

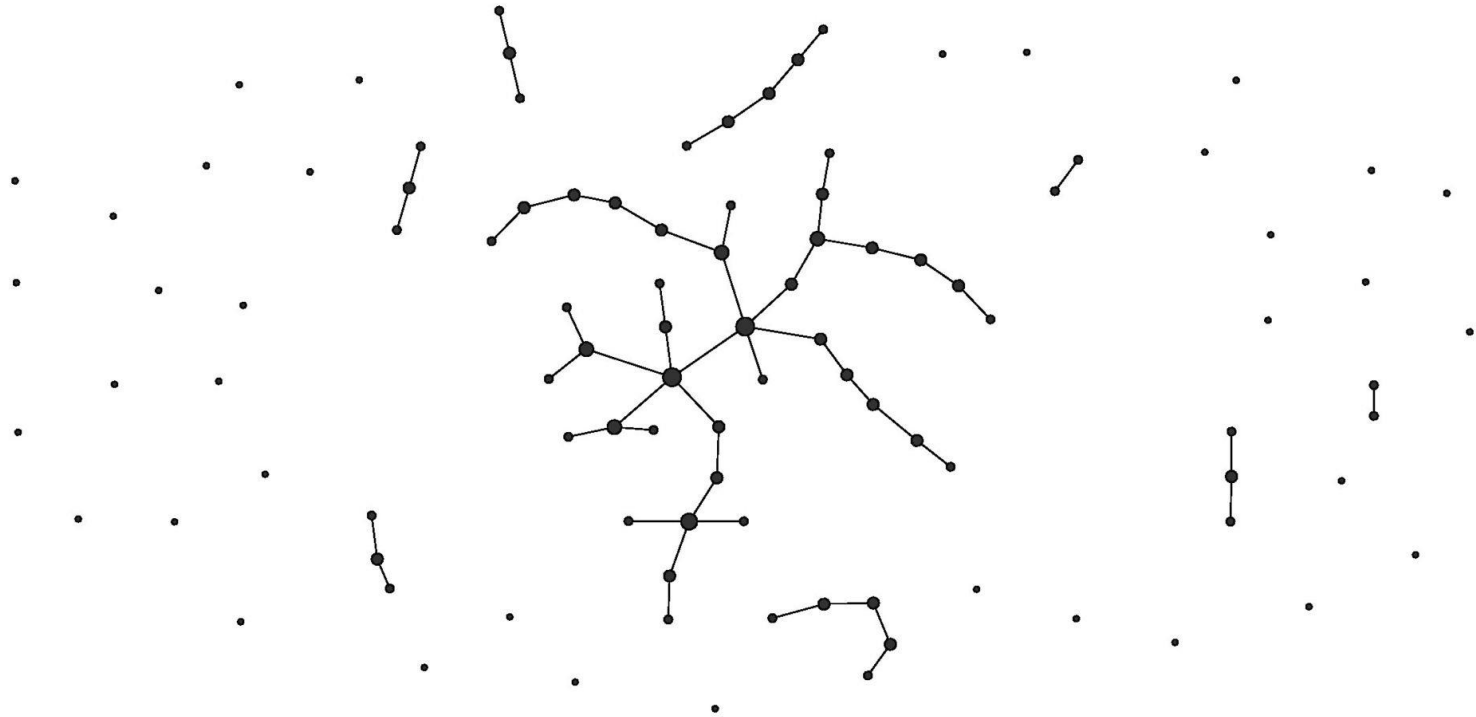
CSE 258 – Lecture 13

Web Mining and Recommender Systems

Triadic closure; strong & weak ties

Monday...

Random models of networks: Erdos Renyi random graphs



(picture from Wikipedia http://en.wikipedia.org/wiki/Erd%C5%91s%E2%80%93R%C3%A9nyi_model)

Monday...

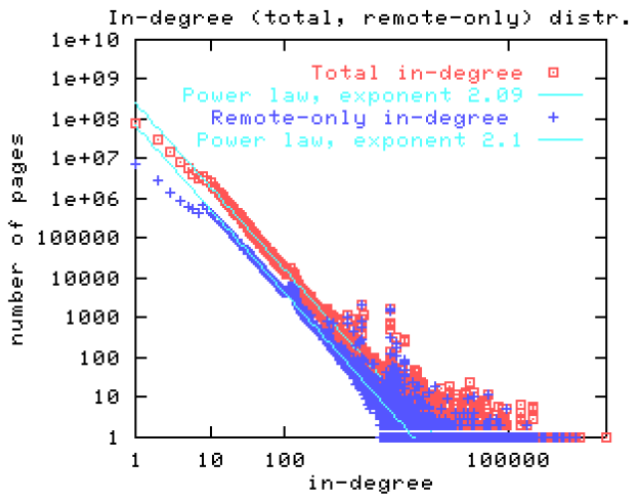
Preferential attachment models of network formation

Consider the following process to generate a network (e.g. a web graph):

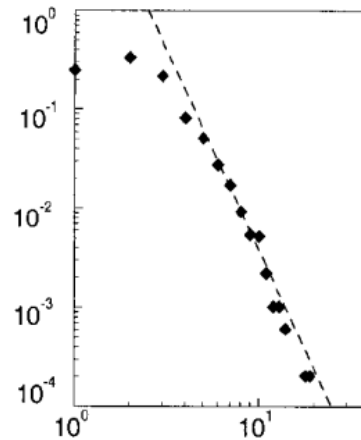
1. Order all of the N pages $1, 2, 3, \dots, N$ and repeat the following process for each page j :
2. Use the following rule to generate a link to another page:
 - a. With probability p , link to a random page $i < j$
 - b. Otherwise, choose a random page i and link to the page ***i links to***

Monday – power laws

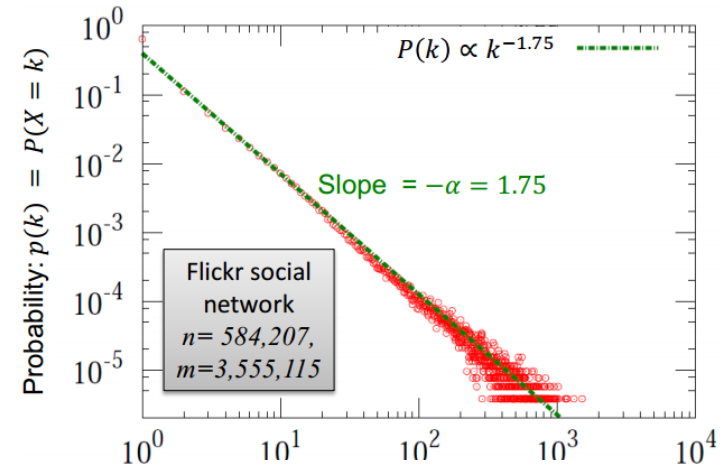
- Social and information networks often follow **power laws**, meaning that a few nodes have **many** of the edges, and many nodes have **a few** edges



e.g. web graph
(Broder et al.)



e.g. power grid
(Barabasi-Albert)



e.g. Flickr
(Leskovec)

How can we **characterize, model, and reason about** the structure of social networks?

