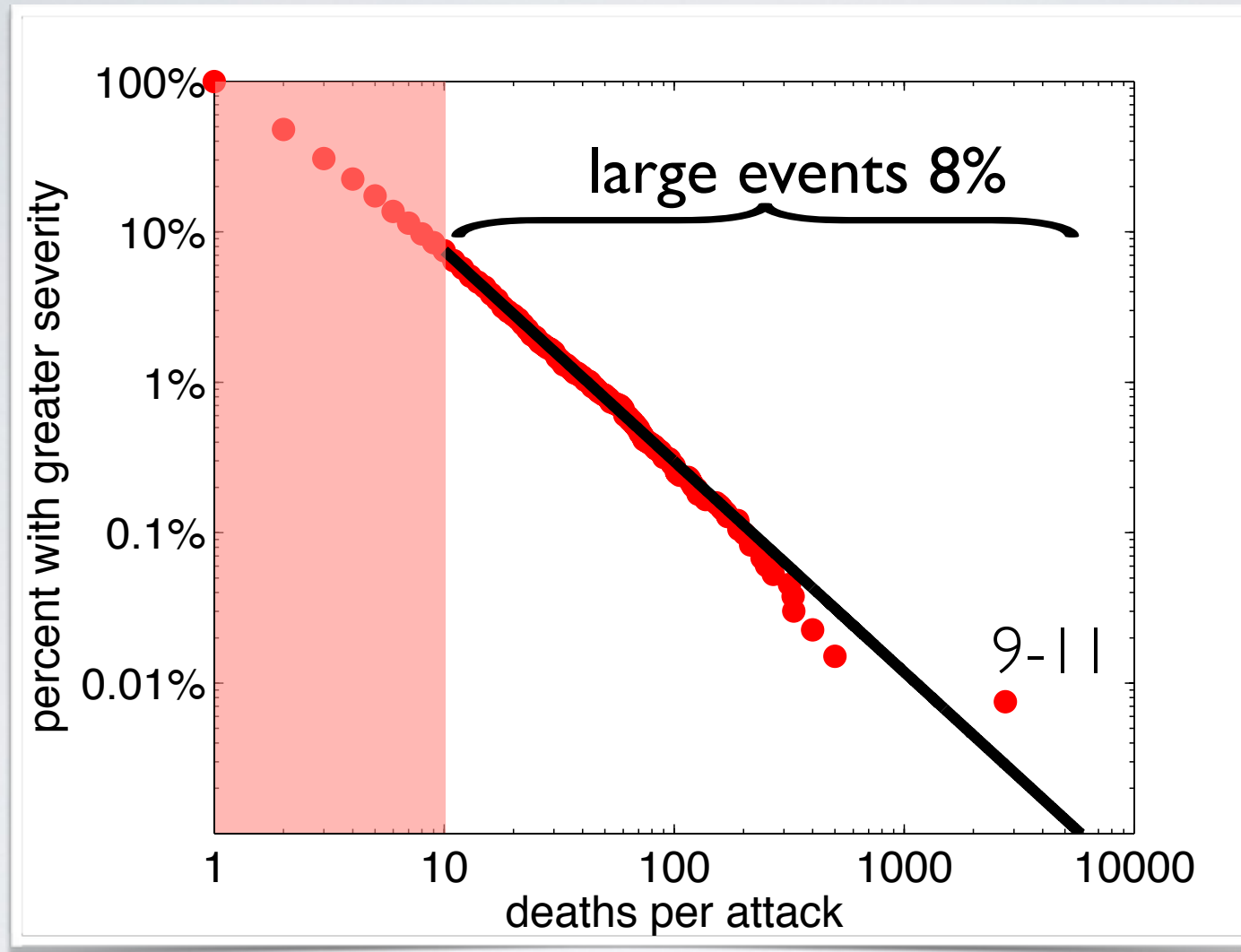
deadly terrorist events, 1968-2008



deadly terrorist events, 1968-2008



deadly terrorist events, 1968-2008



it follows a power-law distribution

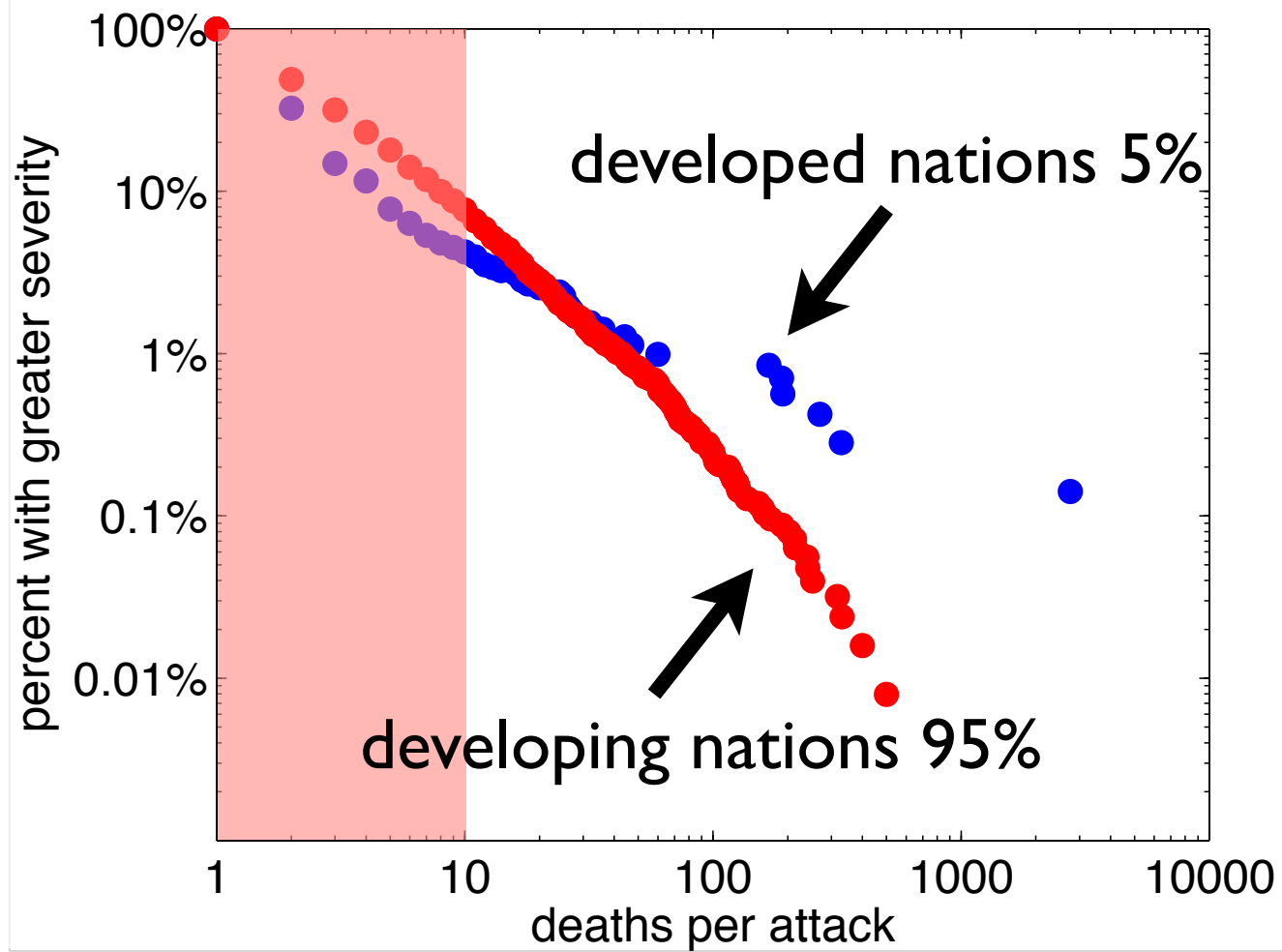
up by 10x in severity = down by 250x in frequency

terrorism

- do big events happen everywhere equally?

data analysis:

1. *divide events by (i) developed nation or (ii) developing nation*
2. *for each type, count # times 1 death, 2 deaths, etc.*
3. *visualize the two distributions*



