



Big Data <-> Big Networks

Albert Díaz Guilera

- “Donat el caràcter i la finalitat exclusivament docent i eminentment il·lustrativa de les explicacions a classe d'aquesta presentació, l'autor s'acull a l'article 32 de la Llei de propietat intel·lectual vigent respecte de l'ús parcial d'obres alienes com ara imatges, gràfics o altre material contingudes en les diferents diapositives”
- “Dado el carácter y la finalidad exclusivamente docente y eminentemente ilustrativa de las explicaciones en clase de esta presentación, el autor se acoge al artículo 32 de la Ley de Propiedad Intelectual vigente respecto al uso parcial de obras ajenas como imágenes, gráficos u otro material contenidos en las diferentes diapositivas”.

1



UNIVERSITAT
BARCELONA

Complex Networks

Albert Díaz Guilera
<http://diaz-guilera.net>
@anduviera

C lab B complexity lab barcelona

 COMPLEXITAT

2

Big Data <-> Big Networks Albert Díaz Guilera

Macroscale

- Shortest-paths
- Clustering (of the network)
- Distributions (degree, betweenness, ...)
- Statistical properties: clustering
- Correlations

3

Big Data <-> Big Networks Albert Díaz Guilera

Distance between two nodes

- Number of links that make up the shortest-path between two nodes

- Centrality: nodes that are “close” to many other nodes in the network.
- Global centrality: defined as the sum of minimum distances to any other nodes in the networks

4

Big Data <-> Big Networks Albert Díaz-Guilera

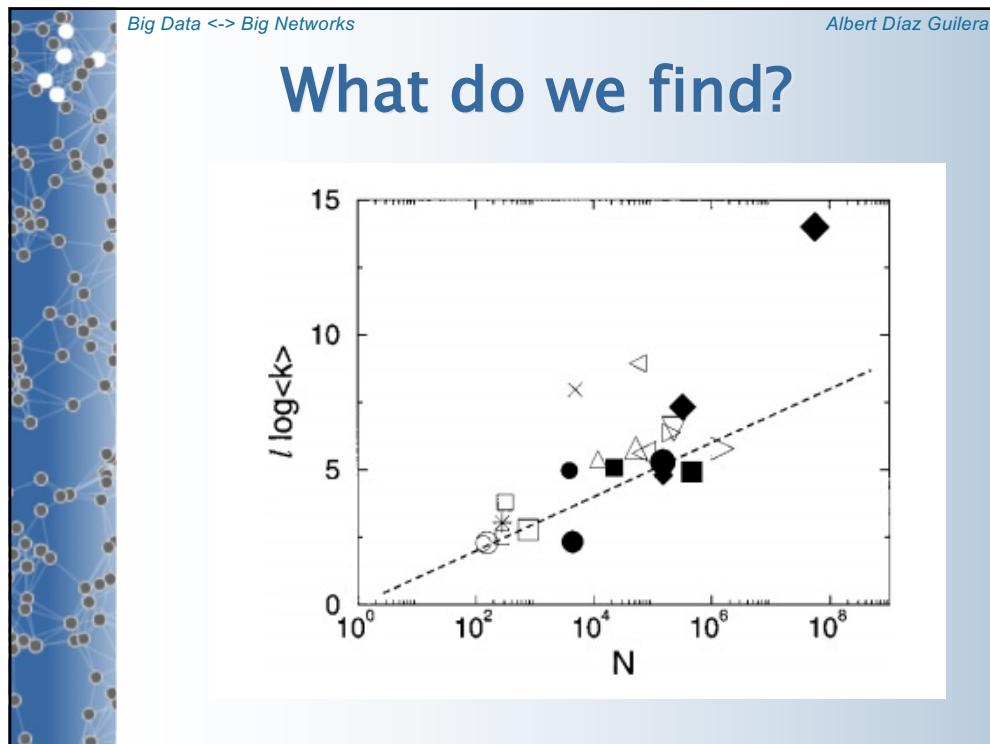
Global centrality of the whole network?

Mean shortest path = average over all pairs of nodes in the network

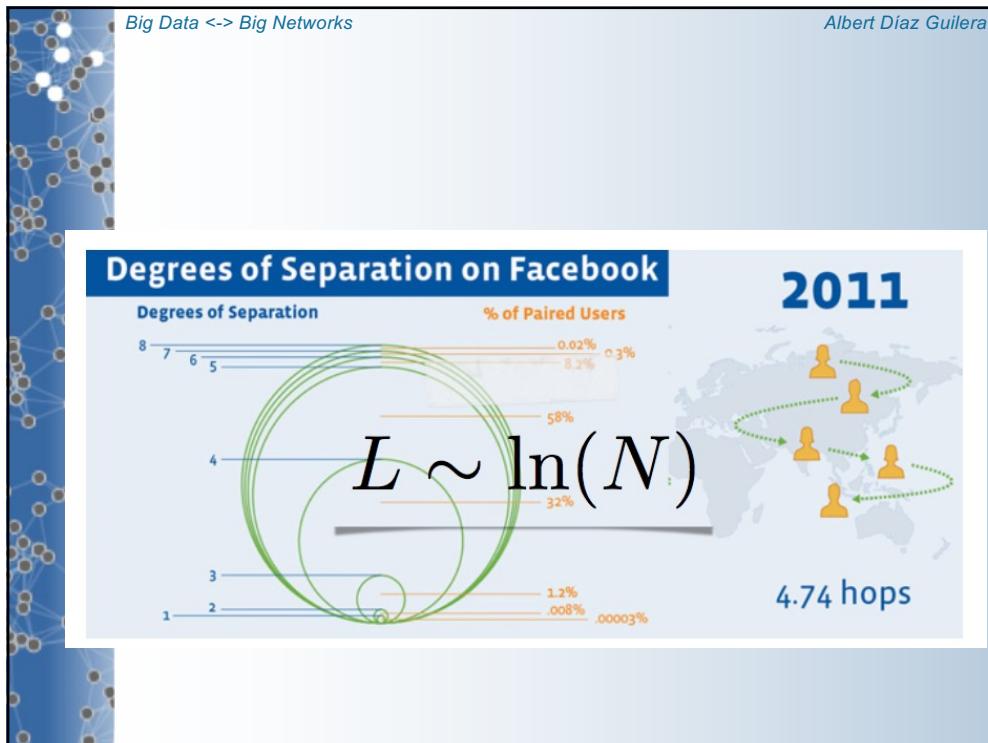
Diameter: largest distance between a pair of nodes in the network



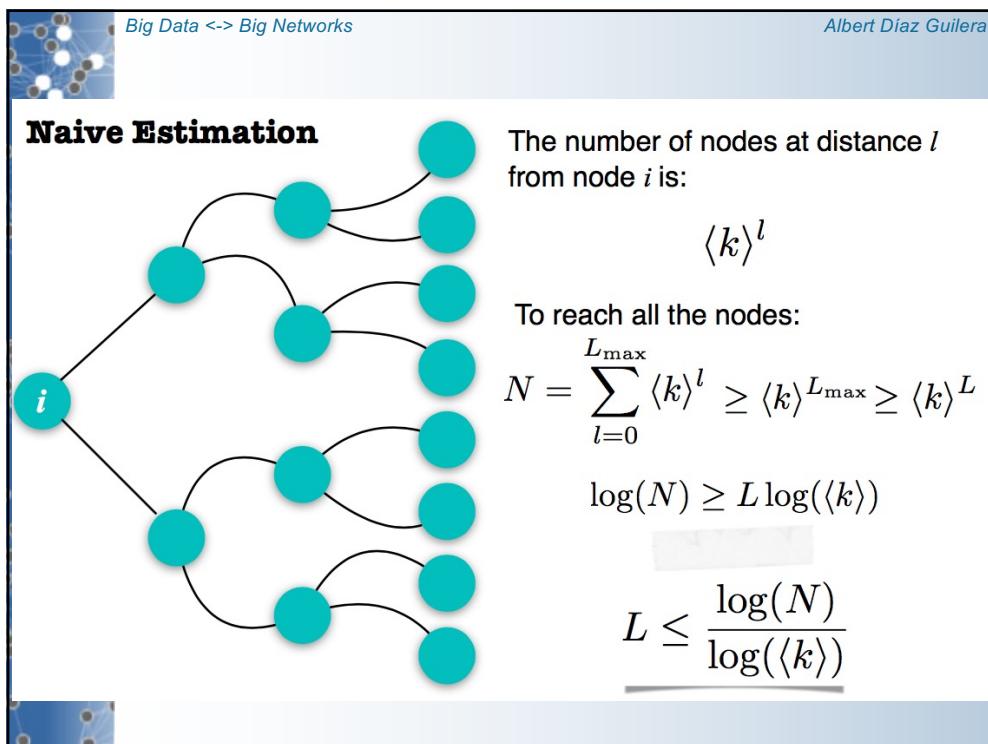
5



6



7



8



Distance Measures

Graph diameter, radius, eccentricity and other properties.

<code>center (G[, e])</code>	Return the center of the graph G.
<code>diameter (G[, e])</code>	Return the diameter of the graph G.
<code>eccentricity (G[, v, sp])</code>	Return the eccentricity of nodes in G.
<code>periphery (G[, e])</code>	Return the periphery of the graph G.
<code>radius (G[, e])</code>	Return the radius of the graph G.

9



center

Albert Díaz Guillera

`center (G, e=None)` [source]

Return the center of the graph G.

