

Анализ сетевых данных

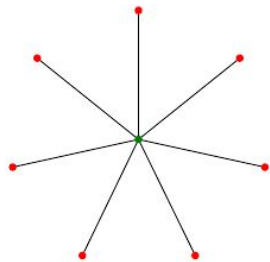
Лекция 2. Random Graph models.

11 Сентября, 2019

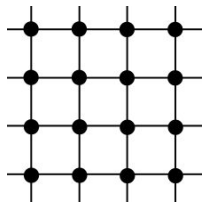
MOOCs

- <http://tuvalu.santafe.edu/~aaronc/courses/5352/> **Aaron Clauset**, lecture notes
- <https://goo.gl/8CghUx> **Leonid Zhukov**, <http://www.leonidzhukov.net/>, Full course videos + lecture notes.
- <http://web.stanford.edu/class/cs224w/> **Jure Leskovec**, lecture notes, bunch of useful materials, *videos are unavailable*.

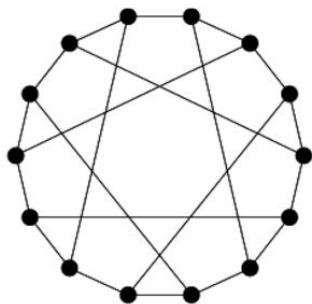
Compare different networks



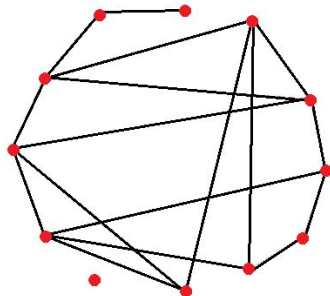
Star graph



Grid graph

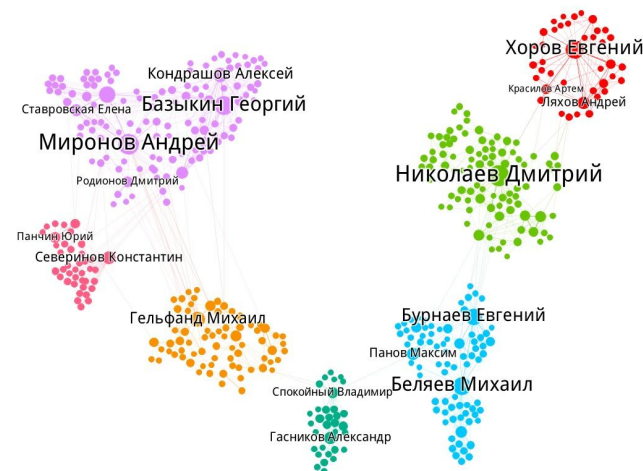


K-regular graph



Random graph

VS



Real network

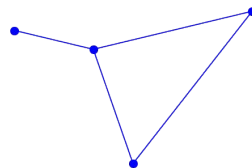
...?

How to compare different networks?

- # nodes, #edges, density
- (Average) node degree
- (Average) clustering coefficient
- (Average) path length
- Node degree distribution
- Centrality measures (next lecture)
- DL Embeddings/Graph convolutions

How to compare different networks?

- # nodes, #edges, density
- (Average) node degree
- (Average) clustering coefficient
- (Average) path length
- Node degree distribution
- Centrality measures (next lecture)
- DL Embeddings/Graph convolutions



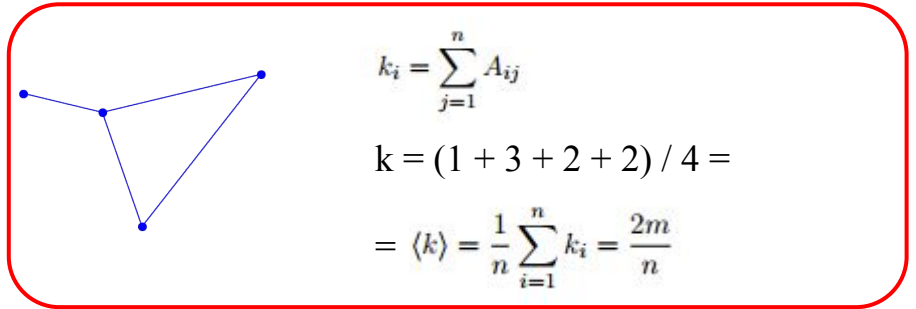
nodes = 4 (n)

edges = 4 (m)

density = $m / [n(n-1)/2] = 2/3$

How to compare different networks?

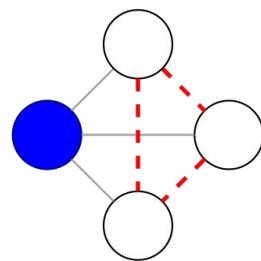
- # nodes, #edges, density
- (Average) node degree
- (Average) clustering coefficient
- (Average) path length
- Node degree distribution
- Centrality measures (next lecture)
- DL Embeddings/Graph convolutions



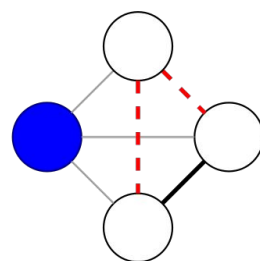
How to compare different networks?