1. Models of network structure
2. Power-laws and scale-free networks, "rich-get-richer" phenomena
3. Triadic closure and "the strength of weak ties"
4. Small-world phenomena
5. Hubs & Authorities; PageRank

Triangles

So far we've seen (a little about) how networks can be characterized by their connectivity patterns

What more can we learn by looking at higher-order properties, such as relationships between **triplets** of nodes?

Motivation

Q: Last time you found a job, was it through:

- A complete stranger?
 - A close friend?
 - An acquaintance?

A: Surprisingly, people often find jobs through **acquaintances** rather than through close friends (Granovetter, 1973)

Motivation

- Your friends (hopefully) would seem to have the greatest motivation to help you
- But! Your closest friends have limited information that you don't already know about
- Alternately, acquaintances act as a "bridge" to a different part of the social network, and expose you to new information

This phenomenon is known as **the strength of weak ties**

Motivation

- To make this concrete, we'd like to come up with some notion of "tie strength" in networks
- To do this, we need to go beyond just looking at edges in isolation, and looking at how an edge connects one part of a network to another

Refs:

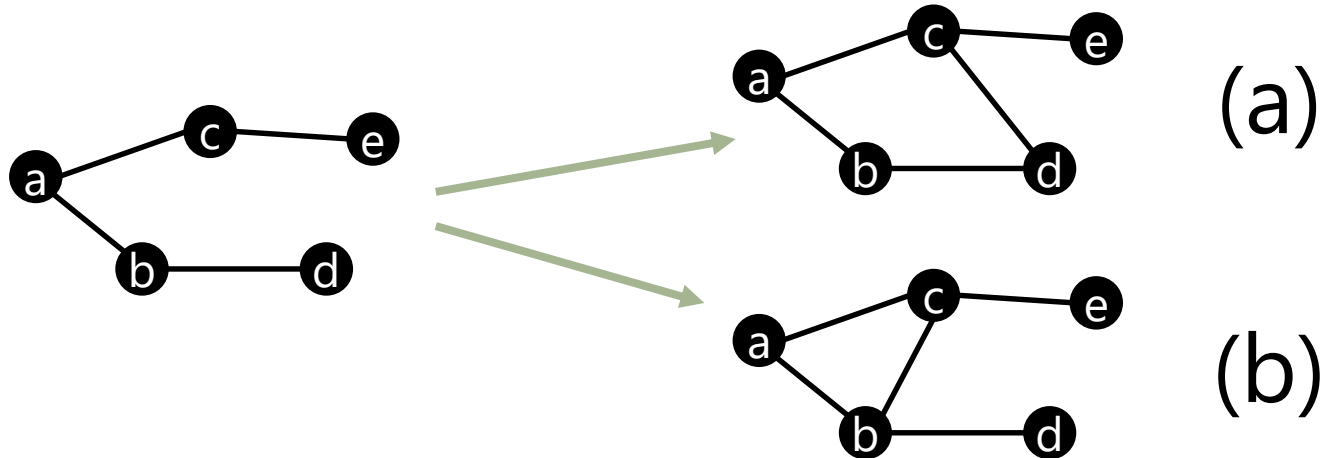
"The Strength of Weak Ties", Granovetter (1973): <http://goo.gl/wVJVIN>

"Getting a Job", Granovetter (1974)

Triangles

Triadic closure

Q: Which edge is most likely to form **next** in this (social) network?



A: (b), because it creates a **triad** in the network

Triangles

“If two people in a social network have a friend in common, then there is an increased likelihood that they will become friends themselves at some point in the future” (Ropoport, 1953)

Three reasons (from Heider, 1958; see Easley & Kleinberg):

- Every mutual friend *a* between *bob* and *chris* gives them an **opportunity** to meet
- If *bob* is friends with *ashton*, then knowing that *chris* is friends with *ashton* gives *bob* a reason to **trust** *chris*
- If *chris* and *bob* don't become friends, this causes stress for *ashton* (having two friends who don't like each other), so there is an **incentive** for them to connect

Triangles

The extent to which this is true is measured by the (local)
clustering coefficient:

- The clustering coefficient of a node i is the probability that two of i 's friends will be friends with each other:

neighbours of i pairs of neighbours that are edges

$$C_i = \frac{\sum_{j,k \in \Gamma(i)} \delta((j,k) \in E)}{k_i(k_i - 1)}$$

(edges (j,k) and (k,j) are both counted for undirected graphs)

degree of node i

- This ranges between 0 (none of my friends are friends with each other) and 1 (all of my friends are friends with each other)

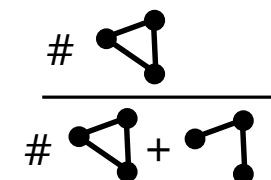
Triangles

The extent to which this is true is measured by the (local) **clustering coefficient**:

- The clustering coefficient of the **graph** is usually defined as the average of local clustering coefficients

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

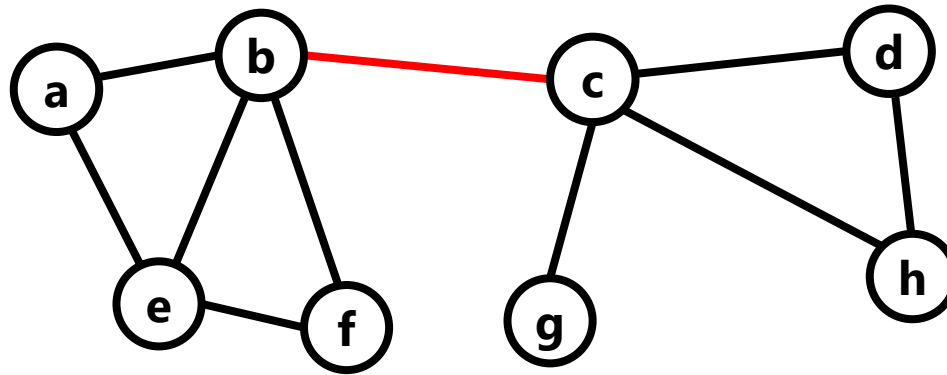
- Alternately it can be defined as the fraction of connected triplets in the graph that are closed (these do not evaluate to the same thing!):

$$C = \frac{\# \text{ of closed triplets}}{\# \text{ of connected triplets}}$$


Bridges

Next, we can talk about the role of edges in relation to the rest of the network, starting with a few more definitions

1. Bridge edge

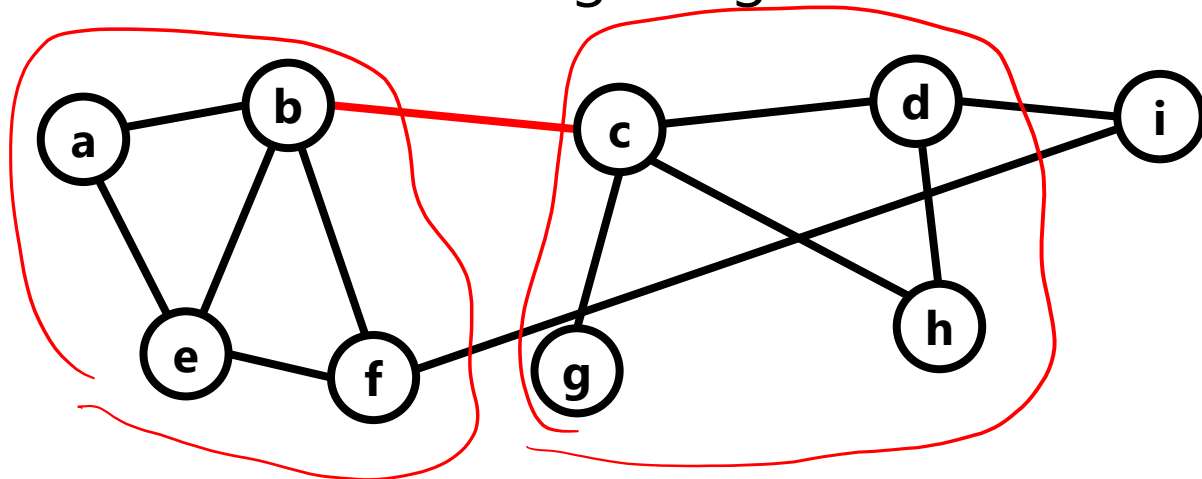


An edge (b,c) is a **bridge** edge if removing it would leave no path between b and c in the resulting network