The center is the set of nodes with eccentricity equal to radius.

Parameters:

- `G (NetworkX graph)` – A graph
- `e (eccentricity dictionary, optional)` – A precomputed dictionary of eccentricities.

Returns:

`c` – List of nodes in center

Return type:

list

diameter

`diameter (G, e=None)` [source]

Return the diameter of the graph G.

The diameter is the maximum eccentricity.

Parameters:

- `G (NetworkX graph)` – A graph
- `e (eccentricity dictionary, optional)` – A precomputed dictionary of eccentricities.

Returns:

`d` – Diameter of graph

Return type:

integer

10

Big Data <-> Big Networks

Albert Díaz Guilera



eccentricity

`eccentricity (G, v=None, sp=None) [source]`

Return the eccentricity of nodes in G.

The eccentricity of a node v is the maximum distance from v to all other nodes in G.

Parameters:

- G (NetworkX graph) – A graph
- v (node, optional) – Return value of specified node
- sp (dict of dicts, optional) – All pairs shortest path lengths as a dictionary of dictionaries

Returns: ecc – A dictionary of eccentricity values keyed by node.

Return type: dictionary

periphery

`periphery (G, e=None) [source]`

Return the periphery of the graph G.

The periphery is the set of nodes with eccentricity equal to the diameter.

Parameters:

- G (NetworkX graph) – A graph
- e (eccentricity dictionary, optional) – A precomputed dictionary of eccentricities.

Returns: p – List of nodes in periphery

Return type: list

11

Big Data <-> Big Networks

Albert Díaz Guilera



radius

`radius (G, e=None) [source]`

Return the radius of the graph G.

The radius is the minimum eccentricity.

Parameters:

- G (NetworkX graph) – A graph
- e (eccentricity dictionary, optional) – A precomputed dictionary of eccentricities.

Returns: r – Radius of graph

Return type: integer

12

Big Data <-> Big Networks *Albert Díaz Guilera*



dijkstra_path

`dijkstra_path (G, source, target, weight='weight')` [\[source\]](#)

Returns the shortest path from source to target in a weighted graph G.

Parameters:

- `G (NetworkX graph)` –
- `source (node)` – Starting node
- `target (node)` – Ending node
- `weight (string, optional (default='weight'))` – Edge data key corresponding to the edge weight

Returns: `path` – List of nodes in a shortest path.

Return type: `list`

14

Big Data <-> Big Networks *Albert Díaz Guilera*



Clustering

- Cycles in social network analysis language
- Circles of friends in which every member knows each other

15

Big Data <-> Big Networks Albert Díaz-Guilera

Clustering of a node:

$$C_i = \frac{\text{\#triangles connected to } i}{\text{\#possible triangles connected to } i} = \frac{2E_i}{k_i(k_i - 1)}$$

Clustering of the Network:

$$C = \frac{1}{N} \sum_{i=1}^N C_i$$



- Alternative definition: ratio between total number of triangles and possible

16

Big Data <-> Big Networks Albert Díaz-Guilera

What happens in real networks?

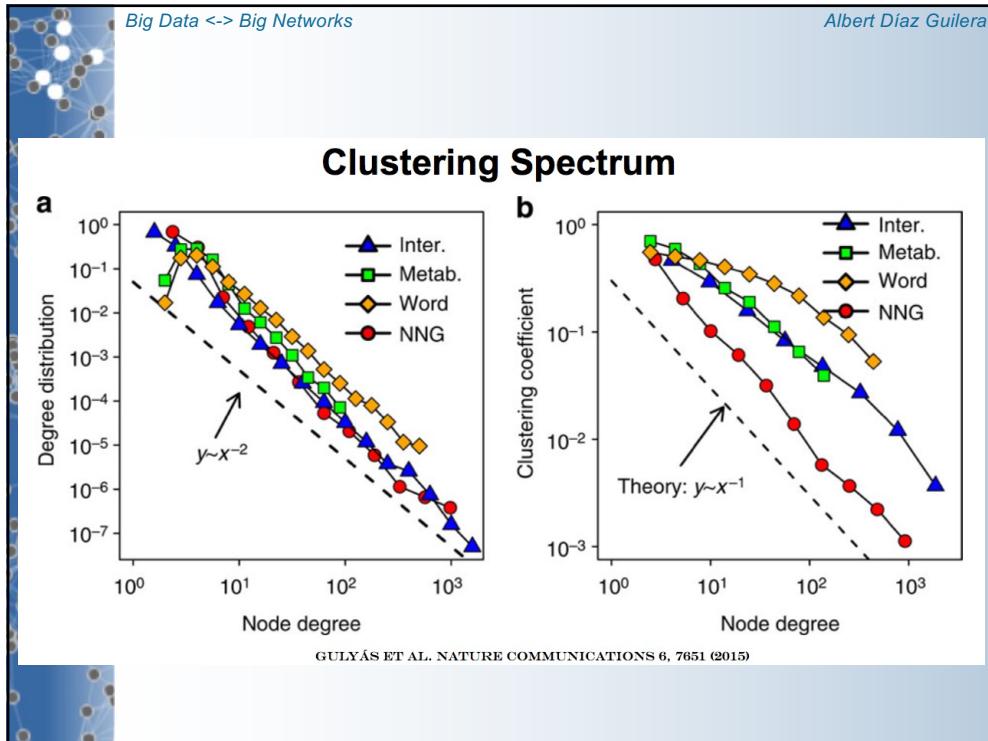



$C/k \propto N^{-1}$

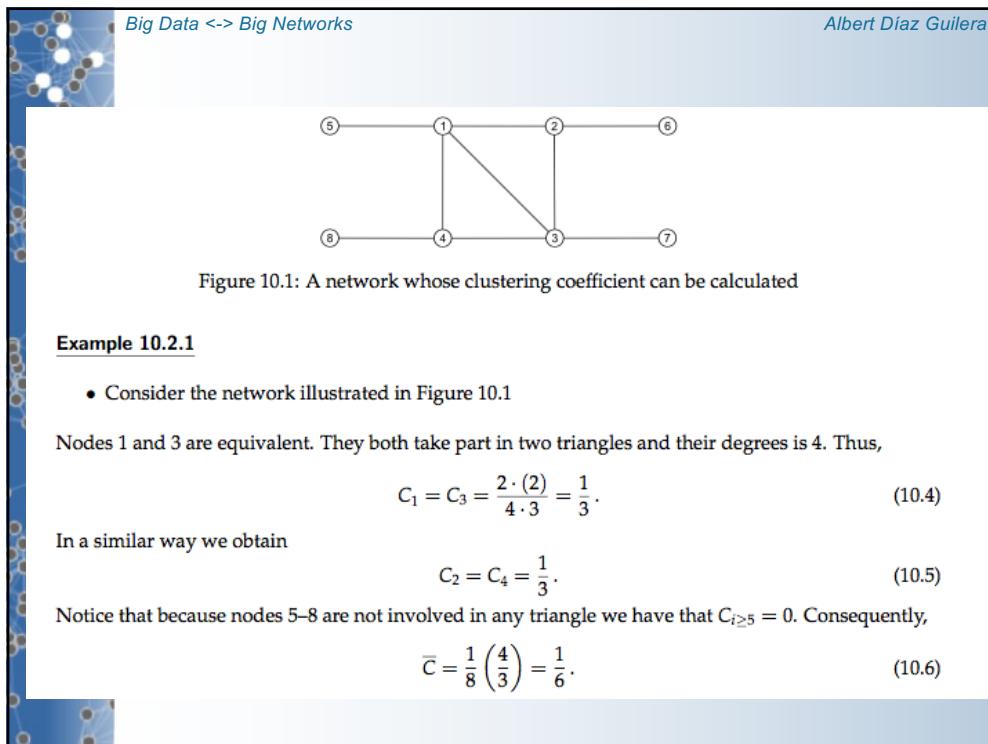
$C/k \propto N^{-1}$

- The clustering coefficient is much larger than it is in an equivalent random network

17



18



19

Big Data <-> Big Networks Albert Díaz Guilera

10.3 The Newman Clustering Coefficient

Another way of quantifying the global clustering of a network is by means of the Newman clustering coefficient, also known as the **transitivity index** of the network. Let $t = |C_3|$ be the total number of triangles, and let $|P_2|$ be the number of paths of length 2 in the network (representing all potential three-way relationships). Then,

$$C = \frac{3t}{|P_2|} = \frac{3|C_3|}{|P_2|}. \quad (10.7)$$

- Consider again the network illustrated in Figure 10.1. We can obtain the number of triangles in that network by using the spectral properties of the adjacency matrix. That is,

$$t = \frac{1}{6} \text{tr}(A^3) = 2. \quad (10.8)$$

The number of paths of length 2 in the network can be obtained using the following formula (which we will justify later).

$$|P_2| = \sum_{i=1}^n \binom{k_i}{2} = \sum_{i=1}^n \frac{k_i(k_i - 1)}{2} = 18. \quad (10.9)$$

Thus,

$$C = \frac{3 \times 2}{18} = \frac{1}{3}. \quad (10.10)$$

20

Big Data <-> Big Networks Albert Díaz Guilera

In general, the Watts–Strogatz index quantifies how clustered a network is locally, while Newman one indicates how clustered the network is as a whole. In general there is a good correlation between both indices for real-world networks as illustrated in Figure 10.4.