- # nodes, #edges, density
- (Average) node degree
- (Average) clustering coefficient
- (Average) path length
- Node degree distribution
- Centrality measures (next lecture)
- DL Embeddings/Graph convolutions

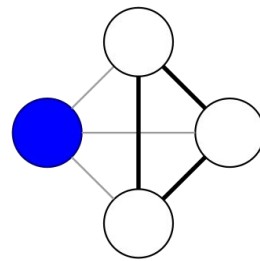
$$C_i = \frac{(\text{number of pairs of neighbors of } i \text{ that are connected})}{(\text{number of pairs of neighbors of } i)}$$
$$= \sum_{jk} A_{ij} A_{jk} A_{ki} / \binom{k_i}{2} = \frac{2 n_i}{k_i (k_i - 1)}, n_i = \# \text{ neighbours of node } i$$



$$c = 0$$



$$c = 1/3$$

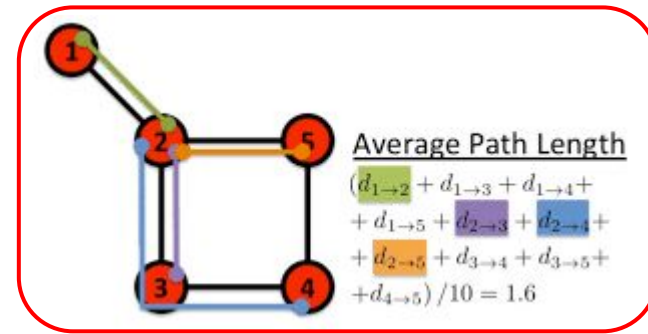


$$c = 1$$

$$C = \frac{(\text{number of triangles}) \times 3}{(\text{number of connected triples})}$$
$$= \sum_{ijk} A_{ij} A_{jk} A_{ki} / \sum_{ijk} A_{ij} A_{jk},$$

How to compare different networks?

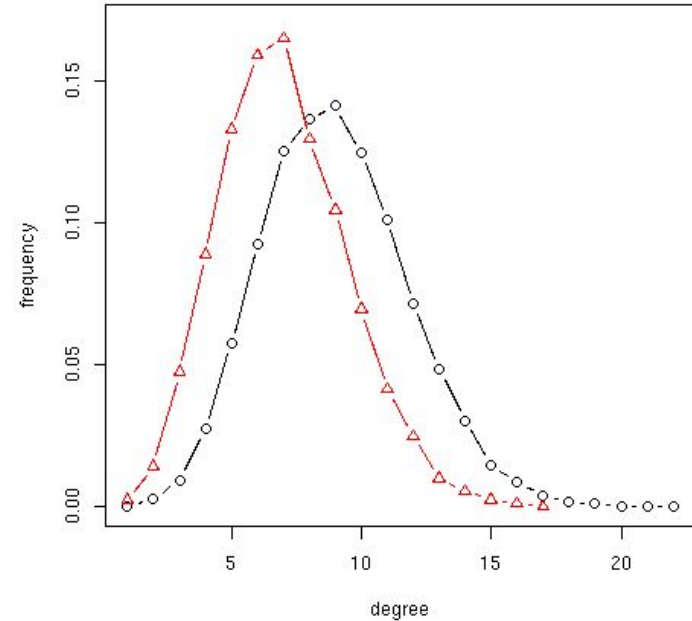
- # nodes, #edges, density
- (Average) node degree
- (Average) clustering coefficient
- (Average) path length
- Node degree distribution
- Centrality measures (next lecture)
- DL Embeddings/Graph convolutions



[source](#)

How to compare different networks?

- # nodes, #edges, density
- (Average) node degree
- (Average) clustering coefficient
- (Average) path length
- **Node degree distribution**
- Centrality measures (next lecture)
- DL Embeddings/Graph convolutions



Erdos-Renyi random graph

- $G_{n,p}$ model:

Graph with n nodes and for each pair of nodes the probability of an edge between them is equal to p .

- $G_{n,m}$ model:

A randomly selected graph from the set of C_N^m graphs, with $N = n(n-1)/2$, where $n = \text{\#nodes}$ and $m = \text{\#edges}$

Random graph model (Erdos & Renyi, 1959)

Erdos-Renyi random graph

- $G_{n,p}$ model:

Graph with n nodes and for each pair of nodes the probability of an edge between them is equal to p .

- $G_{n,m}$ model:

A randomly selected graph from the set of C_N^m graphs, with $N = n(n-1)/2$, where $n = \text{\#nodes}$ and $m = \text{\#edges}$

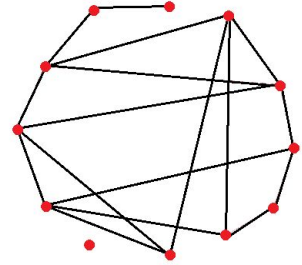
Random graph model (Erdos & Renyi, 1959)

Эти две модели -
эквивалентны,
Приведем набросок
доказательства

Erdos-Renyi properties

$G_{n,p}$ model:

- $\langle m \rangle = p * n * (n-1) / 2$
- $\langle k \rangle = (n-1) * p \approx n * p$

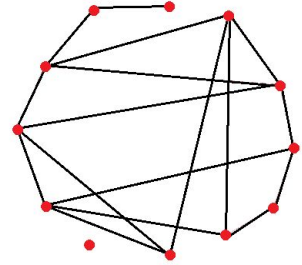


Erdos-Renyi properties

$G_{n,p}$ model:

- $\langle m \rangle = p * n * (n-1) / 2$
- $\langle k \rangle = (n-1) * p \approx n * p$

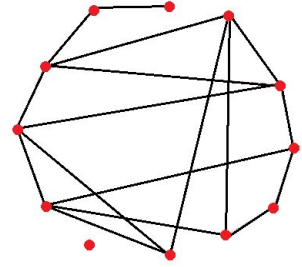
What is node degree distribution?