terrorism

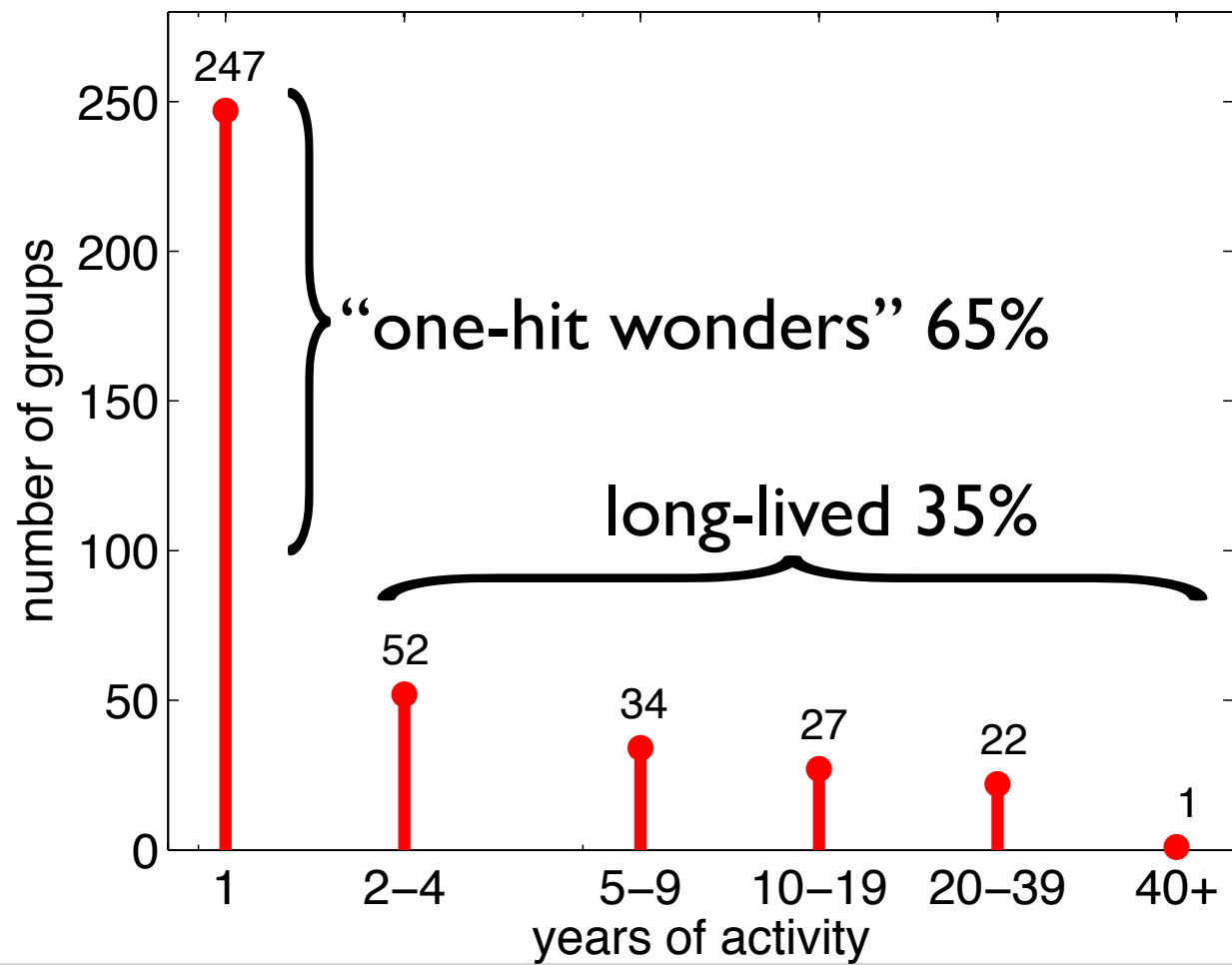
- how long do terrorist groups last?

data analysis:

1. *for each unique terrorist group*
2. *count # years between oldest and newest event*
3. *then count # groups with 1 year, 2 years, etc. of activity*
4. *visualize the distribution*