Bridges

In practice, "bridges" aren't a very useful definition, since there will be very few edges that completely isolate two parts of the graph

2. **Local** bridge edge

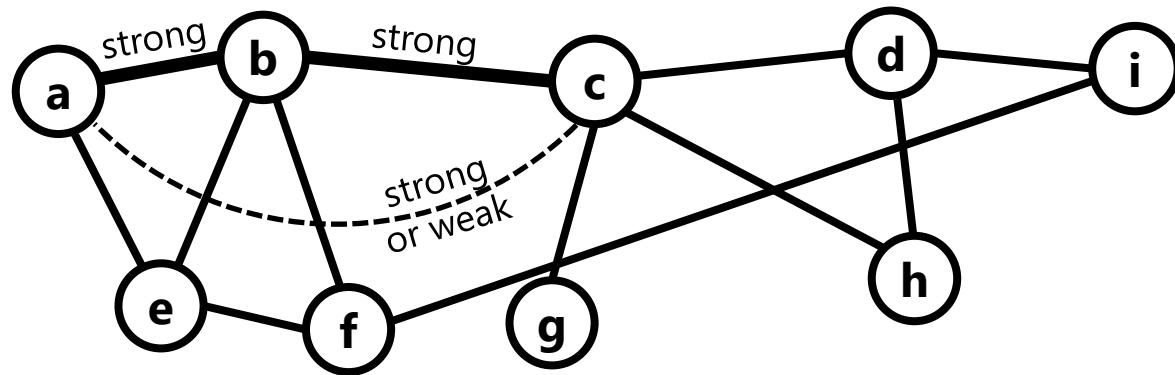


An edge (b,c) is a **local bridge** if removing it would leave no edge between b's friends and c's friends (though there could be more distant connections)

Strong & weak ties

We can now define the concept of “strong” and “weak” ties (which roughly correspond to notions of “friends” and “acquaintances”)

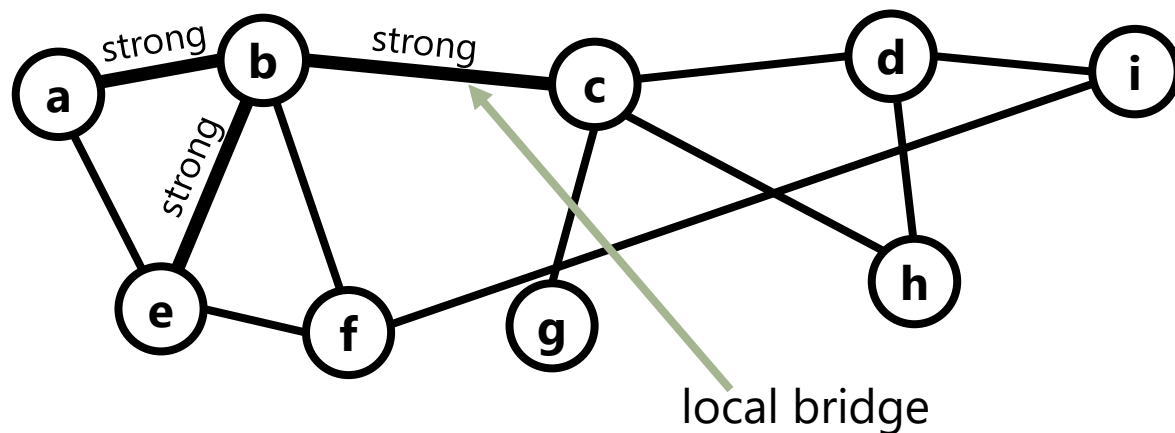
3. Strong triadic closure property



If (a,b) and (b,c) are connected by **strong** ties, there must be at least a **weak** tie between a and c

Strong & weak ties

Granovetter's theorem: if the strong triadic closure property is satisfied for a node, and that node is involved in two strong ties, then any incident local bridge must be a **weak tie**



Proof (by contradiction): (1) b has two strong ties (to a and e); (2) suppose it has a **strong** tie to c via a local bridge; (3) but now a tie must exist between c and a (or c and e) due to strong triadic closure; (4) so $b \rightarrow c$ cannot be a bridge

Strong & weak ties

Granovetter's theorem: so, if we're receiving information from distant parts of the network (i.e., via "local bridges") then we must be receiving it via **weak ties**

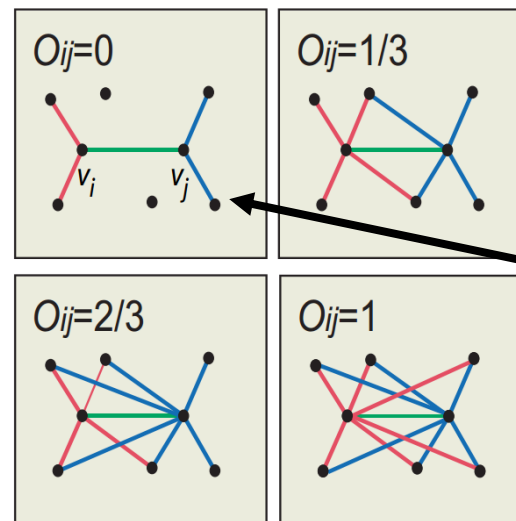
Q: How to test this theorem empirically on real data?