Erdos-Renyi properties

$G_{n,p}$ model:

- $\langle m \rangle = p * n * (n-1) / 2$
- $\langle k \rangle = (n-1) * p \approx n * p$



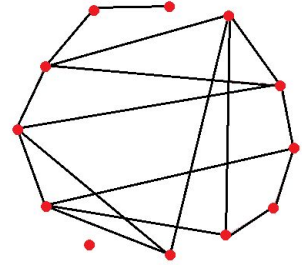
What is node degree distribution?

Probability that given node i has degree $k_i = k$

Erdos-Renyi properties

$G_{n,p}$ model:

- $\langle m \rangle = p * n * (n-1) / 2$
- $\langle k \rangle = (n-1) * p \approx n * p$



What is node degree distribution?

Probability that given node i has degree $k_i = k$

$$P(k_i = k) = P(k) = C_{n-1}^k p^k (1 - p)^{n-1-k}$$

(Bernoulli distribution)

Erdos-Renyi degree distribution

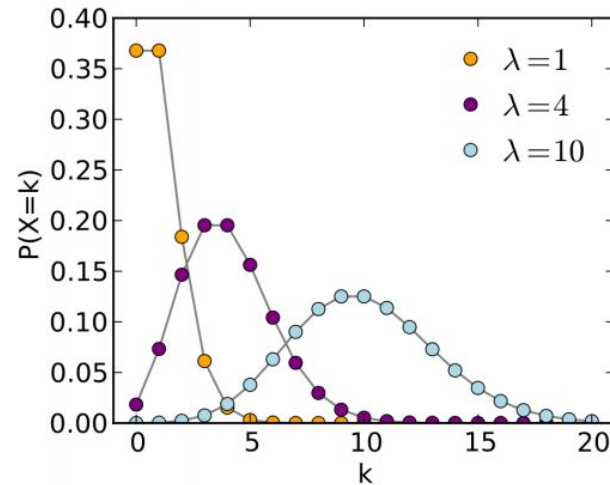
Limiting case of Bernoulli distribution (when n goes to infinity) - Poisson distribution (with parameter $\lambda = \langle k \rangle = np$)

$$P(k) = \frac{\langle k \rangle^k e^{-\langle k \rangle}}{k!} = \frac{\lambda^k e^{-\lambda}}{k!}$$

Erdos-Renyi degree distribution

Limiting case of Bernoulli distribution (when n goes to infinity) - Poisson distribution (with parameter $\lambda = \langle k \rangle = np$)

$$P(k) = \frac{\langle k \rangle^k e^{-\langle k \rangle}}{k!} = \frac{\lambda^k e^{-\lambda}}{k!}$$

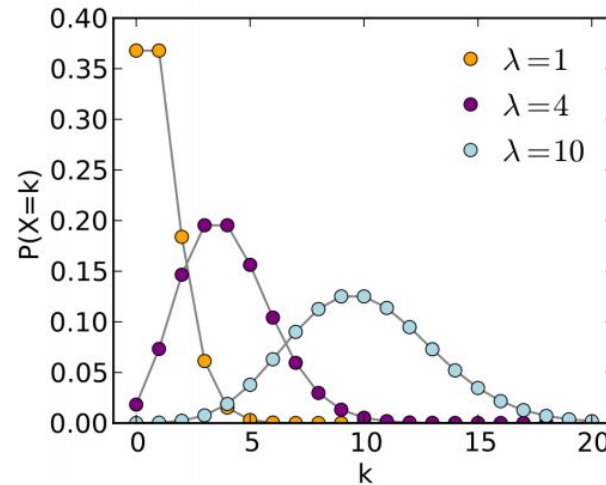


$$P(k_i = k) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad \lambda = pn$$

Erdos-Renyi degree distribution

Limiting case of Bernoulli distribution (when n goes to infinity) - Poisson distribution (with parameter $\lambda = \langle k \rangle = (n-1)p$)

$$P(k) = \frac{\langle k \rangle^k e^{-\langle k \rangle}}{k!} = \frac{\lambda^k e^{-\lambda}}{k!}$$



Which in case of $np \rightarrow \infty$,
goes to Gaussian

$$P(k_i = k) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad \lambda = pn$$

$G_{n,p}$ VS $G_{n,m}$

Bernoulli distribution

- Mean = $\langle k \rangle = (n-1)p$
- Variance = $\sigma^2 = p(1-p)(n-1)$

With fixed p and $n \rightarrow \infty$, distribution becomes *narrow*:

$G_{n,p}$ VS $G_{n,m}$

Bernoulli distribution

- Mean = $\langle k \rangle = (n-1)p$
- Variance = $\sigma^2 = p(1-p)(n-1)$

With fixed p and $n \rightarrow \infty$, distribution becomes *narrow*:

$$\sigma / \langle k \rangle = [(1-p) / p(n-1)]^{1/2} \approx 1 / (n-1)^{1/2}$$

thus we are increasingly confident that the degree of a node is equal to $\langle k \rangle$

$G_{n,p}$ vs $G_{n,m}$

Bernoulli distribution

- Mean = $\langle k \rangle = (n-1)p$
- Variance = $\sigma^2 = p(1-p)(n-1)$

With fixed p and $n \rightarrow \infty$, distribution becomes *narrow*:

$$\sigma / \langle k \rangle = [(1-p) / p(n-1)]^{1/2} \approx 1 / (n-1)^{1/2}$$

thus we are increasingly confident that the degree of a node is equal to $\langle k \rangle$