major terrorist organizations | 1968-2008

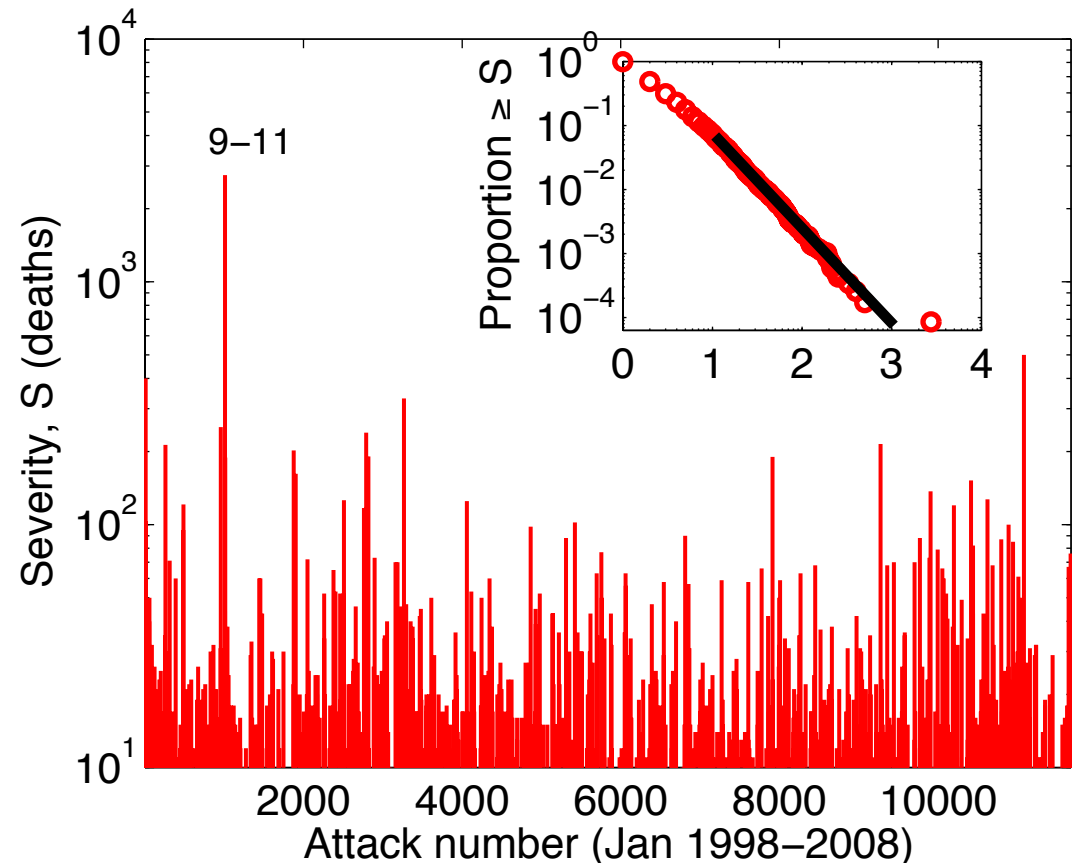
1. Revolutionary Armed Forces of Colombia (FARC)
2. Hamas
3. Taliban
4. Basque Fatherland and Freedom (ETA)
5. Communist Party of Nepal-Maoist (CPN-M)
6. National Liberation Army (Colombia)
7. Palestinian Islamic Jihad (PIJ)
8. Liberation Tigers of Tamil Eelam (LTTE)
9. al-Fatah
10. Communist Party of India-Maoist
11. al-Qaeda Organization in the Land of the Two Rivers
12. Anti-Castro Cubans
13. Hezbollah
14. Fronte di Liberazione Naziunale di a Corsica (FLNC)
15. Shining Path
16. Islamic State of Iraq
17. Popular Front for the Liberation of Palestine (PFLP)
18. United Liberation Front of Assam (ULFA)
19. al-Aqsa Martyrs Brigades
20. Kurdistan Workers' Party (PKK)
21. Tupac Amaru Revolutionary Movement
22. Ansar al-Sunnah Army
23. Black September
24. New People's Army (NPA)
25. Abu Nidal Organization (ANO)
26. Mujahideen Shura Council
27. Armenian Secret Army for the Liberation of Armenia
28. Irish Republican Army (IRA)
29. Revolutionary People's Liberation Party/Front (DHKP/C)
30. People's War Group (PWG)
31. United Self-Defense Forces of Colombia (AUC)
32. Jewish Defense League (JDL)
33. Amal
34. Armed Islamic Group
35. Palestine Liberation Organization (PLO)
36. Earth Liberation Front (ELF)
37. Abu Sayyaf Group (ASG)
38. Popular Resistance Committees
39. Manuel Rodriguez Patriotic Front
40. Revolutionary Organization 17 November (RO-N17)
41. al-Qaeda Organization in the Islamic Maghreb
42. Baloch Liberation Army (BLA)
43. Revolutionary People's Struggle
44. Red Army Faction
45. Islamic Army in Iraq
46. Democratic Front for the Liberation of Palestine (DFLP)
47. UNITA
48. Revolutionary Nuclei
49. al-Gama'a al-Islamiyya (GAI)
50. Free Aceh Movement (GAM)
51. Kurdistan Freedom Hawks
52. April 19 Movement
53. Lord's Resistance Army (LRA)
54. Moro Islamic Liberation Front (MILF)
55. Real Irish Republican Army (RIRA)
56. al-Qaeda
57. Tawhid and Jihad
58. Popular Liberation Army
59. Eritrean Liberation Front (ELF)
60. Montoneros
61. Turkish Communist Party Marxist-Leninist (TKP/ML-TIKKO)
62. Mozambique National Resistance Movement
63. Ulster Defence Association/Ulster Freedom Fighters
64. Purbo Banglar Communist Party (PBCP)
65. National Liberation Front of Tripura (NLFT)