C (x-axis)	C̄ (y-axis)
0.05	0.05
0.08	0.12
0.10	0.15
0.12	0.20
0.15	0.25
0.18	0.30
0.20	0.35
0.22	0.40
0.25	0.45
0.28	0.50
0.30	0.55
0.35	0.60
0.40	0.65
0.45	0.70
0.50	0.75
0.55	0.80
0.60	0.85

Figure 10.4: Correlation between Watts–Strogatz (\bar{C}) and Newman (C) clustering coefficients for 20 real-world networks.

21

Big Data <-> Big Networks Albert Díaz Guilera



Clustering

Algorithms to characterize the number of triangles in a graph.

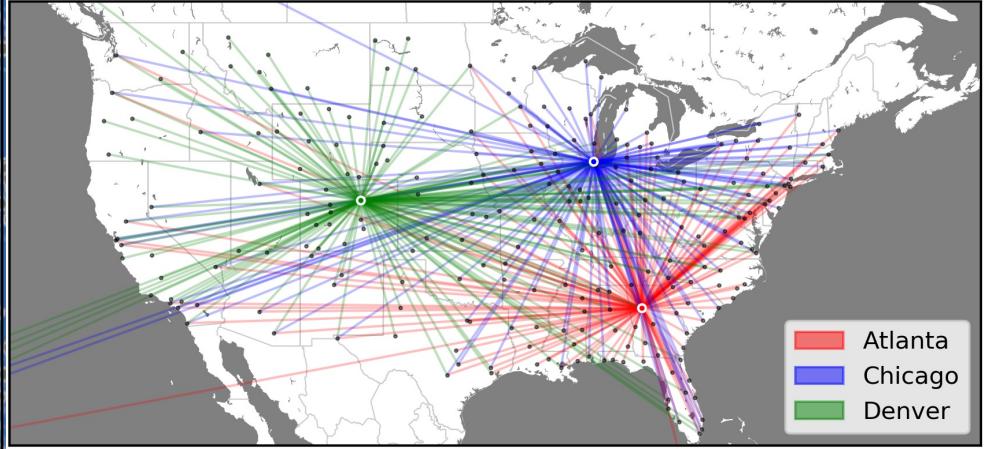
<code>triangles (G[, nodes])</code>	Compute the number of triangles.
<code>transitivity (G)</code>	Compute graph transitivity, the fraction of all possible triangles present in G.
<code>clustering (G[, nodes, weight])</code>	Compute the clustering coefficient for nodes.
<code>average_clustering (G[, nodes, weight, ...])</code>	Compute the average clustering coefficient for the graph G.
<code>square_clustering (G[, nodes])</code>	Compute the squares clustering coefficient for nodes.



23

Big Data <-> Big Networks Albert Díaz Guilera

Degree: real networks are heterogeneous



Atlanta Chicago Denver

Some nodes are more important than others



24

Big Data <-> Big Networks

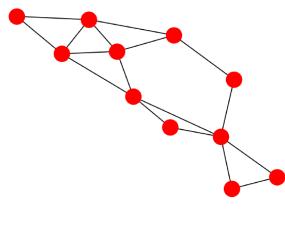
Albert Díaz-Guilera

Degree distribution (macroscopic scale)

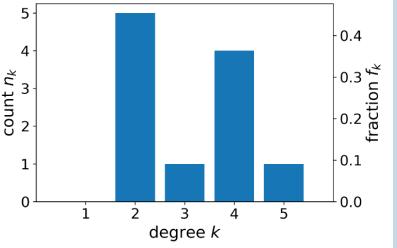
- Gives an idea of the spread in the number of links the nodes have
- $P(k)$ is the probability that a randomly selected node has k links
- Statistical sense

25

Degree distributions



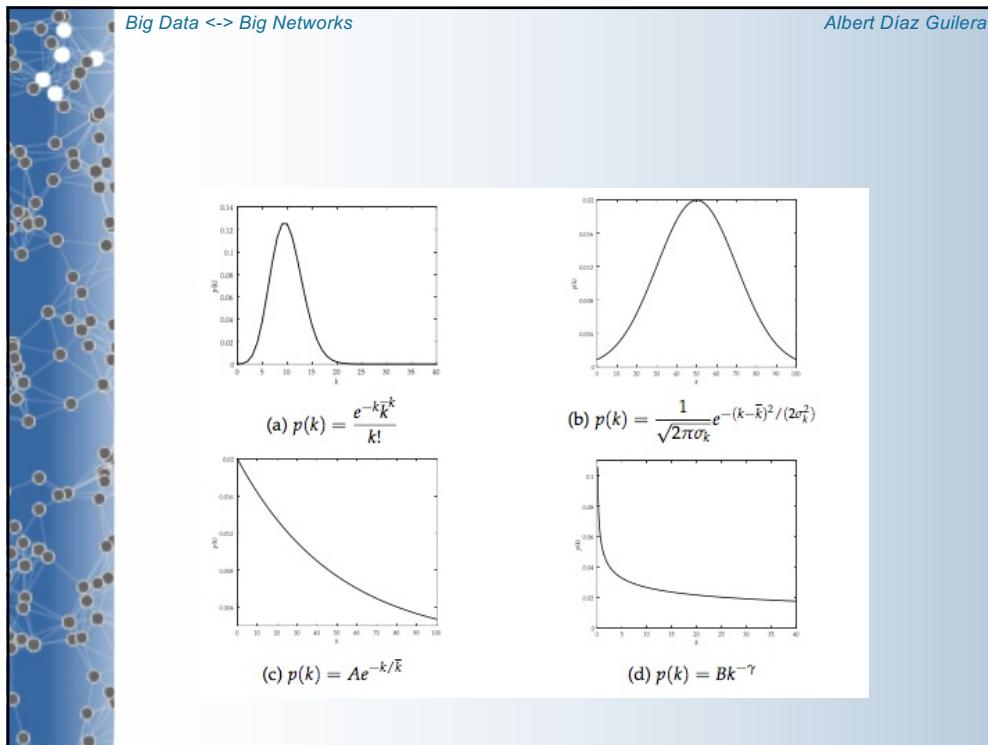
Histogram



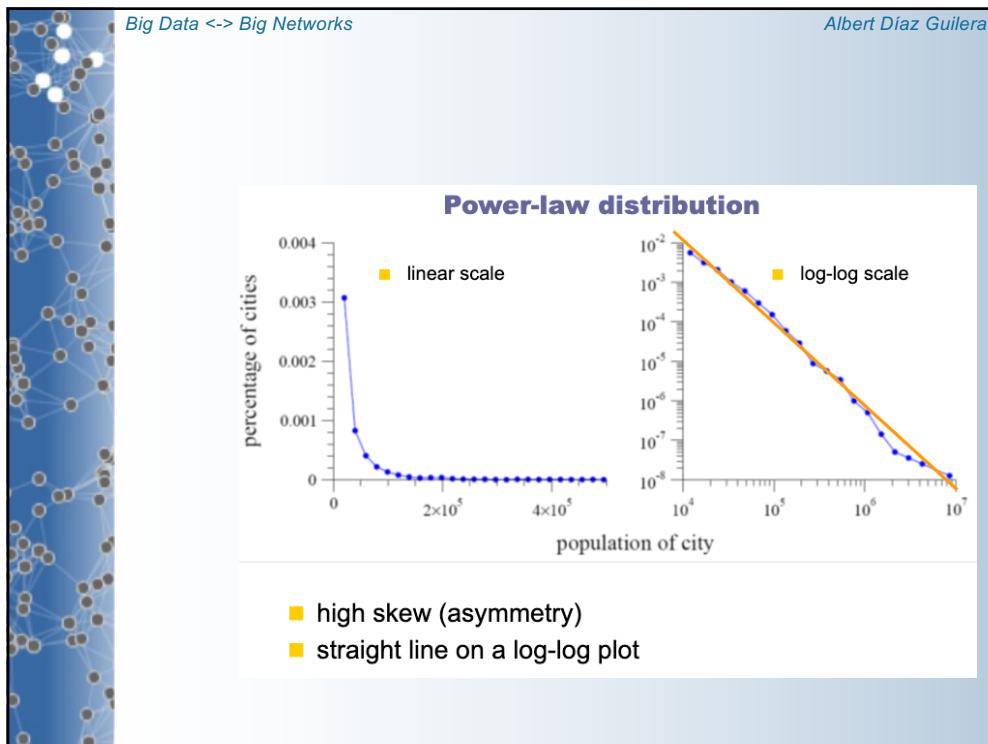
degree k	count n_k	frequency f_k
1	0	0.0
2	5	0.4
3	1	0.1
4	4	0.3
5	1	0.1

- n_k = number of nodes with degree k
- $f_k = \frac{n_k}{N}$ = frequency of degree k