$G_{n,p}$ and $G_{n,m}$ are the same

Erdos-Renyi clustering coefficient

$$C_i = \frac{2 n_i}{k_i (k_i - 1)}$$

since edges appear i.i.d. with probability p :

$$n_i = p * k_i (k_i - 1) / 2$$

then $C_i \approx \langle k \rangle / n$

Erdos-Renyi clustering coefficient

$$C_i = \frac{2 n_i}{k_i (k_i - 1)}$$

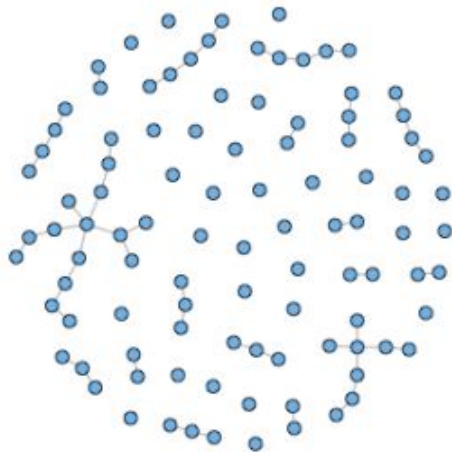
since edges appear i.i.d. with probability p :

$$n_i = p * k_i (k_i - 1) / 2$$

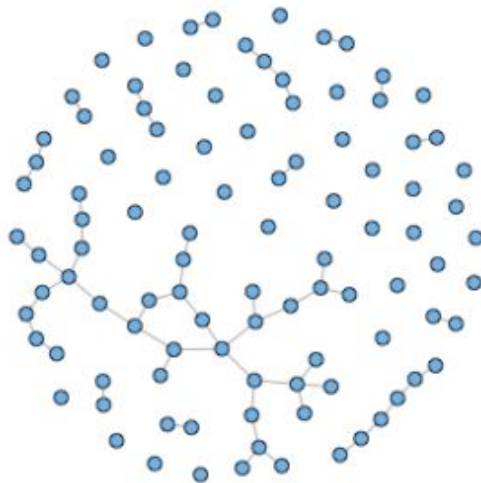
$$\text{then } C_i \approx \langle k \rangle / n$$

This means that with n goes to infinity clustering coefficient of a random graph goes to 0

What about connectivity?



$$p < p_c$$



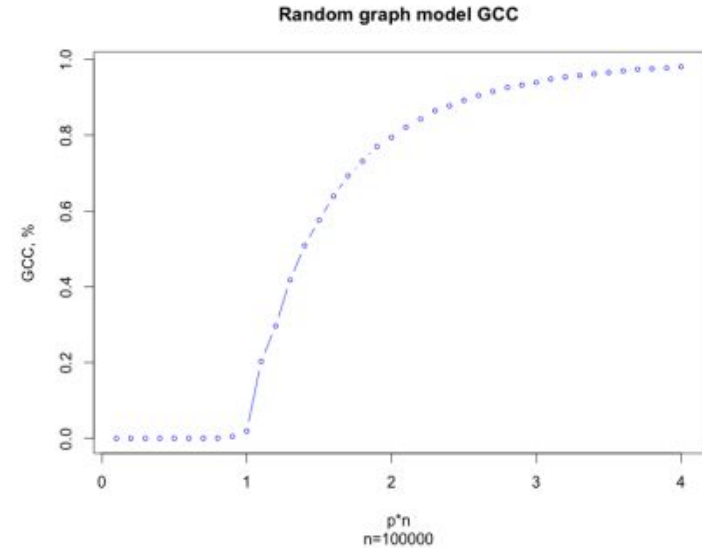
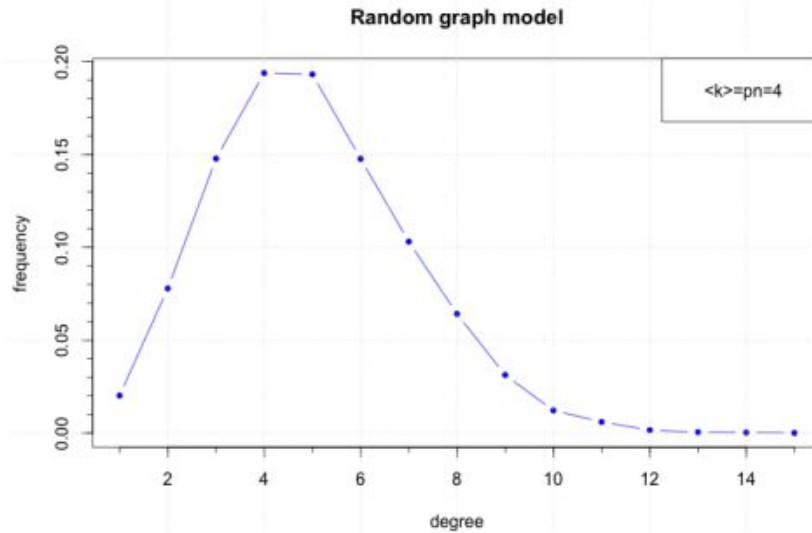
$$p = p_c$$