66. First of October Antifascist Resistance Group (GRAPO)
67. Red Hand Defenders (RHD)
68. Lashkar-e-Taiba (LeT)
69. Hizbul Mujahideen (HM)
70. Mujahideen Youth Movement
71. Bersatu
72. People's Revolutionary Army (Argentina)
73. Farabundo Marti National Liberation Front
74. Jaish-e-Mohammad (JeM)
75. Islamic Jihad Jerusalem
76. Peronist Armed Forces
77. Khmer Rouge
78. Justice Commandos for the Armenian Genocide
79. Continuity Irish Republican Army (CIRA)
80. PKK/KONGRA-GEL
81. National Democratic Front of Bodoland (NDFB)
82. Lautaro Youth Movement
83. Action Directe
84. Polisario Front
85. Mujahedin-e-Khalq (MeK)
86. Maoist Communist Center (MCC)
87. Popular Forces of April 25
88. Third of October Group
89. Baader-Meinhof Group
90. Breton Revolutionary Army (ARB)
91. Orly Organization
92. People's Liberation Forces (El Salvador)
93. Front for the Liberation of the Cabinda Enclave
94. Abu al-Rish Brigades
95. African National Congress (South Africa)
96. Moro National Liberation Front (MNLF)
97. Islamic Great Eastern Raiders Front
98. Palestinian Revolution Forces General Command
99. Chukakuha
100. Communist Combatant Cells
101. Popular Front for the Liberation of Palestine -- General Command (PFLP-GC)
102. Red Brigades
103. Japanese Red Army (JRA)
104. Animal Liberation Front (ALF)
105. Committee of Solidarity with Arab and Middle East Political Prisoners (CSPPA)
106. Front for the Liberation of Lebanon from Foreigners (FLLF)
107. Jamatul Mujahedin Bangladesh
108. Informal Anarchist Federation
109. Sudan People's Liberation Army
110. Ninth of June Organization
111. Guerrilla Army of the Poor
112. Loyalist Volunteer Force (LVF)
113. Anti-Imperialist International Brigade
114. All Tripura Tiger Force (ATTf)
115. People's Revolutionary Army (Colombia)
116. Social Resistance
117. Arab Communist Organization (ACO)
118. Anti-Terrorist Liberation Group
119. Riyadh us-Saliheyn Martyrs' Brigade
120. Kosovo Liberation Army (KLA)
121. Lashkar-e-Jhangvi (LeJ)
122. Revolutionary United Front (RUF)
123. Jamiat ul-Mujahedin (JuM)
124. Alex Boncayao Brigade (ABB)
125. Pattani United Liberation Organization (PULO)
126. Group Bakunin Gdansk Paris Guatemala Salvador
127. Irish National Liberation Army (INLA)
128. Revolutionary Struggle
129. Lebanese Armed Revolutionary Faction
130. Ananda Marga
131. Tupamaros
132.