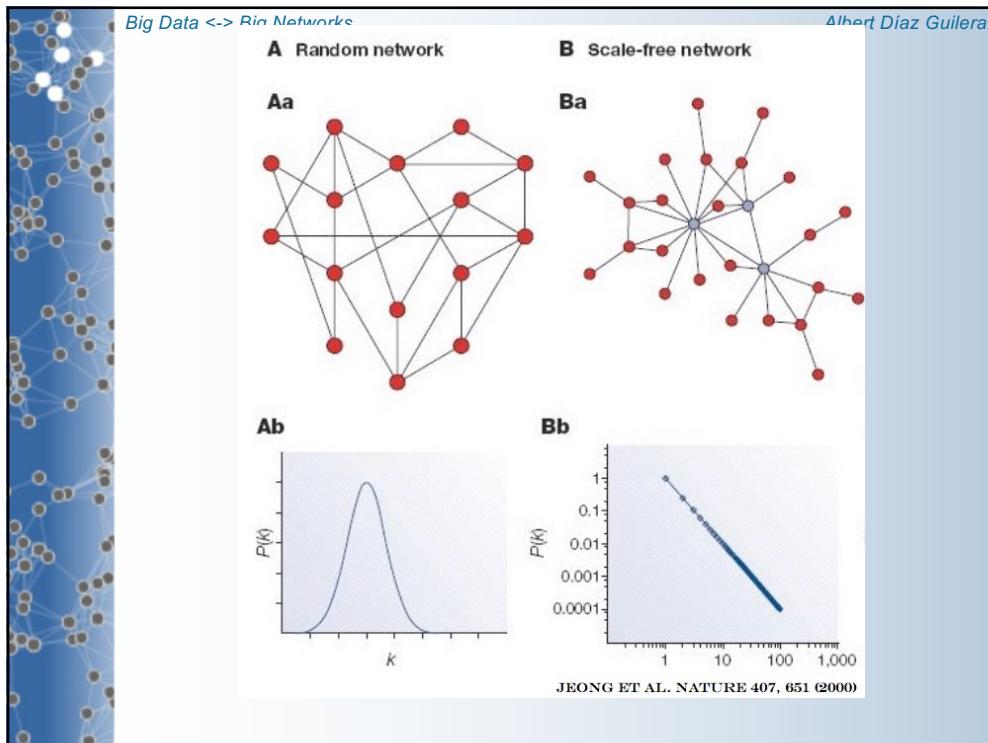
26



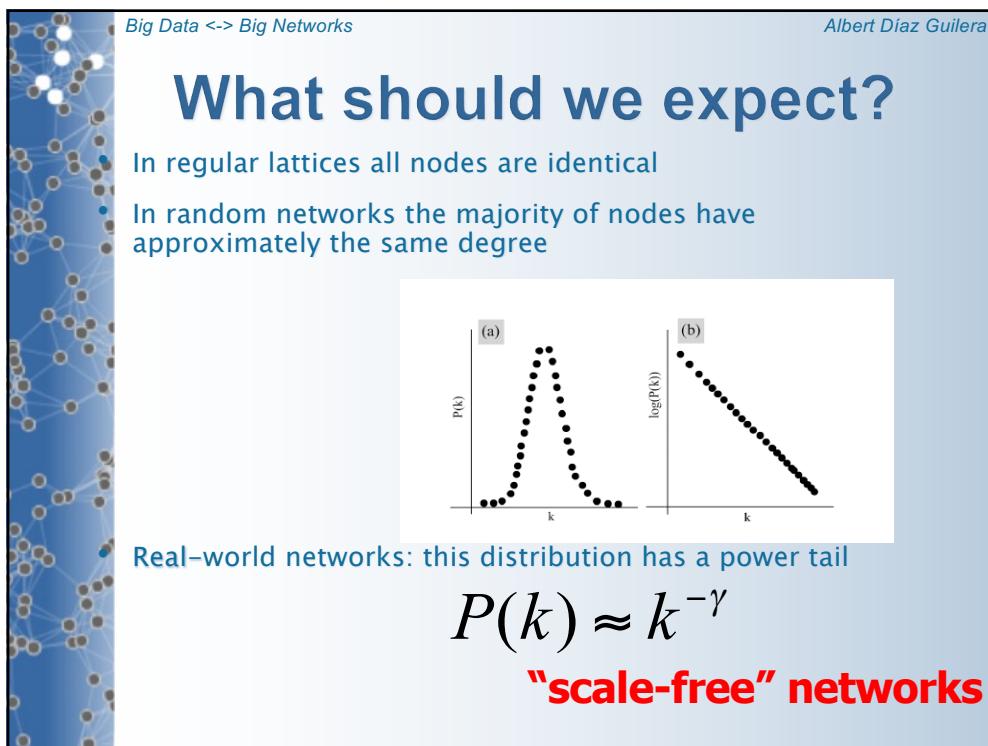
27



28



29

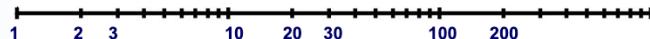


30



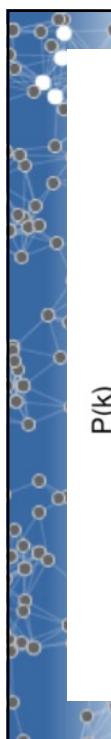
Logarithmic scale

- **Question:** how to plot a probability distribution if the variable spans a large range of values, from small to (very) large?
- **Answer:** use the **logarithmic scale**
- **How to do it:** report the logarithms of the values on the x- and y-axes
 - powers of a number will be uniformly spaced



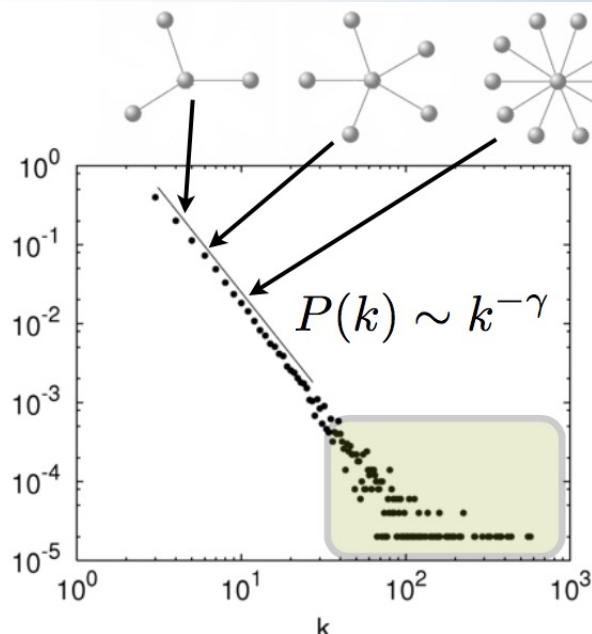
■ $2^0=1, 2^1=2, 2^2=4, 2^3=8, 2^4=16, 2^5=32, 2^6=64, \dots$