A: Onnela et al. 2007 studied networks of mobile phone calls

Defn. 1: Define the "overlap" between two nodes to be the Jaccard similarity between their connections

$$O_{i,j} = \frac{\Gamma(i) \cap \Gamma(j)}{\Gamma(i) \cup \Gamma(j)}$$

neighbours of i



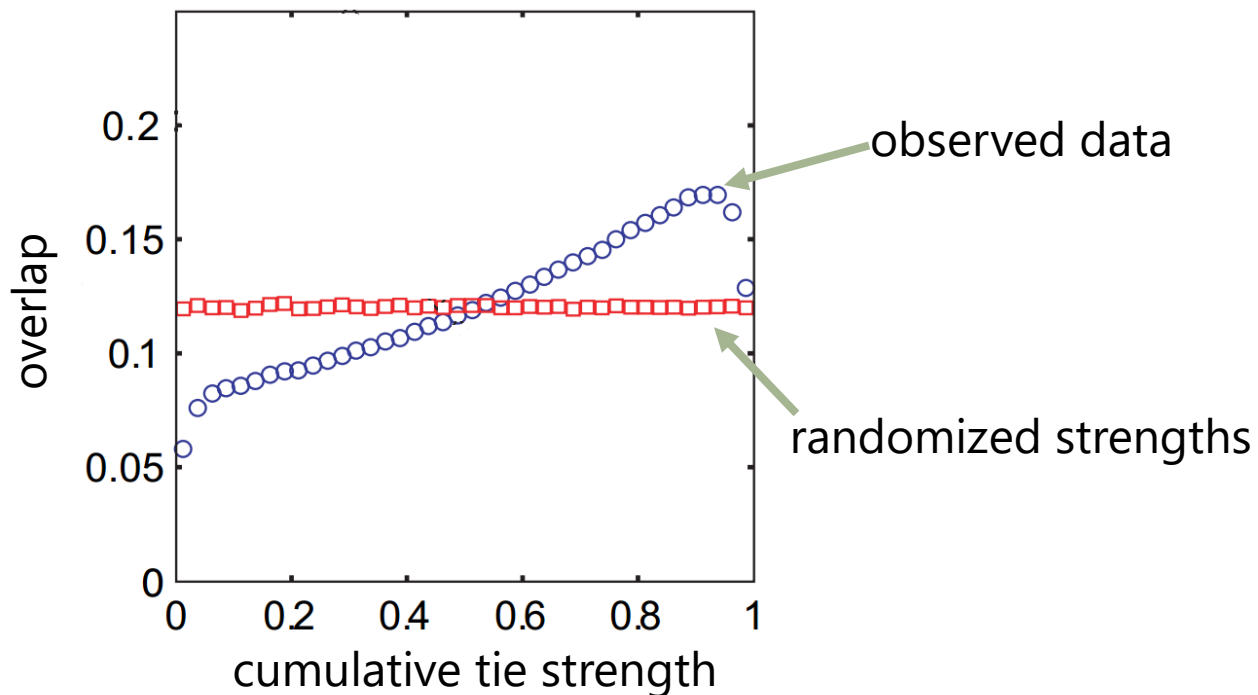
"local bridges" have overlap 0

(picture from Onnela et al., 2007)

Strong & weak ties

Secondly, define the “strength” of a tie in terms of the number of phone calls between i and j

finding: the “stronger” our tie, the more likely there are to be additional ties between our mutual friends



Strong & weak ties

Another case study (Ugander et al., 2012)

Suppose a user receives four e-mail invites to join facebook from users who are already on facebook. Under what conditions are we most likely to accept the invite (and join facebook)?

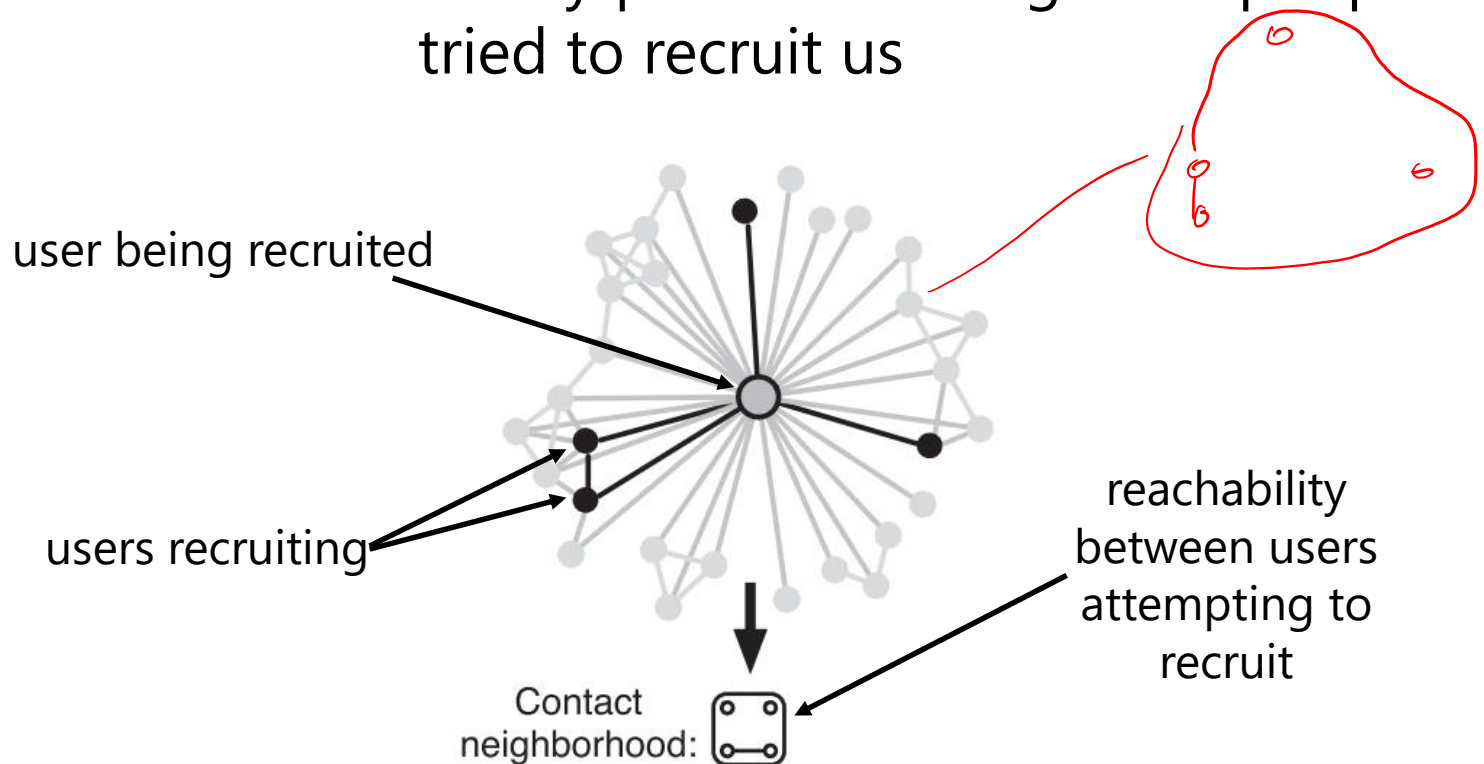
1. If those four invites are from four close friends?
2. If our invites are from four acquaintances?
3. If the invites are from a combination of friends, acquaintances, work colleagues, and family members?

hypothesis: the invitations are most likely to be adopted if they come from **distinct groups** of people in the network

Strong & weak ties

Another case study (Ugander et al., 2012)

Let's consider the connectivity patterns amongst the people who tried to recruit us

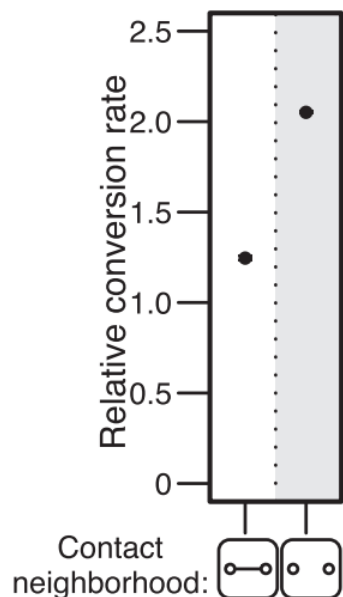


(picture from Ugander et al., 2012)