terrorism

terrorism

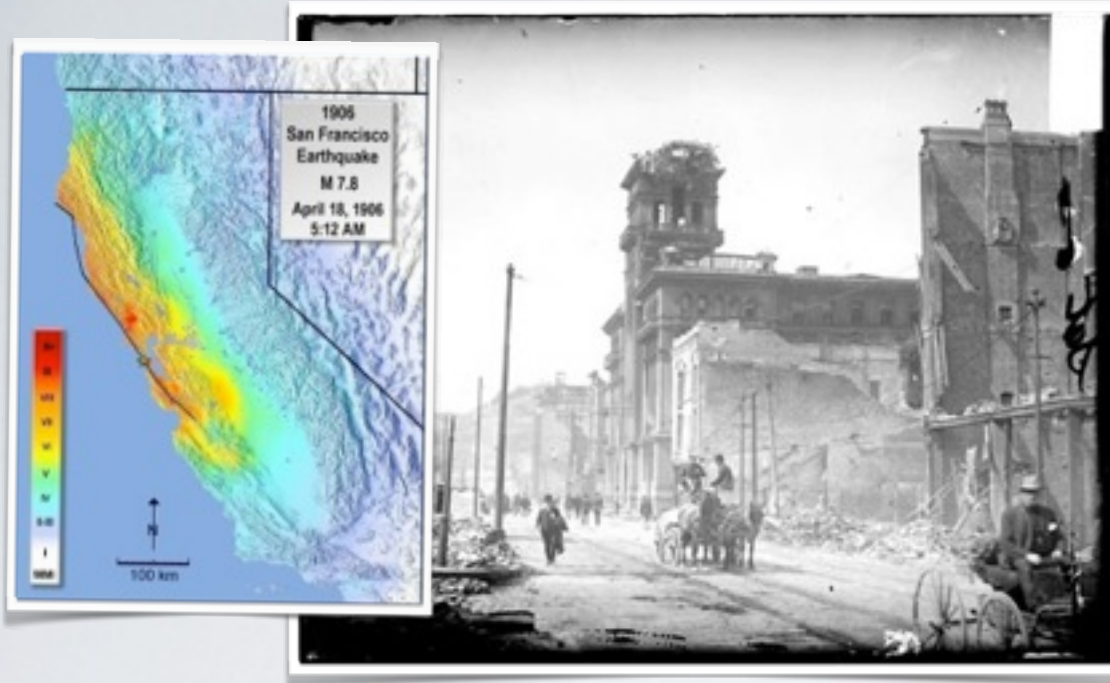
- terrorism occurs mostly in global “hot spots”
- deaths follow a power law
- big events often in developed countries
- most terrorist groups short-lived