$$p > p_c$$

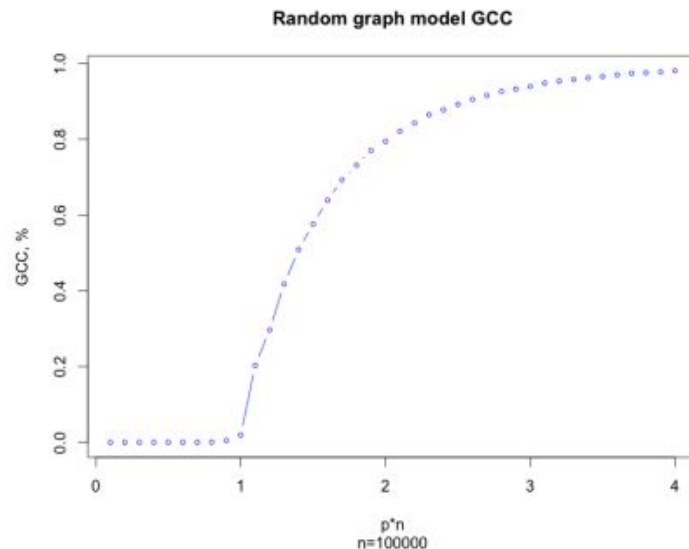
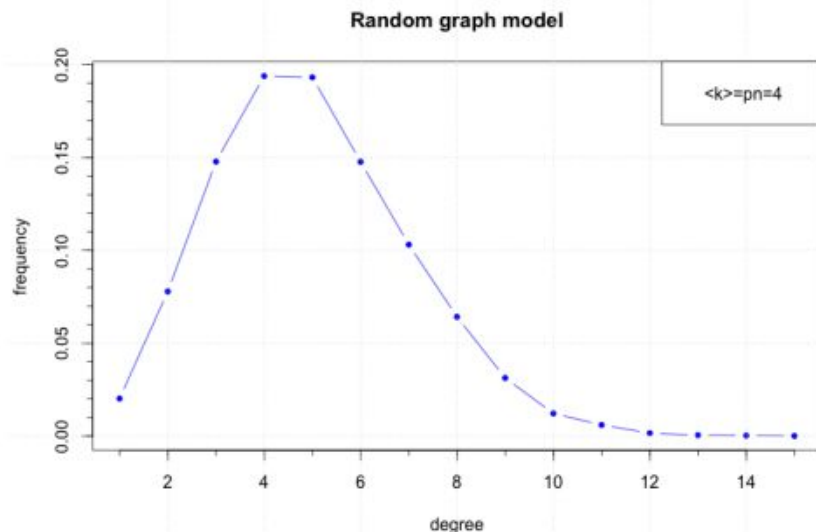
What about connectivity?

$$\langle k \rangle = pn$$



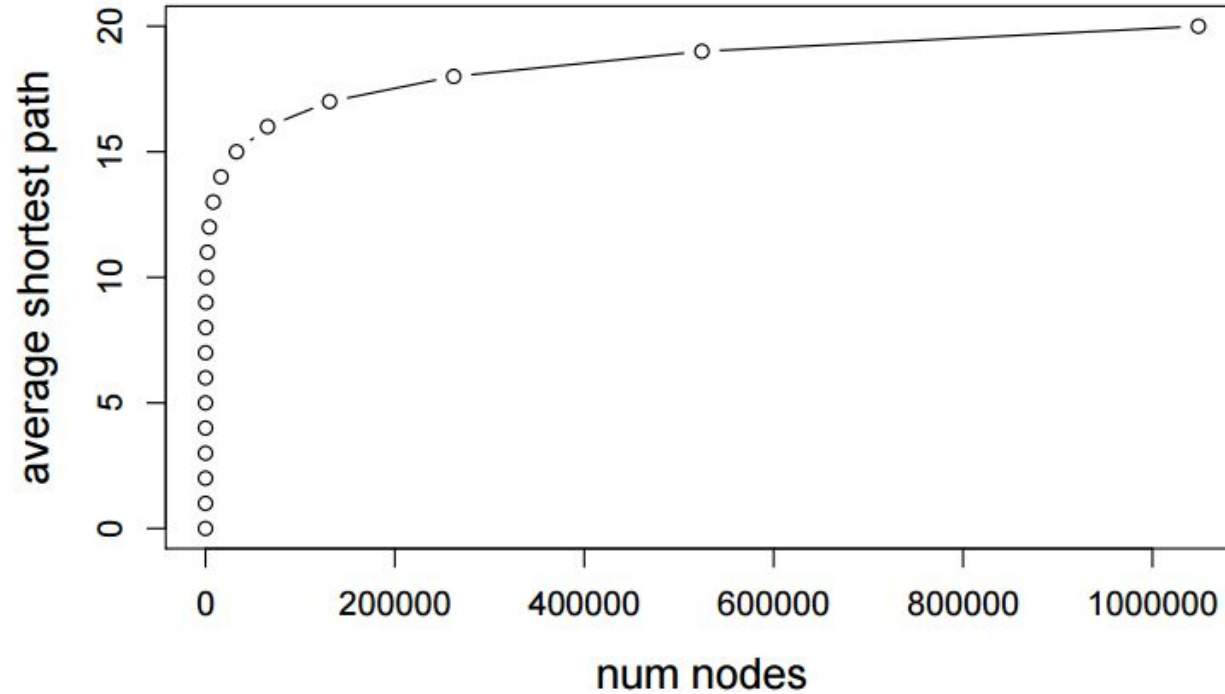
What about connectivity?

$$\langle k \rangle = pn$$



- It could be shown that with $\langle k \rangle = 1$ the largest connected component contains $O(n^{2/3})$ nodes.
- With $\langle k \rangle > 1$ it quickly has all the nodes.