Strong & weak ties

Another case study (Ugander et al., 2012)

Let's consider the connectivity patterns amongst the people who tried to recruit us

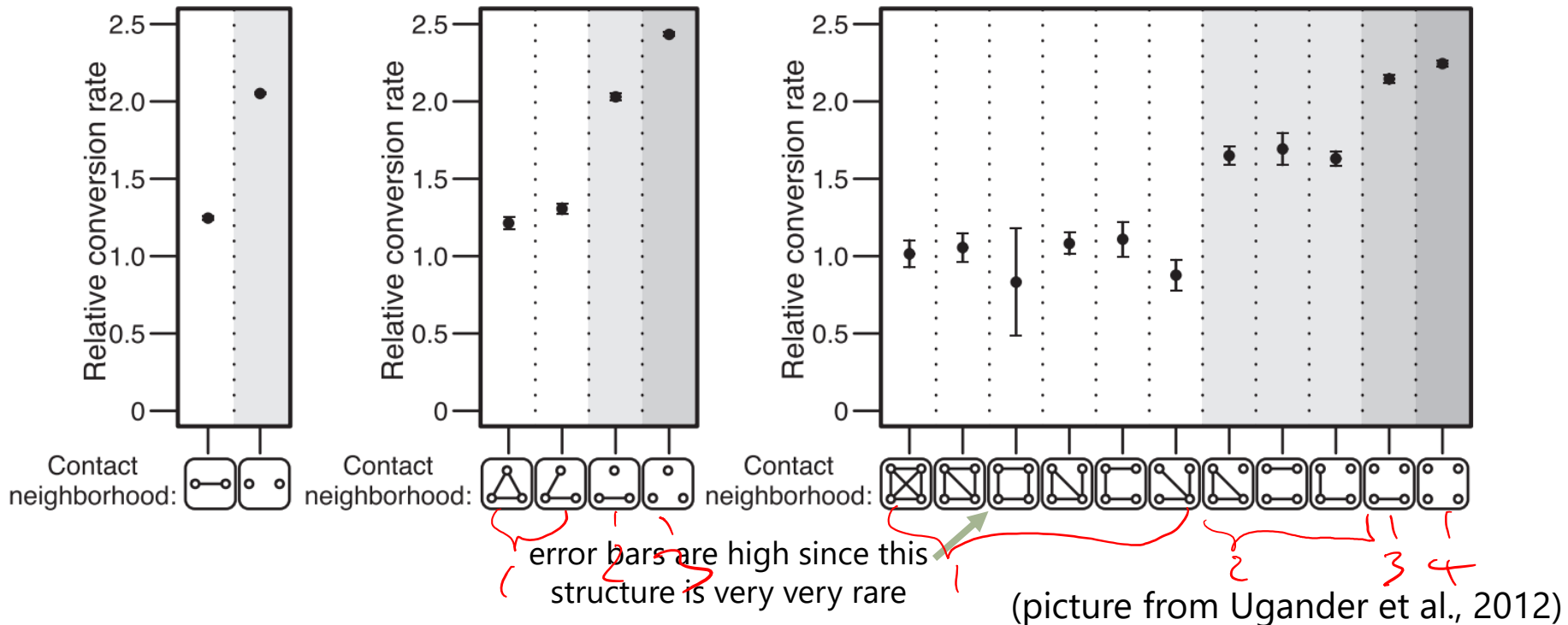


- **Case 1:** two users attempted to recruit
 - **y-axis:** relative to recruitment by a single user
- **finding:** recruitments are **more likely to succeed** if they come from friends who are **not connected to each other**

Strong & weak ties

Another case study (Ugander et al., 2012)

Let's consider the connectivity patterns amongst the people who tried to recruit us



Strong & weak ties

So far:

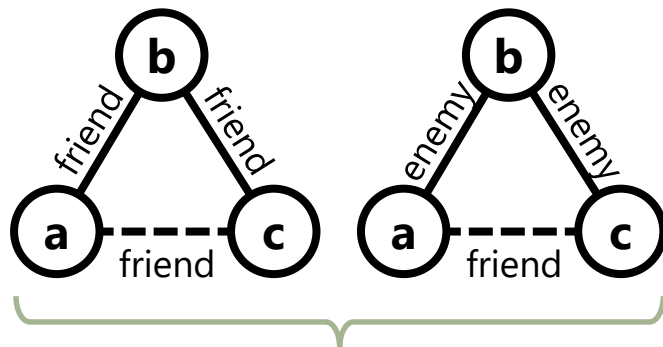
Important aspects of network structure can be explained by the way an edge connects two parts of the network to each other:

- Edges tend to close open triads (clustering coefficient etc.)
- It can be argued that edges that bridge different parts of the network somehow correspond to “weak” connections (Granovetter; Onnela et al.)
- Disconnected parts of the networks (or parts connected by local bridges) expose us to distinct sources of information (Granovetter; Ugander et al.)

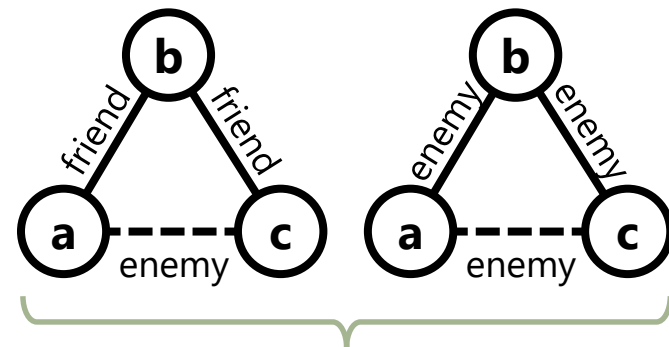
See also...

Structural balance

Some of the assumptions that we've seen today may not hold if edges have **signs** associated with them



balanced: the edge $a \rightarrow c$ is **likely** to form



imbalanced: the edge $a \rightarrow c$ is **unlikely** to form

(see e.g. Heider, 1946)

Questions?

Further reading:

- Easley & Kleinberg, Chapter 3
 - The strength of weak ties
(Granovetter, 1973)
<http://goo.gl/wJVIN>
 - Bearman & Moody

“Suicide and friendships among American adolescents”

http://www.soc.duke.edu/~jmoody77/suicide_ajph.pdf

- Onnela et al.’s mobile phone study

“Structure and tie strengths in mobile communication networks”

http://www.hks.harvard.edu/davidlazer/files/papers/Lazer_PNAS_2007.pdf

- Ugander et al.’s facebook study

“Structural diversity in social contagion”

<file:///C:/Users/julian/Downloads/PNAS-2012-Ugander-5962-6.pdf>

CSE 258 – Lecture 13

Web Mining and Recommender Systems

Small-world phenomena

Small worlds

- We've seen random graph models that reproduce the **power-law** behaviour of real-world networks
- But what about other types of network behaviour, e.g. can we develop a random graph model that reproduces small-world phenomena? Or which have the correct ratio of closed to open triangles?

Six degrees of separation

Another famous study...