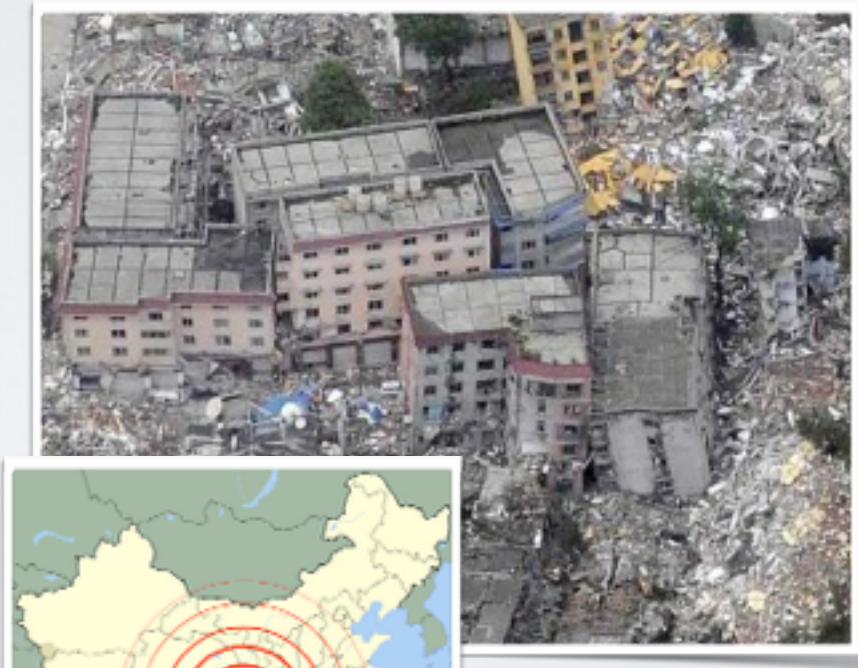
$$(\text{frequency}) \propto (\text{deaths})^{-\alpha}$$

what else follows power laws?