31

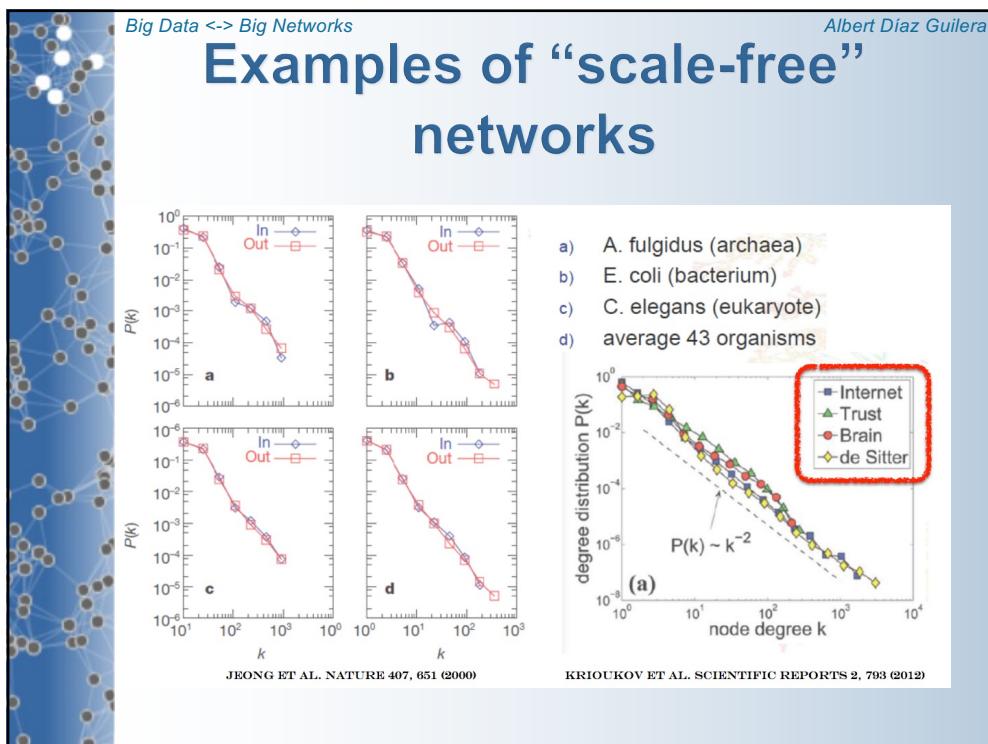


Big Data <-> Big Networks

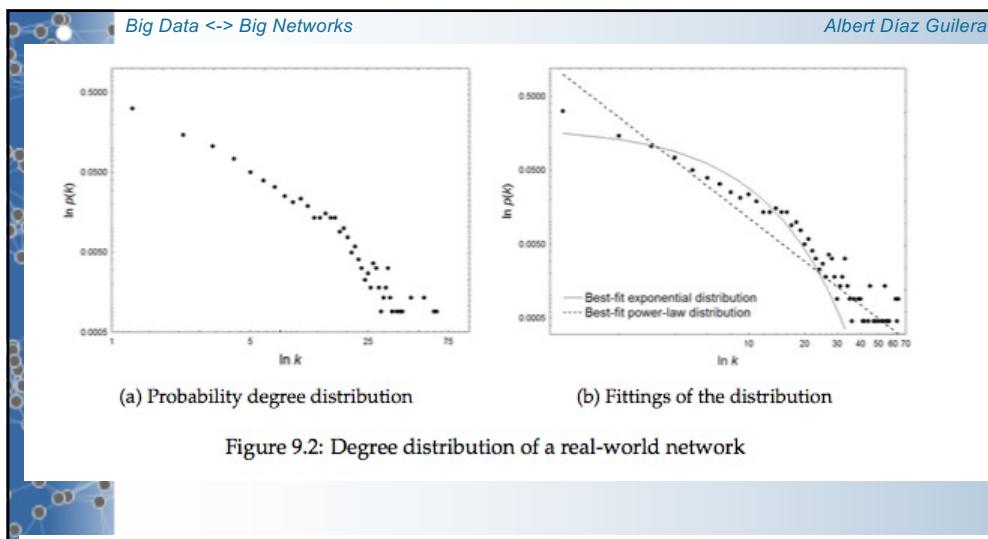
Albert Díaz-Guilera



32

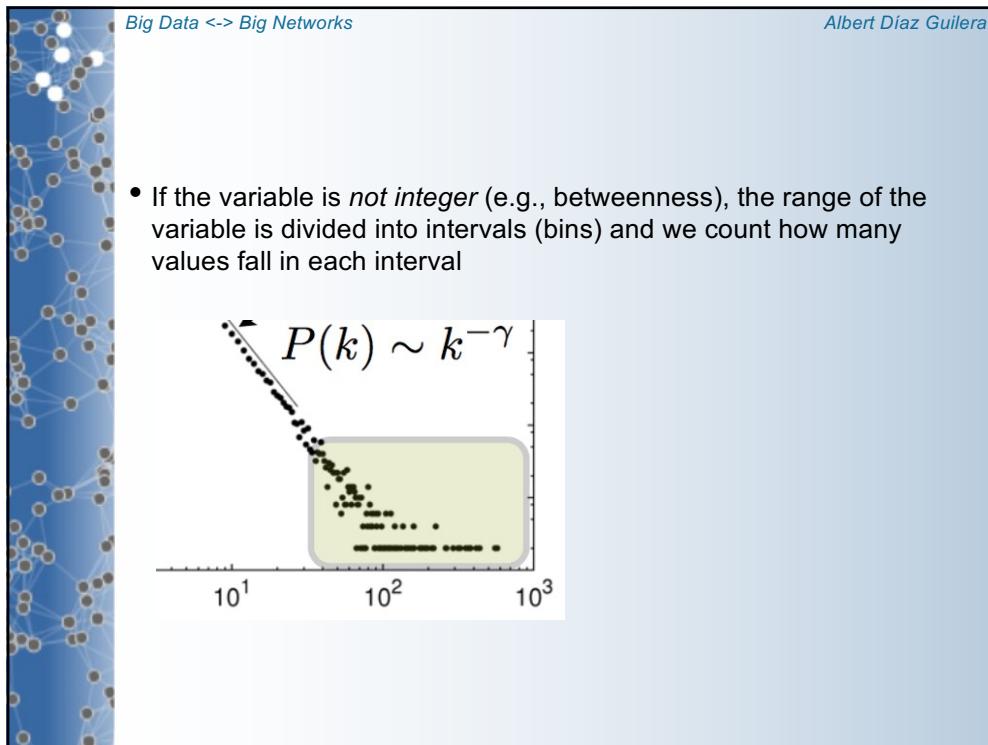


33

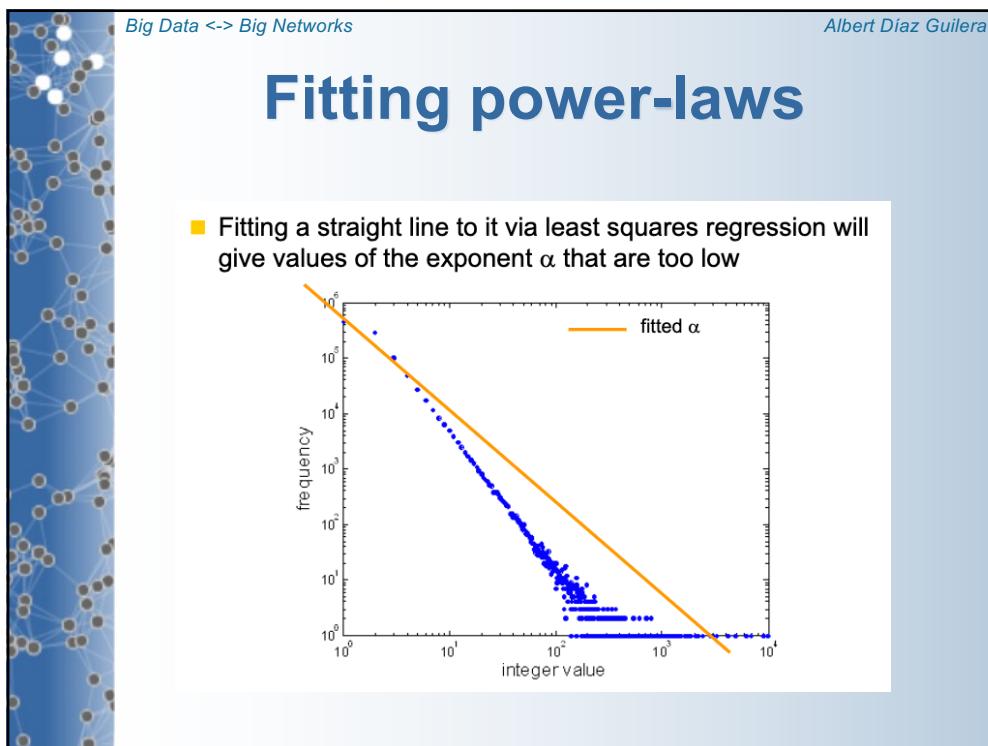


The degree distributions of real-world networks, however, do not look so smooth as the ones illustrated in Figure 9.1 and there are a number of difficulties in fitting the best model to describe the experimental data. For example, the relationship between $p(k)$ and k can be erratic (see Figure 9.2a); there may not be sufficient data to fit a statistically significant model; or there may be many statistical distributions to which the same dataset can be fitted (see Figure 9.2b).

34



35



36

Big Data <-> Big Networks Albert Díaz Guilera

Logarithmic binning

- Most common and not very accurate method:
 - Bin the different values of x and create a frequency histogram

$\ln(\# \text{ of times } x \text{ occurred})$

$\ln(x)$

$\ln(x)$ is the natural logarithm of x , but any other base of the logarithm will give the same exponent of a because $\log_{10}(x) = \ln(x)/\ln(10)$

37

Big Data <-> Big Networks Albert Díaz Guilera

Logarithmic binning

- bin data into exponentially wider bins:
 - 1, 2, 4, 8, 16, 32, ...
- normalize by the width of the bin

evenly spaced datapoints

$\alpha = 2.41$

less noise in the tail of the distribution

38

Big Data <-> Big Networks Albert Díaz Guilera

Cumulative distributions

The second is to consider the cumulative distribution function (CDF) defined as

$$P'(k) = \sum_{k'=1}^k p(k'). \quad (9.1)$$

This represents the probability of choosing at random a node with degree smaller than or equal to k . Common CDFs for degree distributions are

Poisson :	$P'(k) = e^{-\bar{k}} \sum_{i=1}^{[k]} \frac{\bar{k}^i}{i!}$
Exponential :	$P'(k) = 1 - e^{-\bar{k}/k}$.
Powerlaw :	$P'(k) = 1 - k^{-\gamma+1}$.

39

Big Data <-> Big Networks Albert Díaz Guilera

Cumulative distributions

- fitted exponent (2.43) much closer to actual (2.5)

A log-log plot showing the frequency sample (y-axis, ranging from 10^0 to 10^6) versus v (x-axis, ranging from 10^0 to 10^4). The plot contains blue square data points and a red line representing a fit with slope $\alpha-1 = 1.43$. The data points follow a power-law-like decay.

40

Big Data <-> Big Networks

Albert Díaz Guilera

Where to start fitting?