- Stanley Milgram wanted to test the (already popular) hypothesis that people in social networks are separated by only a small number of “hops”
 - He conducted the following experiment:

1. “Random” pairs of users were chosen, with start points in Omaha & Wichita, and endpoints in Boston
2. Users at the start point were sent a letter describing the study: they were to get the letter to the endpoint, but only by contacting somebody with whom they had a direct connection
3. So, either they sent the letter directly, or they wrote their name on it and passed it on to somebody they believed had a high likelihood of knowing the target (they also mailed the researchers so that they could track the progress of the letters)



Six degrees of separation

Another famous study...

Of those letters that reached their destination, the average path length was between 5.5 and 6 (thus the origin of the expression). At least two facts about this study are somewhat remarkable:

- First, that short paths appear to be abundant in the network
- Second, that people are capable of discovering them in a “decentralized” fashion, i.e., they’re somehow good at “guessing” which links will be closer to the target

Six degrees of separation

Such small-world phenomena turn out to be abundant in a variety of network settings

e.g. Erdos numbers:

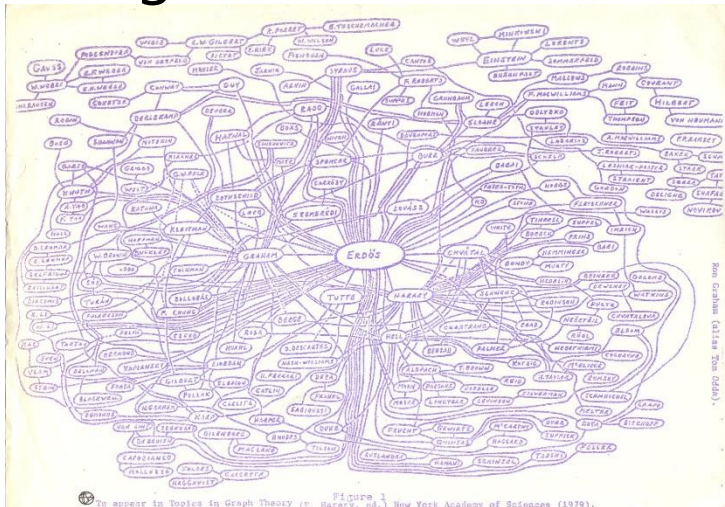


Figure 1
To appear in Topics in Graph Theory (P. Harary, ed.), New York Academy of Sciences (1979).

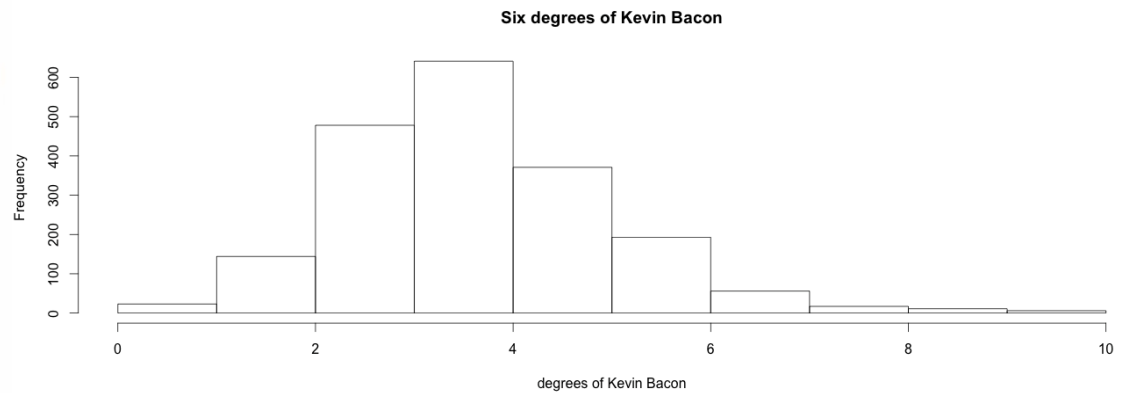
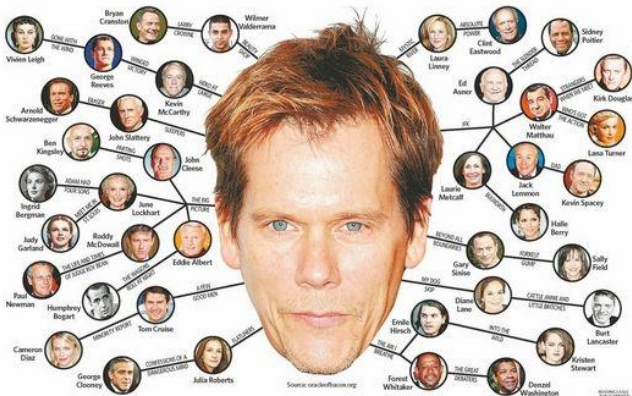
Erdős # 0	-	1 person
Erdős # 1	-	504 people
Erdős # 2	-	6593 people
Erdős # 3	-	33605 people
Erdős # 4	-	83642 people
Erdős # 5	-	87760 people
Erdős # 6	-	40014 people
Erdős # 7	-	11591 people
Erdős # 8	-	3146 people
Erdős # 9	-	819 people
Erdős #10	-	244 people
Erdős #11	-	68 people
Erdős #12	-	23 people
Erdős #13	-	5 people



Six degrees of separation

Such small-world phenomena turn out to be abundant in a variety of network settings

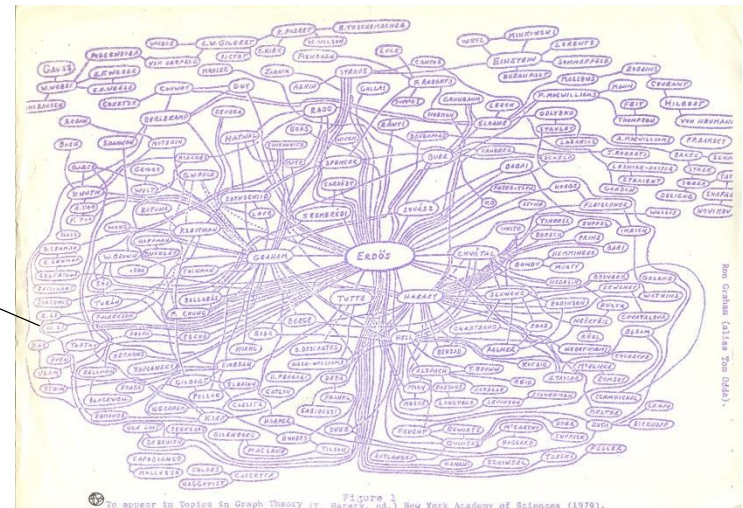
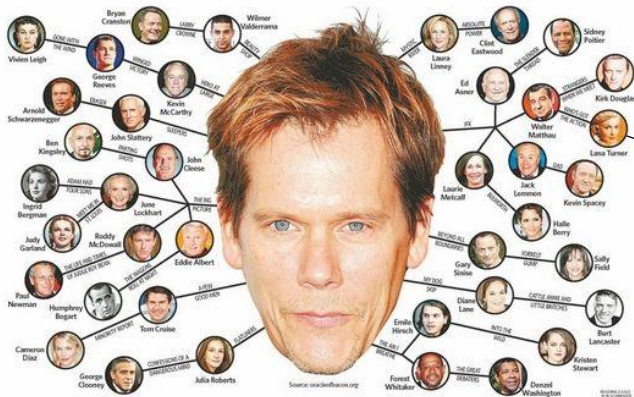
e.g. Bacon numbers:



Six degrees of separation

Such small-world phenomena turn out to be abundant in a variety of network settings

Bacon/Erdos numbers:

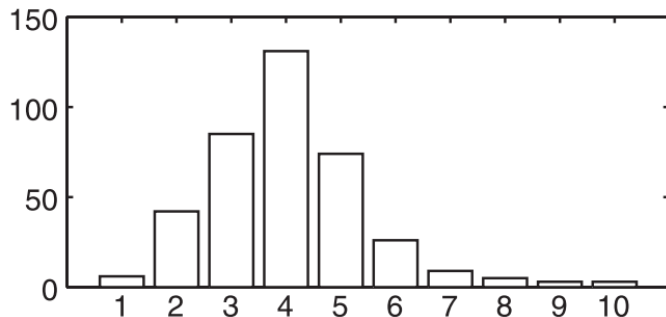


Kevin Bacon→Sarah Michelle Gellar→**Natalie Portman**→Abigail Baird→Michael Gazzaniga→J. Victor→Joseph Gillis→Paul Erdos

Six degrees of separation

Dodds, Muhamed, & Watts repeated Milgram's experiments using e-mail

- 18 "targets" in 13 countries
- 60,000+ participants across 24,133 chains
- Only 384 (!) reached their targets



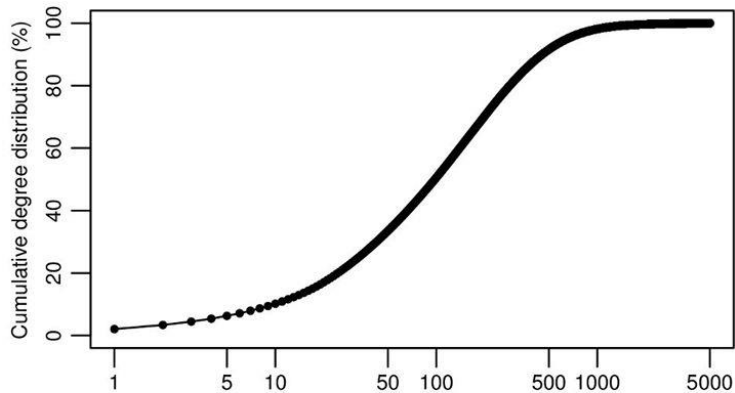
Histogram of (completed) chain lengths – average is just **4.01!**

<i>L</i>	<i>N</i>	Location	Travel	Family	Work	Education	Friends	Cooperative	Other
1	19,718	33	16	11	16	3	9	9	3
2	7,414	40	11	11	19	4	6	7	2
3	2,834	37	8	10	26	6	6	4	3
4	1,014	33	6	7	31	8	5	5	5
5	349	27	3	6	38	12	6	3	5
6	117	21	3	5	42	15	4	5	5
7	37	16	3	3	46	19	8	5	0

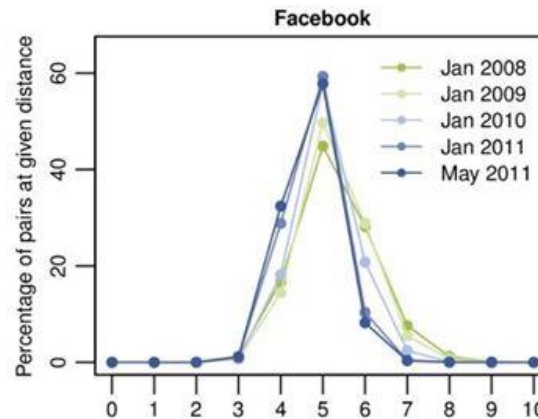
Reasons for choosing the next recipient at each point in the chain

Six degrees of separation

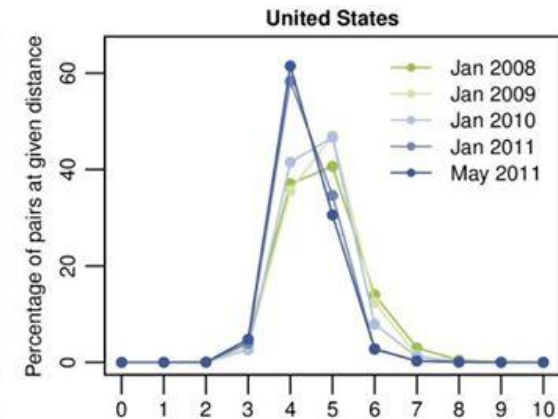
Actual shortest-path distances are similar to those in Dodds' experiment:



Cumulative degree distribution (# of friends) of Facebook users



Hop distance between Facebook users



Hop distance between users in the US

This suggests that people choose a reasonably good heuristic when choosing shortest paths in a decentralized fashion (assuming that FB is a good proxy for "real" social networks)

Six degrees of separation

Q: is this result surprising?

- **Maybe not:** We have ~ 100 friends on Facebook, so 100^2 friends-of-friends, 10^6 at length three, 10^8 at length four, **everyone** at length 5
- **But:** Due to our previous argument that people close triads, the **vast majority** of new links will be between friends of friends (i.e., we're increasing the **density** of our local network, rather than making distant links more reachable)
- In fact **92%** of new connections on Facebook are to a friend of a friend (Backstrom & Leskovec, 2011)

Six degrees of separation

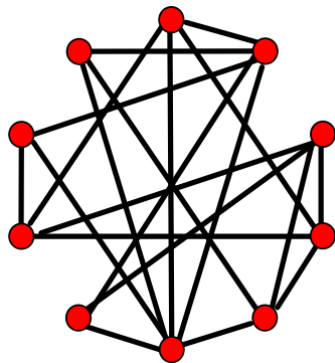
Definition: Network diameter

- A network's diameter is the length of its **longest shortest path**
- **Note:** iterating over all pairs of nodes i and j and then running a shortest-paths algorithm is going to be prohibitively slow
- Instead, the "all pairs shortest paths" algorithm computes all shortest paths simultaneously, and is more efficient ($O(N^2 \log N)$ to $O(N^3)$, depending on the graph structure)
- In practice, one doesn't **really** care about the diameter, but rather the distribution of shortest path lengths, e.g., what is the average/90th percentile shortest-path distance
- This latter quantity can be computed just by randomly sampling pairs of nodes and computing their distance
 - When we say that a network exhibits the "small world phenomenon", we are really saying this latter quantity is small

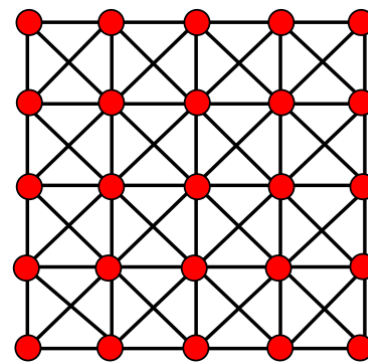
Six degrees of separation

Q: is this a contradiction?