1906 San Francisco, M7.8



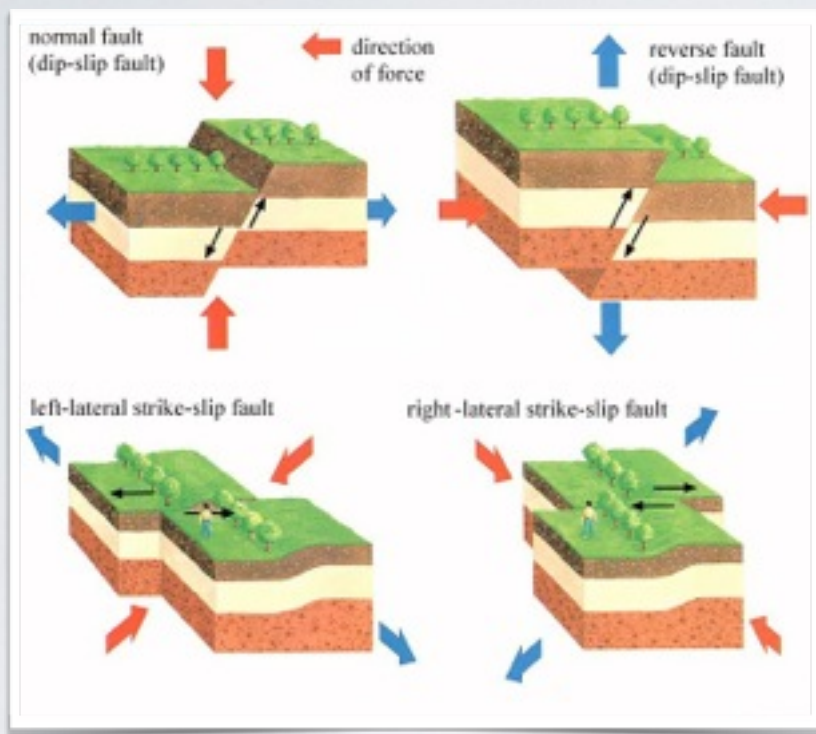
2008 Sichuan, M7.9



2011 Japan, M8.9

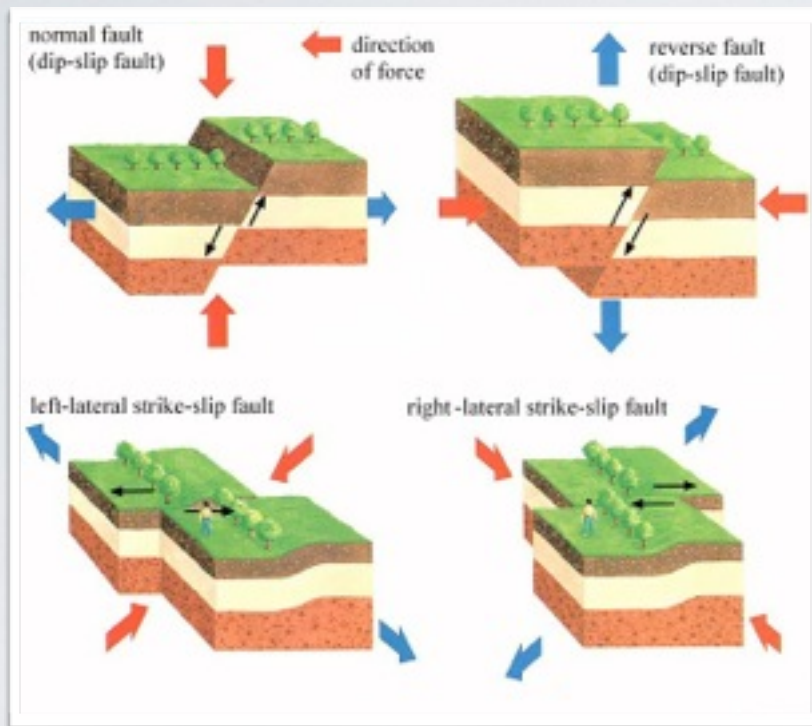


earthquakes

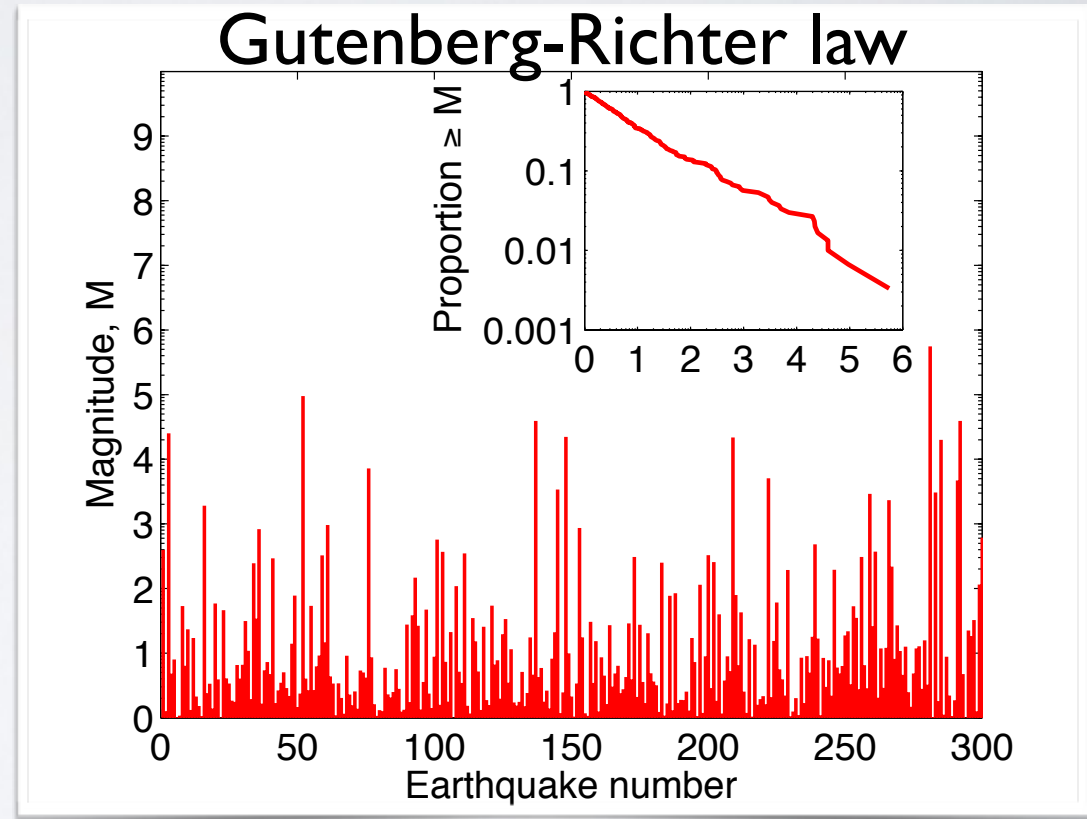


earthquake physics

earthquakes



earthquake physics



$$(\text{frequency}) \propto (\text{seismic moment})^{-\alpha}$$

earthquakes

Gutenberg-Richter law

$$F \propto M^{-\alpha}$$

physics largely *known*

processes *fixed*

forecasting possible
(years of successes)

prediction very hard
(years of failures)

terrorism

Richardson's law

$$F \propto S^{-\alpha}$$

processes largely *unknown*

processes *dynamic, adaptive*

how do we forecast?

what can we predict?
what can we not predict?

data, data, data

data, data, data

some of my projects published in 2013

- scoring dynamics in professional team sports
- scoring dynamics in the video game *Halo*
- identifying patterns in malaria gene networks
- detecting friendships in online social networks
- body size evolution of horses over the past 55 million years
- social networks in c1400 American Southwest
- forecasting large events in terrorism
- how large should whales be?
- ...

data, data, data

Scoring dynamics across professional team sports: tempo, balance and predictability