- Some data exhibit a power law only in the tail
- After binning or taking the cumulative distribution you can fit to the tail

41

Big Data <-> Big Networks

Albert Díaz Guilera

Example:

- Distribution of citations to papers
- power law is evident only in the tail
 - $x_{\min} > 100$ citations

Source: MEJ Newman, 'Power laws, Pareto distributions and Zipf's law'

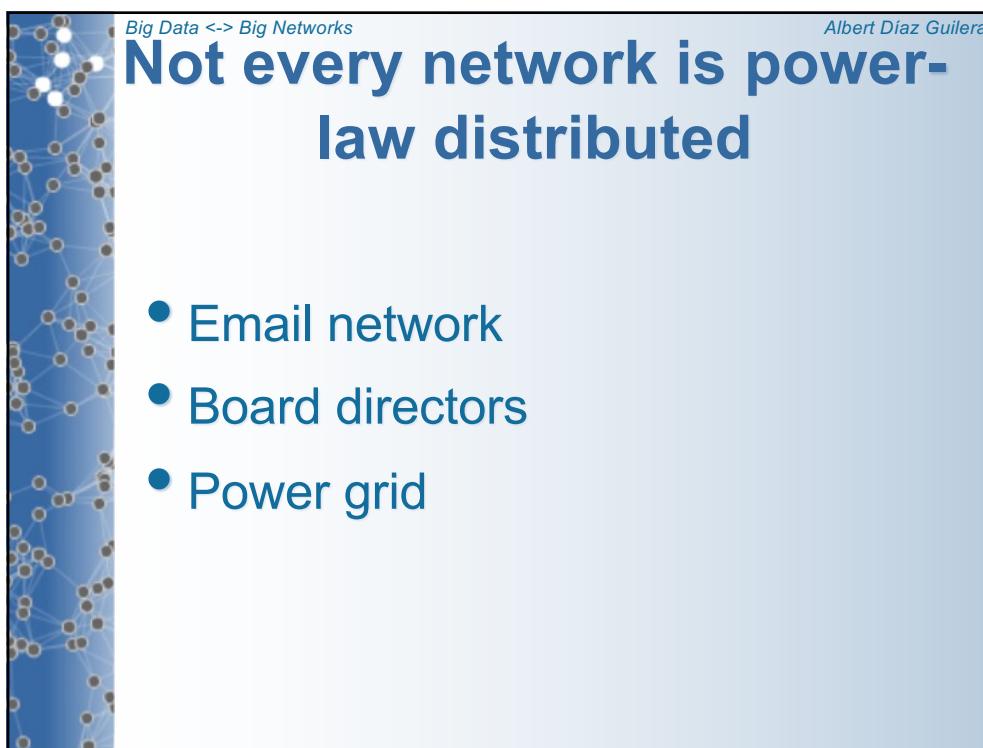
42

Big Data <-> Big Networks Albert Díaz Guilera

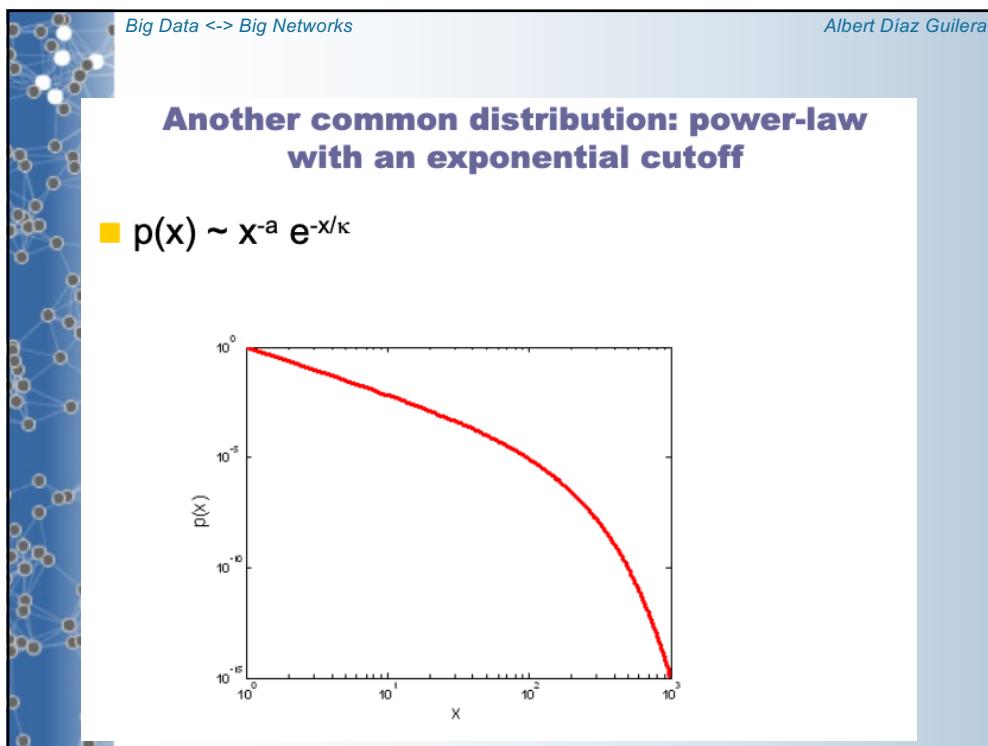


Many real world networks are power law	
	exponent α (in/out degree)
film actors	2.3
telephone call graph	2.1
email networks	1.5/2.0
sexual contacts	3.2
WWW	2.3/2.7
internet	2.5
peer-to-peer	2.1
metabolic network	2.2
protein interactions	2.4

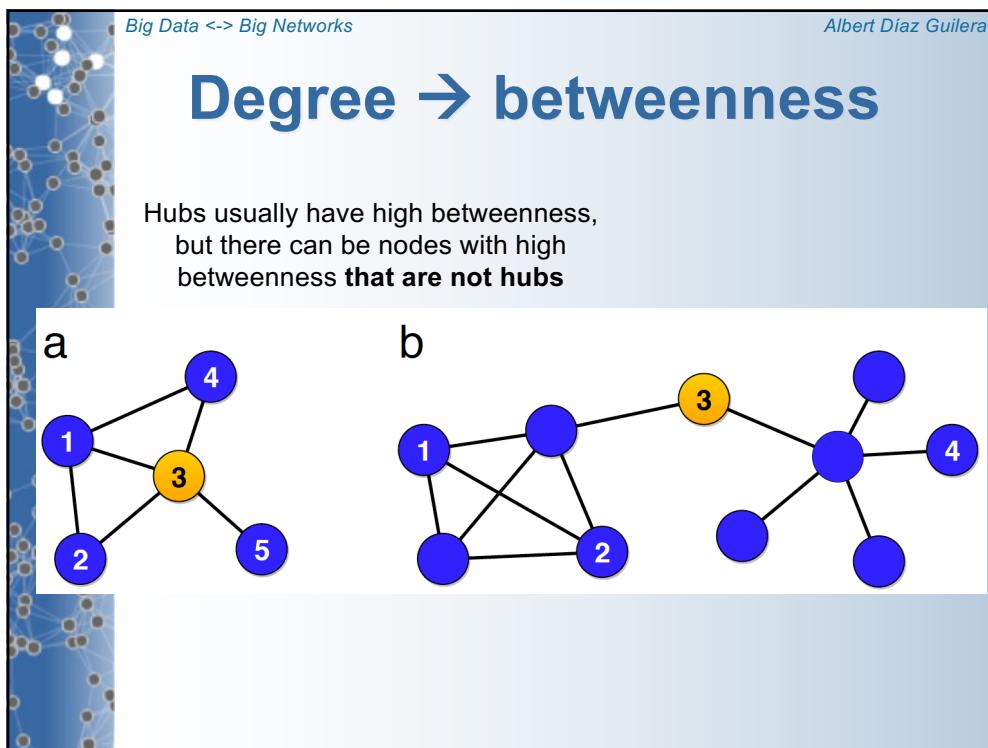
43



44



45



46

Big Data <-> Big Networks

Albert Díaz Guilera

Universality

- We find the same macroscopic behavior in systems which have completely different origins
- Engineered or self-organized
- Biology, transportation, social, ...
- Which are the basic underlying principles?
- Models (simple models) that can explain the basic trends

47

Big Data <-> Big Networks

Albert Díaz Guilera

Robustness

OPEN ACCESS Freely available online

PLOS ONE

Attack Robustness and Centrality of Complex Networks

Swami Iyer¹, Timothy Killingback^{2*}, Bala Sundaram³, Zhen Wang⁴

¹ Computer Science Department, University of Massachusetts, Boston, Massachusetts, United States of America, ² Mathematics Department, University of Massachusetts, Boston, Massachusetts, United States of America, ³ Physics Department, University of Massachusetts, Boston, Massachusetts, United States of America, ⁴ Physics Department, University of Massachusetts, Boston, Massachusetts, United States of America

Scale-free network

Random Graph

Simultaneous

Sequential

48

Big Data <-> Big Networks

Albert Díaz Guilera



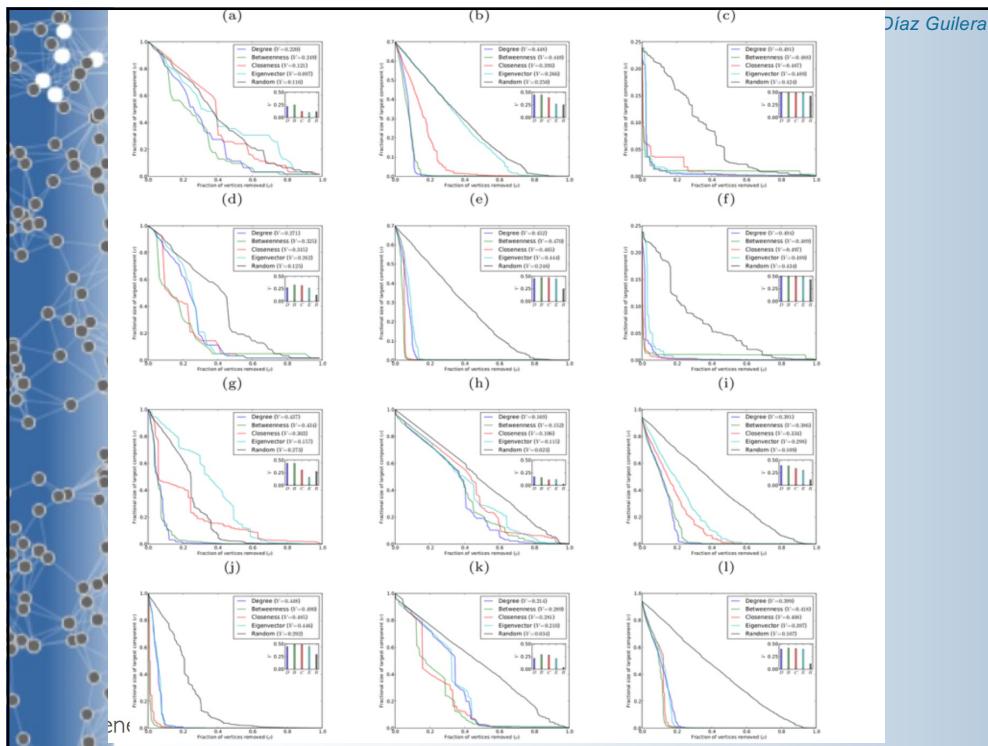
Once a suitable centrality measure has been fixed, we can compute σ as a function of ρ for removing vertices in decreasing order of that centrality measure. The robustness of a network under this type of vertex removal can be quantified by the *R-index*, which is defined by [13]