Average path length



For Erdos-Renyi graph average path length is of order $O(\log n)$

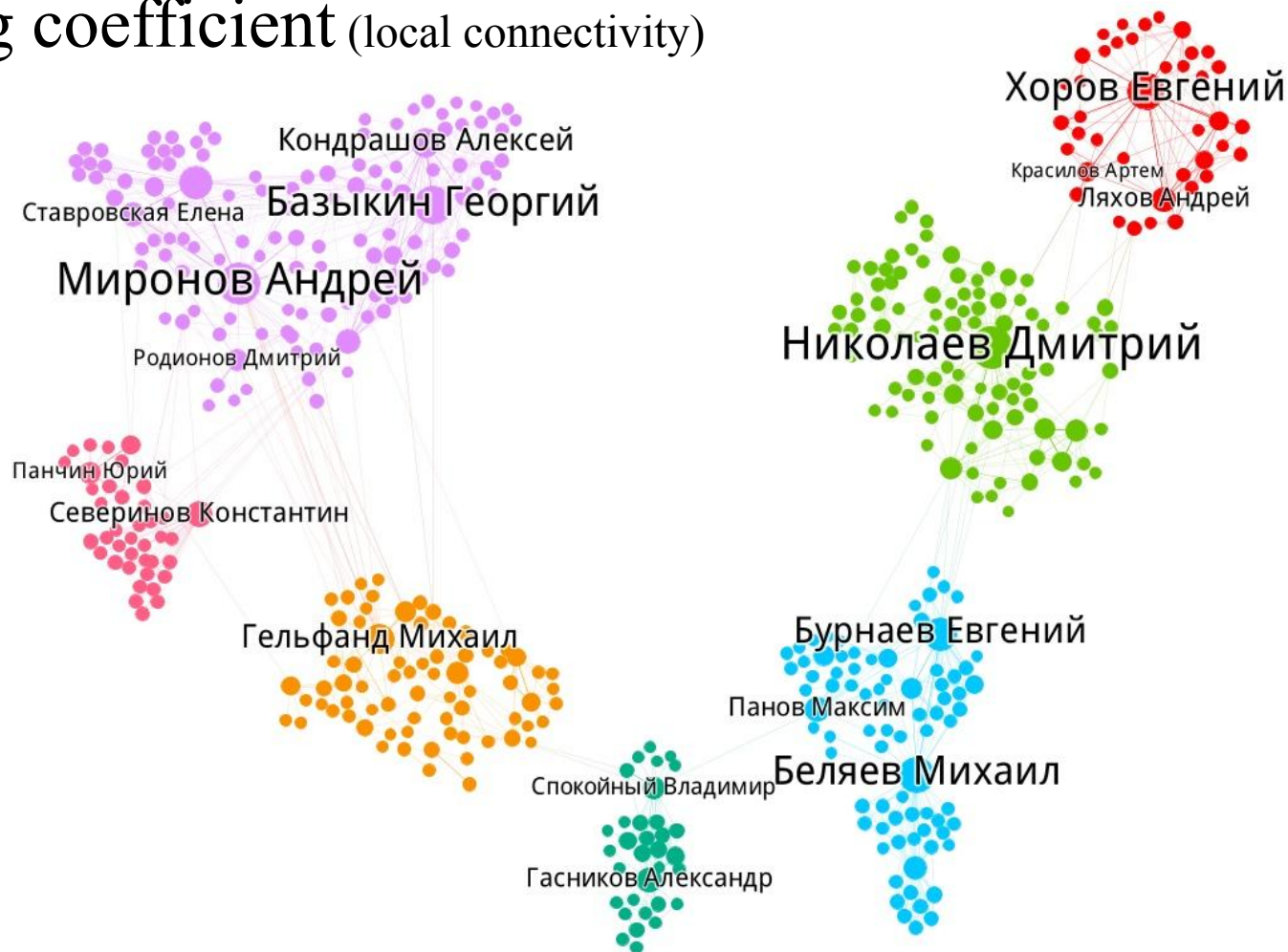
How to compare different networks?

- # nodes, #edges, density
- (Average) node degree
- (Average) clustering coefficient
- (Average) path length
- **Node degree distribution**
- Centrality measures (next lecture)

Key properties

- (Average) clustering coefficient
- (Average) path length
- Node degree distribution

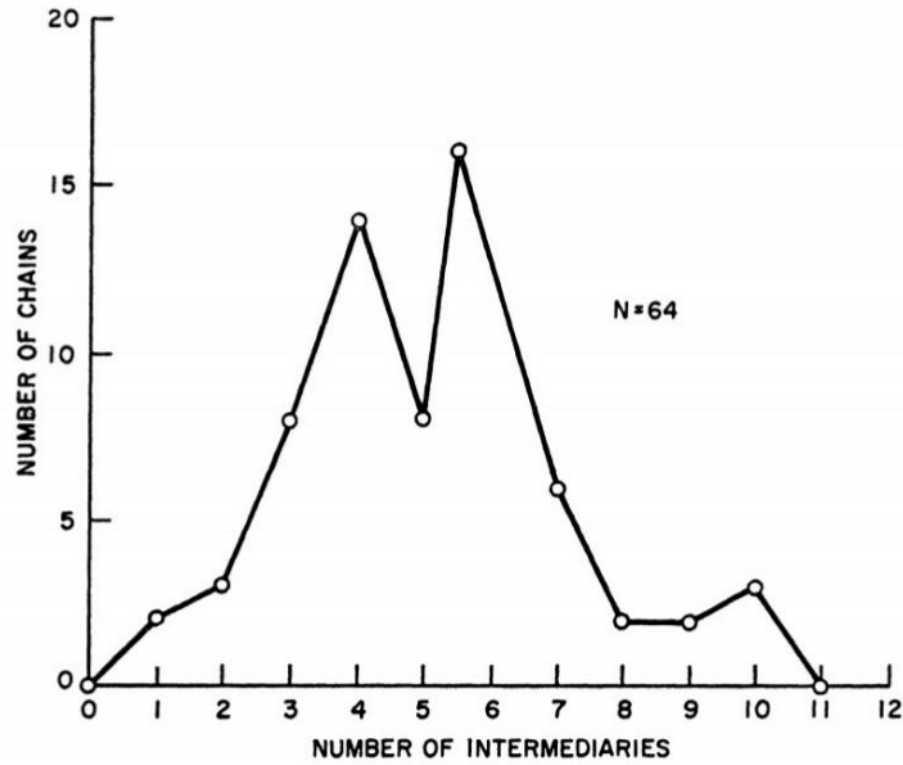
Clustering coefficient (local connectivity)