Sears Merritt^{1,*} and Aaron Clauset^{1,2,3,†}

¹*Department of Computer Science, University of Colorado, Boulder, CO 80309*

²*BioFrontiers Institute, University of Colorado, Boulder, CO 80303*

³*Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe, NM 87501*

sport	abbrv.	seasons	teams	competitions	scoring events
Football (college)	CFB	10, 2000–2009	486	14,588	120,827
Football (pro)	NFL	10, 2000–2009	31	2,654	19,476
Hockey (pro)	NHL	10, 2000–2009	29	11,813	44,989
Basketball (pro)	NBA	9, 2002–2010	31	11,744	1,080,285

data, data, data

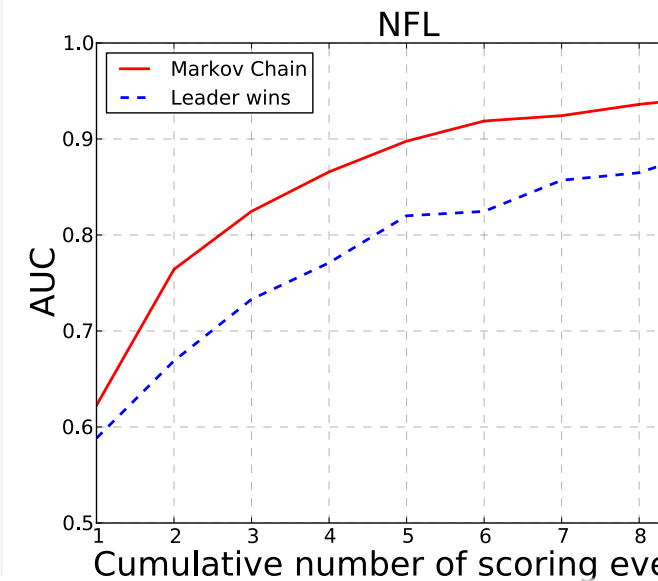
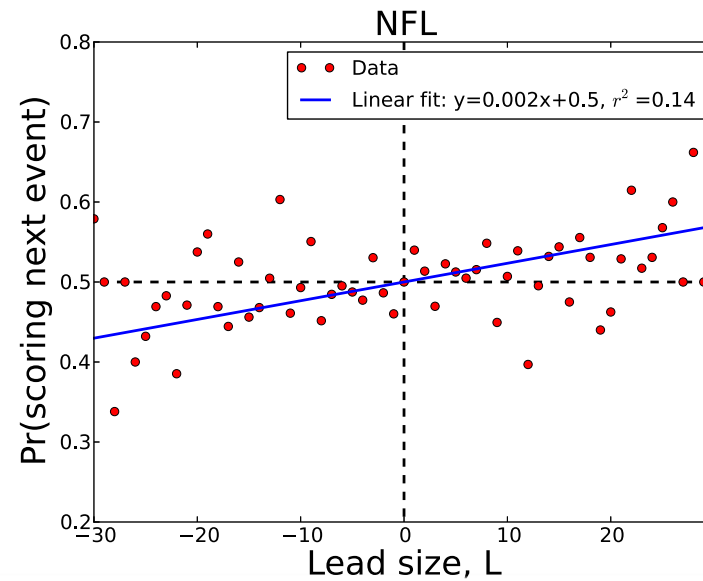
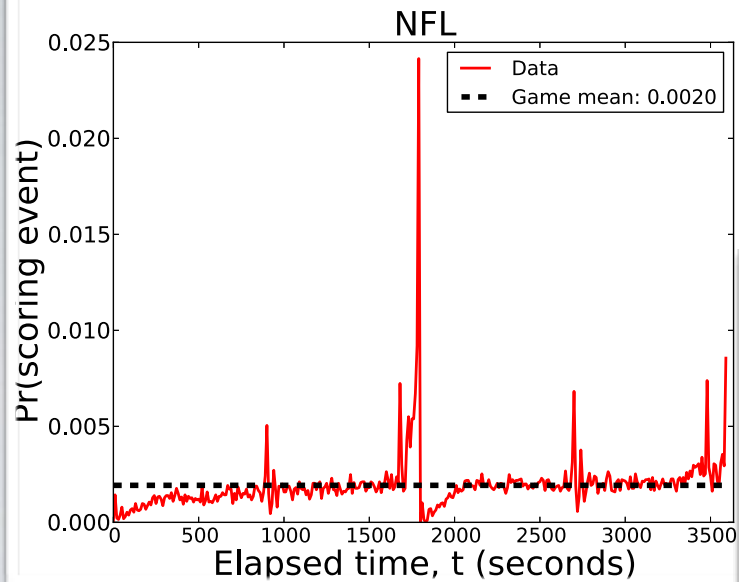
Scoring dynamics across professional team sports: tempo, balance and predictability

Sears Merritt^{1,*} and Aaron Clauset^{1,2,3,†}

¹Department of Computer Science, University of Colorado, Boulder, CO 80309

²BioFrontiers Institute, University of Colorado, Boulder, CO 80303

³Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe, NM 87501



fin

more

follow me on twitter: @aaronclauset

read my blog : www.structureandstrangeness.com

read my papers : www.santafe.edu/~aaronc/