$$R = \frac{1}{N} \sum_{i=1}^N \sigma(i/N).$$

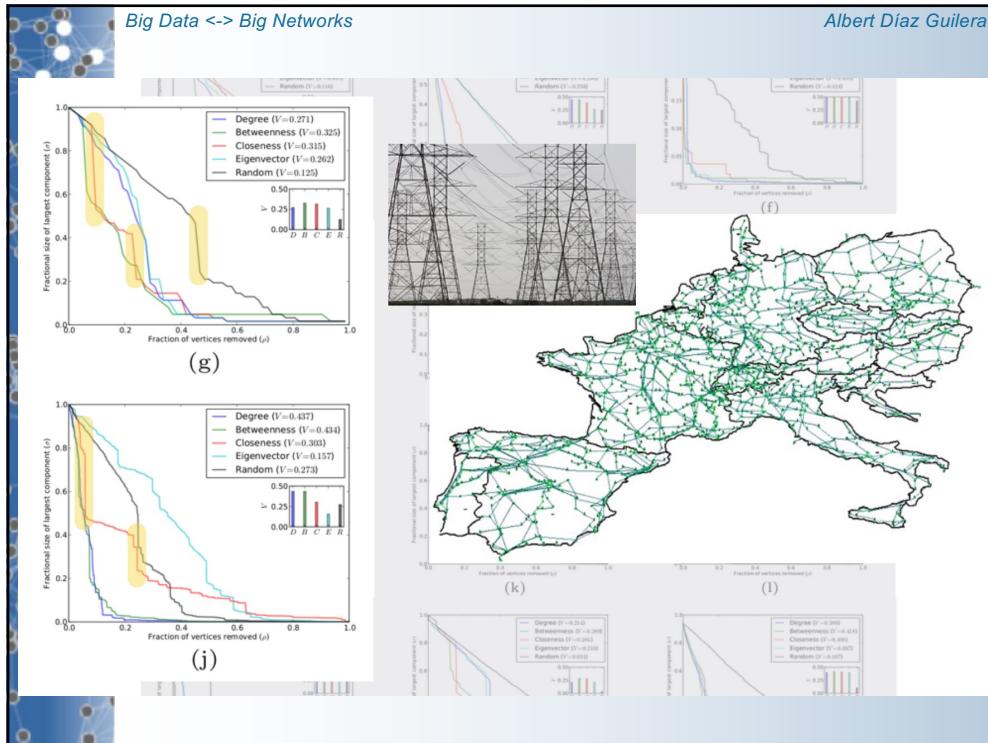
The normalization factor $1/N$ allows the robustness of networks of different sizes to be compared [13]. It is straightforward to show that for any scheme of removing vertices from any network, R attains its minimum value of $1/N$ on the star graph and its maximum value of $\frac{1}{2}(1 - 1/N)$ on the complete graph. Thus, for any network and method of vertex removal, $R \in [0, \frac{1}{2}]$. Consequently, we define the *V-index* V , which measures the *vulnerability* of a network to a given scheme of vertex removal, to be the complementary quantity to R ,

$$V = \frac{1}{2} - R.$$

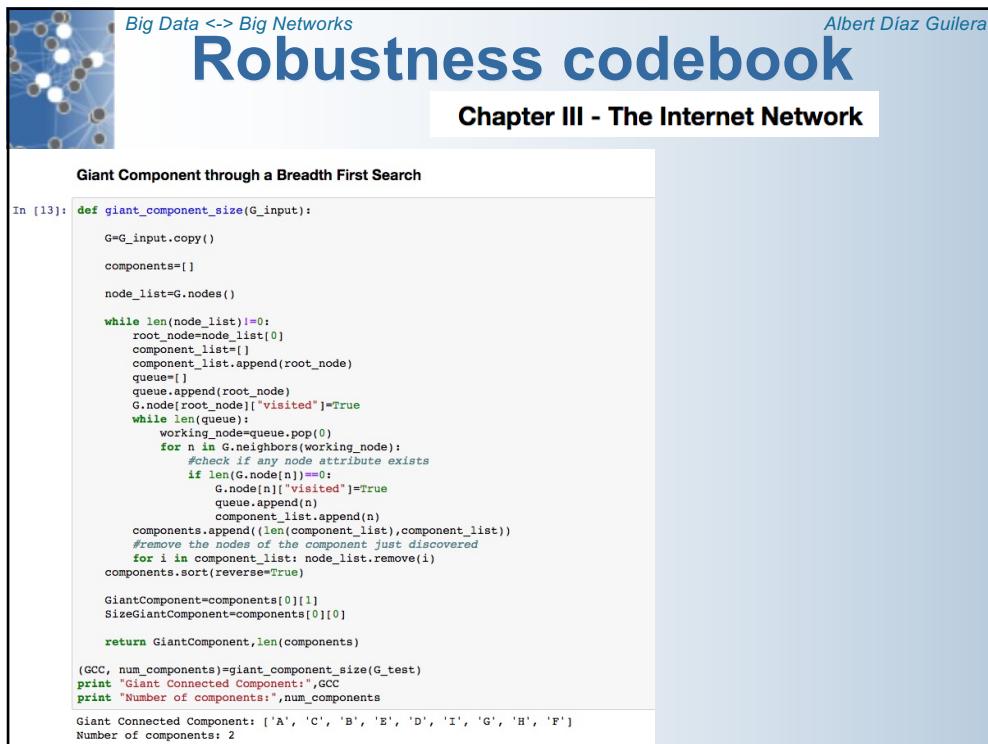
49



50



51



52

Big Data <-> Big Networks Albert Díaz Guilera

Correlations

- Degree correlations: expected degree of the neighbors of a node as a function of its degree

$$k_{\text{nn}}(k) = \sum_{k'} k' P(k'|k)$$

53

Big Data <-> Big Networks Albert Díaz Guilera

The correct mathematical way to quantify such a measure is the *conditioned probability* $p(k_1|k_2)$ to have a vertex with degree k_1 at one side of the edge when at the other site of the edge the degree is k_2 .

We have two constraints on the conditioned probability. The first one is given by normalization condition

$$\sum_{k_1} p(k_1|k_2) = 1. \quad (1.10)$$

For non oriented graphs the same quantity obeys the detailed balance distribution (Boguñá and Pastor-Satorras, 2002)

$$k_2 p(k_1|k_2) P(k_2) = k_1 p(k_2|k_1) P(k_1) \quad (1.11)$$

This balance equation simply states that the number of edges going from vertex k_1 to vertex k_2 must be equal to the number of edges going from vertex k_2 to vertex k_1 .