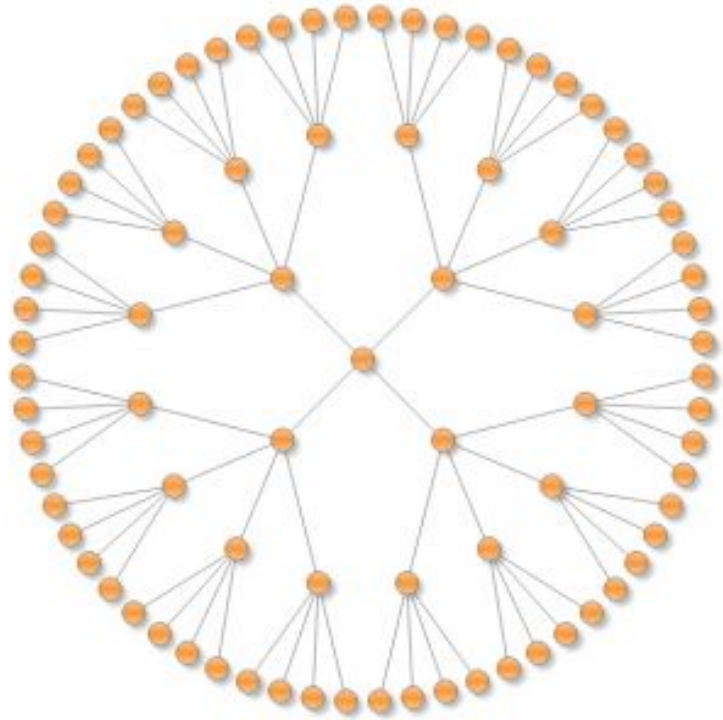
Key properties

- (Average) clustering coefficient
- (Average) path length
- Node degree distribution

Recall Milgram's experiment



Average path length (idea)



Consider a simple model:

Each person has the same number of friends z , total # of people in the world is N , then what is a diameter?

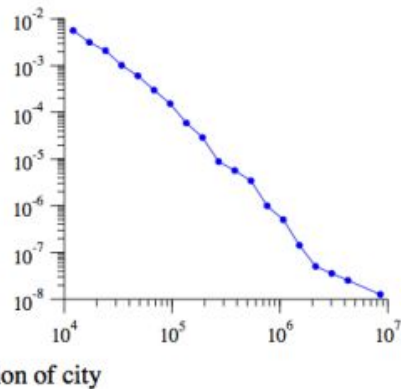
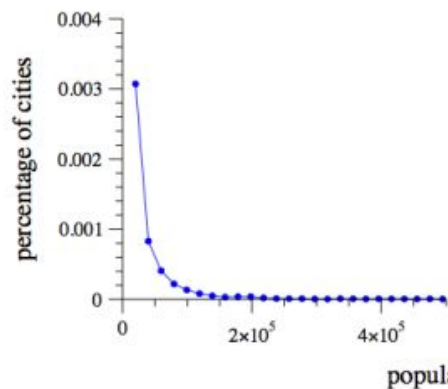
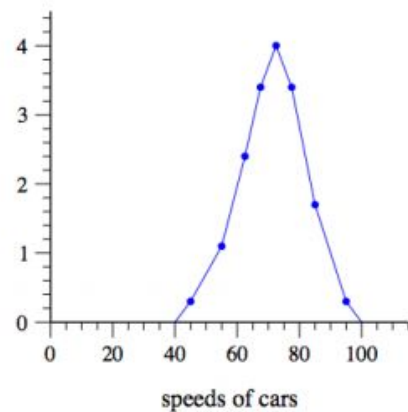
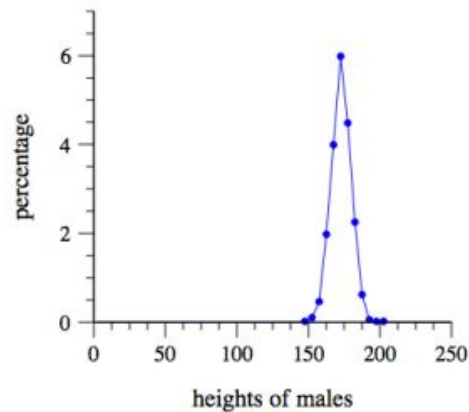
Diameter = the longest of all the calculated shortest paths in a network

An estimate: $z^d = N$, $d = \log N / \log z$
 $N \approx 6.7 \text{ bln}$, $z = 50 \text{ friends}$, $d \approx 5.8$.

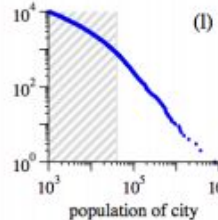
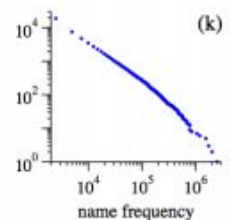
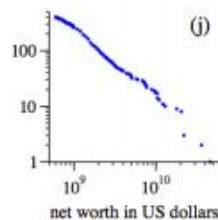
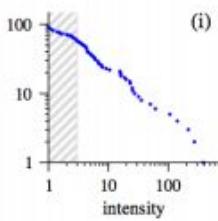
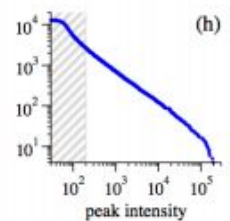
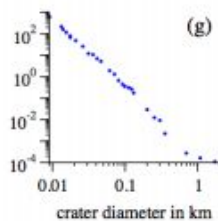
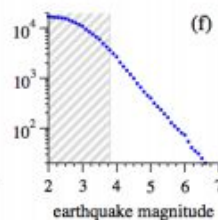
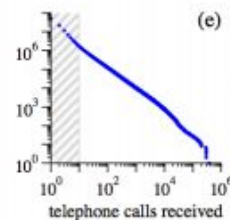
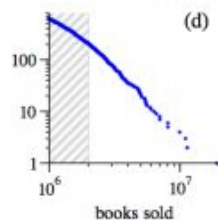
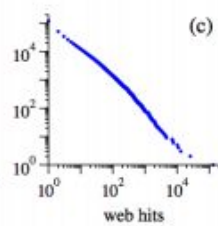
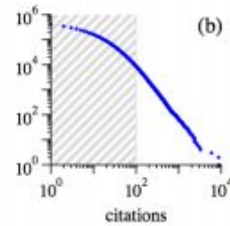
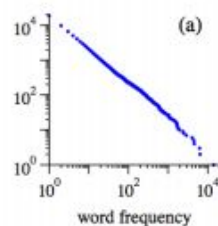
Key properties

- (Average) clustering coefficient
- (Average) path length
- Node degree distribution

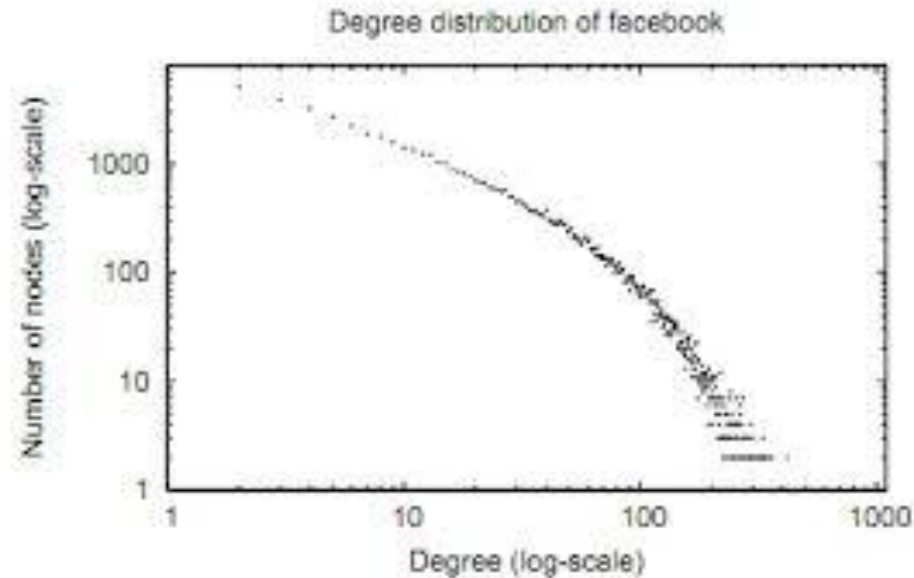
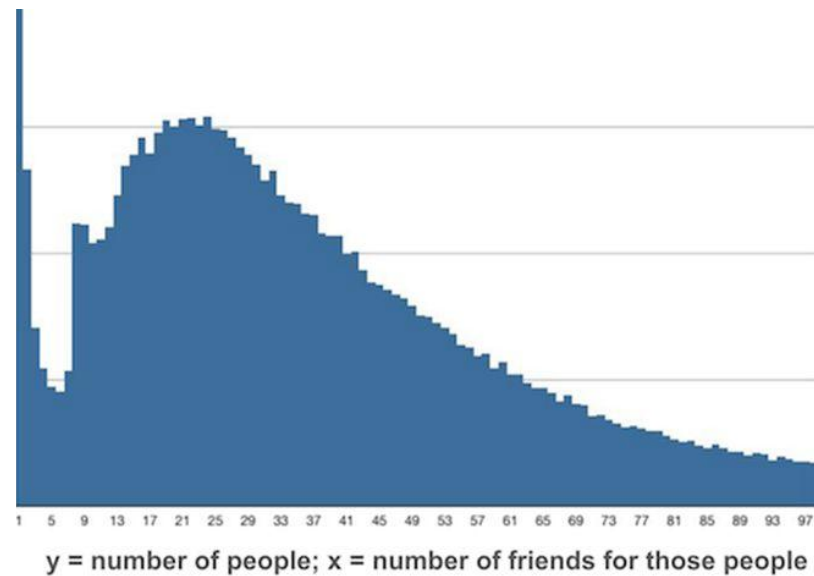
Empirical distributions



log-log scale



Facebook degree distribution



Empirical network features

- Power-law (heavy-tailed) degree distribution
- Small average distance (graph diameter)
- Large clustering coefficient (transitivity)
- Giant connected component, hierarchical structure, etc