54

Big Data <-> Big Networks *Albert Díaz Guilera*



$P(k', k)$: Probability that two nodes of degree k and k' are linked

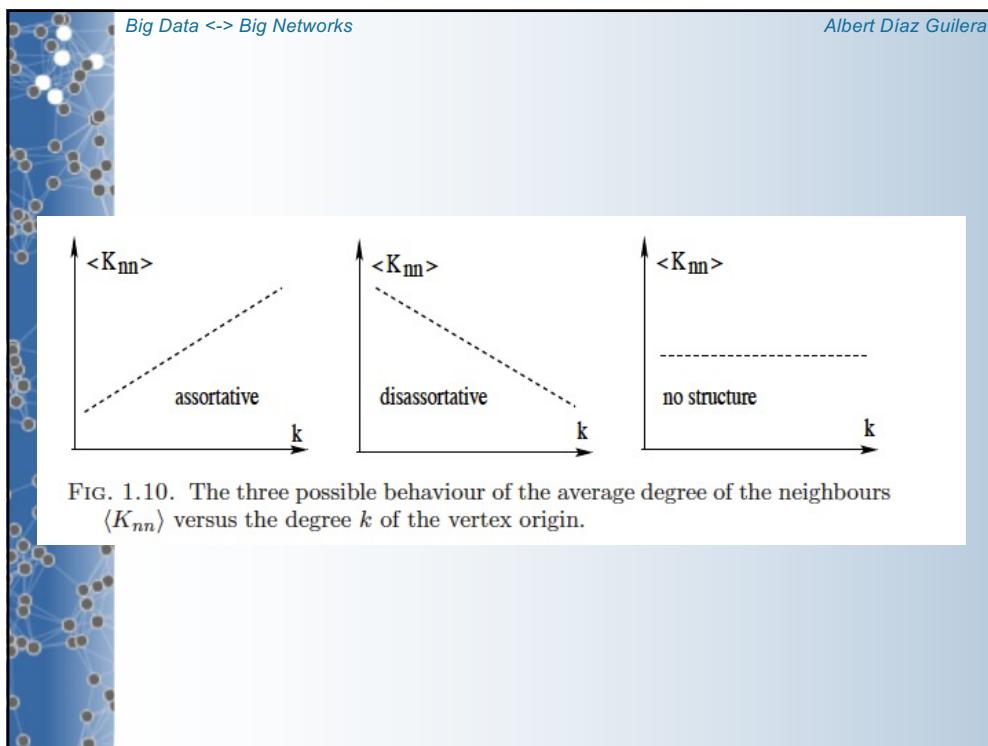
Detailed balance Equation for Networks

$$P(k', k) = kP(k)P(k'|k) = k'P(k')P(k|k')$$

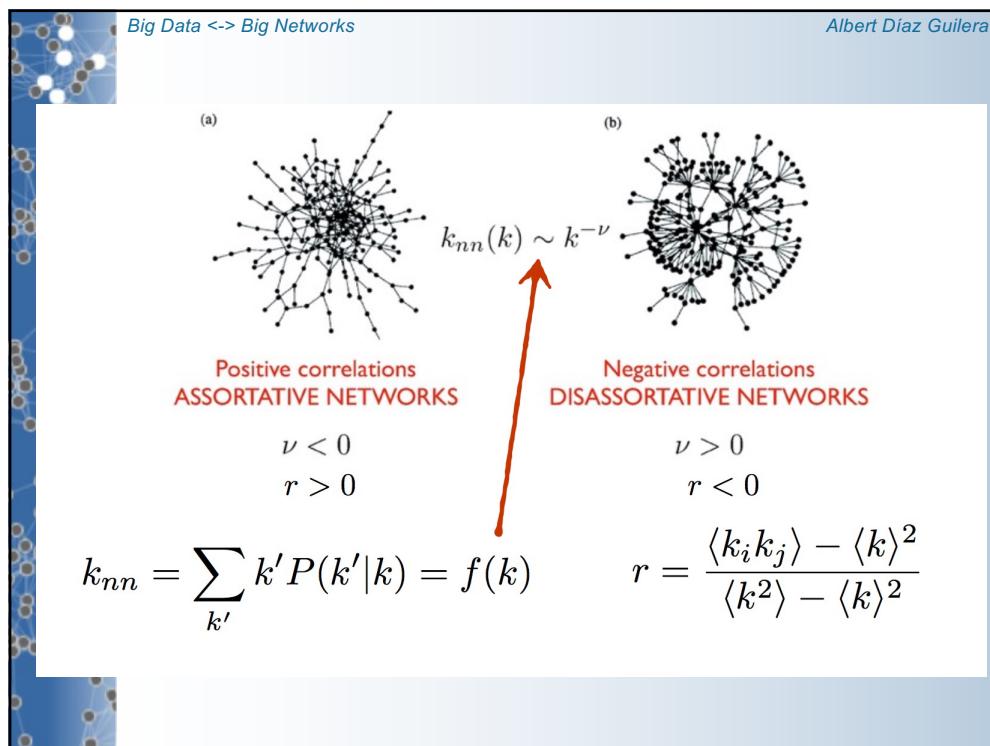
Two ways of measuring

$$k_{nn} = \sum_{k'} k' P(k'|k) = f(k) \quad r = \frac{\langle k_i k_j \rangle - \langle k \rangle^2}{\langle k^2 \rangle - \langle k \rangle^2}$$

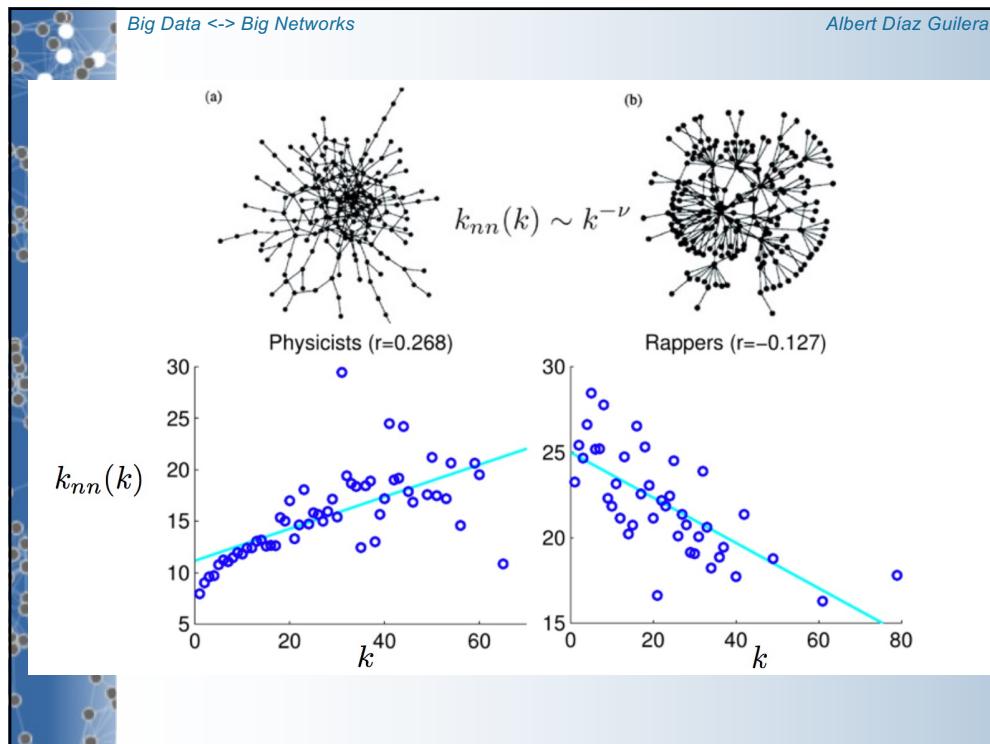

55



56



57



58

Big Data <-> Big Networks *Albert Díaz Guilera*

Network taxonomy

	Network	Type	Nodes	Links	$\langle k \rangle$	L	Clustering	Corr. (r)
Social	film actors	undirected	449 913	25 516 482	113.43	3.48	2.3	0.20
	company directors	undirected	7 673	55 392	14.44	4.60	—	0.59
	math coauthorship	undirected	253 339	496 489	3.92	7.57	—	0.15
	physics coauthorship	undirected	52 909	245 300	9.27	6.19	—	0.45
	biology coauthorship	undirected	1 520 251	11 803 064	15.53	4.92	—	0.088
	telephone call graph	undirected	47 000 000	80 000 000	3.16	—	2.1	0.276
	email messages	directed	59 912	86 300	1.44	4.95	1.5/2.0	0.120
	email address books	directed	16 881	57 029	3.38	5.22	—	0.34
	student relationships	undirected	573	477	1.66	16.01	—	0.363
	sexual contacts	undirected	2 810	—	—	—	3.2	0.127
Information	WWW nd.edu	directed	269 504	1 497 135	5.55	11.27	2.1/2.4	0.11
	WWW Altavista	directed	203 549 046	2 130 000 000	10.46	16.18	2.1/2.7	0.29
	citation network	directed	783 339	6 716 198	8.57	—	3.0/—	—0.067
	Rogot's Thesaurus	directed	1 022	5 103	4.99	4.87	—	0.13
	word co-occurrence	undirected	460 902	17 000 000	70.13	—	2.7	0.44
Technological	Internet	undirected	10 697	31 992	5.98	3.31	2.5	0.035
	power grid	undirected	4 941	6 594	2.67	18.99	—	0.39
	train routes	undirected	587	19 603	66.79	2.16	—	0.080
	software packages	directed	1 439	1 723	1.20	2.42	1.6/1.4	—0.003
	software classes	directed	1 377	2 213	1.61	1.51	—	0.082
	electronic circuits	undirected	24 097	53 248	4.34	11.05	3.0	—0.116
Biological	peer-to-peer network	undirected	880	1 296	1.47	4.28	2.1	—0.154
	metabolic network	undirected	765	3 686	9.64	2.56	2.2	0.012
	protein interactions	undirected	2 115	2 240	2.12	6.80	2.4	0.011
	marine food web	directed	135	598	4.43	2.05	—	0.23
	freshwater food web	directed	92	997	10.84	1.90	—	—0.263
Neural	neural network	directed	307	2 359	7.68	3.97	—	0.48
	—	—	—	—	—	—	—	—0.326

59

Big Data <-> Big Networks *Albert Díaz Guilera*

Assortativity

<code>degree_assortativity_coefficient (G[, x, y, ...])</code>	Compute degree assortativity of graph.
<code>attribute_assortativity_coefficient (G, attribute)</code>	Compute assortativity for node attributes.
<code>numeric_assortativity_coefficient (G, attribute)</code>	Compute assortativity for numerical node attributes.
<code>degree_pearson_correlation_coefficient (G[, ...])</code>	Compute degree assortativity of graph.

Average neighbor degree

<code>average_neighbor_degree (G[, source, target, ...])</code>	Returns the average degree of the neighborhood of each node.
---	--

62



Big Data <-> Big Networks

Albert Díaz-Guilera

<https://cambridgeuniversitypress.github.io/FirstCourseNetworkScience/>

- Chapter 3: Hubs ([Tutorial](#))
- Chapter 4: Direction & Weights ([Tutorial](#))



63