- How can we have a network made up of **dense communities** that is simultaneously a **small world**?
- The shortest paths we could possibly have are $O(\log n)$ (assuming nodes have constant degree)



random connectivity –
low diameter, low
clustering coefficient

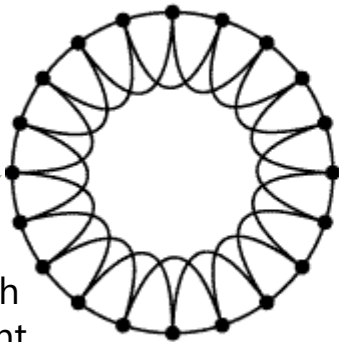


regular lattice – high
clustering coefficient,
high diameter

Six degrees of separation

We'd like a model that reproduces small-world phenomena

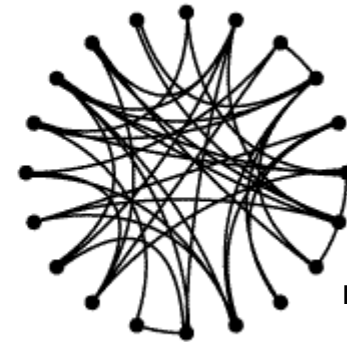
Regular



regular lattice – high clustering coefficient, high diameter

We'd like something "in between" that exhibits both of the desired properties (high cc, low diameter)

Random



random connectivity – low diameter, low clustering coefficient

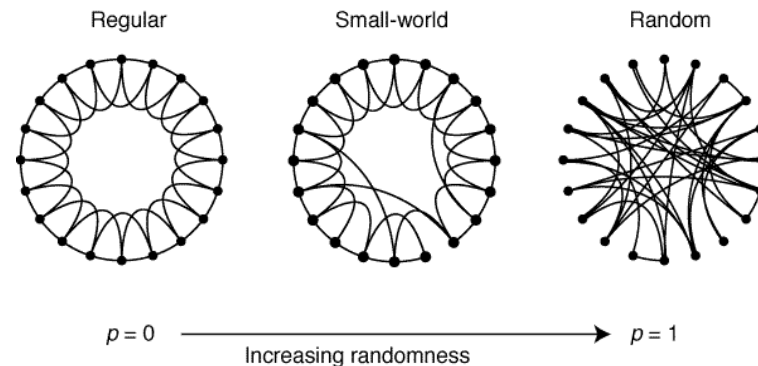
Six degrees of separation

The following model was proposed by Watts & Strogatz (1998)

1. Start with a regular lattice graph (which we know to have high clustering coefficient)

Next – introduce some randomness into the graph

2. For each edge, with prob. p , reconnect one of its endpoints

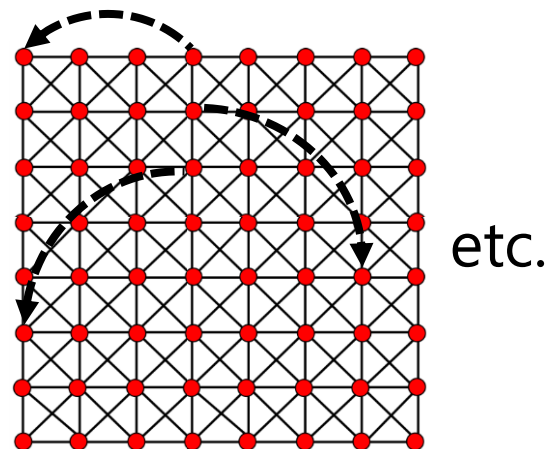


as we increase p , this becomes more like a random graph

Six degrees of separation

Slightly simpler (to reason about formulation) with the same properties

1. Start with a regular lattice graph (which we know to have high clustering coefficient)
2. From each node, add an additional random link



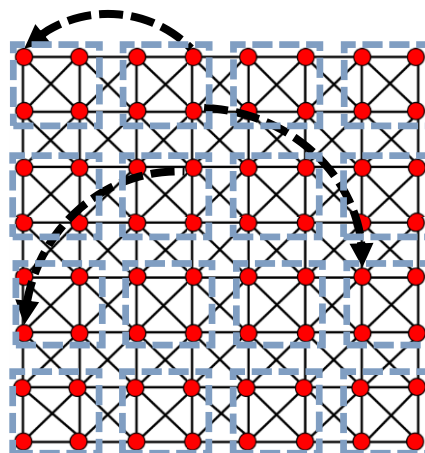
Six degrees of separation

Slightly simpler (to reason about formulation) with the same properties

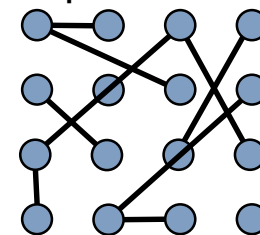
Conceptually, if we combine groups of adjacent nodes into “supernodes”, then what we have formed is a **4-regular** random graph

(very handwavy) proof:

- The clustering coefficient is still high (each node is incident to 12 triangles)
- **4-regular** random graphs have diameter $O(\log n)$ (Bollobas, 2001), so the whole graph has diameter $O(\log n)$



connections between supernodes:



(should be a 4-regular random graph, I didn't finish drawing the edges)

Six degrees of separation

The Watts-Strogatz model

- Helps us to understand the relationship between dense clustering and the small-world phenomenon
- Reproduces the small-world structure of realistic networks
- Does **not** lead to the correct degree distribution (no power laws)

(see Klemm, 2002: "Growing scale-free networks with small-world behavior" <http://ifisc.uib-csic.es/victor/Nets/sw.pdf>)

Six degrees of separation

So far...

- Real networks exhibit **small-world** phenomena: the average distance between nodes grows only logarithmically with the size of the network
- Many experiments have demonstrated this to be true, in mail networks, e-mail networks, and on Facebook etc.
- But we know that social networks are highly **clustered** which is somehow inconsistent with the notion of having low diameter
- To explain this apparent contradiction, we can model networks as some combination of highly-clustered nodes, plus some fraction of “random” connections

Questions?

Further reading:

- Easley & Kleinberg, Chapter 20
 - Milgram's paper

"An experimental study of the small world problem"

<http://www.uvm.edu/~pdodds/files/papers/others/1969/travers1969.pdf>

- Dodds et al.'s small worlds paper

<http://www.cis.upenn.edu/~mkearns/teaching/NetworkedLife/columbia.pdf>

- Facebook's small worlds paper

<http://arxiv.org/abs/1111.4503>

- Watts & Strogatz small worlds model

"Collective dynamics of 'small world' networks"

[file:///C:/Users/julian/Downloads/w s NATURE 0.pdf](file:///C:/Users/julian/Downloads/w_s_NATURE_0.pdf)

- More about random graphs

"Random Graphs" (Bollobas, 2001), Cambridge University Press

CSE 258 – Lecture 13

Web Mining and Recommender Systems

Hubs and Authorities; PageRank

Trust in networks

We already know that there's considerable variation in the connectivity structure of nodes in networks

So how can we find nodes that are in some sense "important" or "authoritative"?

- In links?
- Out links?
- Quality of content?
- Quality of linking pages?
 - etc.

Trust in networks

1. The "HITS" algorithm

Two important notions:

Hubs:

We might consider a node to be of "high quality" if it links to many high-quality nodes. E.g. a high-quality page might be a "hub" for good content
(e.g. Wikipedia lists)

Authorities:

We might consider a node to be of high quality if many high-quality nodes link to it
(e.g. the homepage of a popular newspaper)

Trust in networks

This “self-reinforcing” notion is the idea behind the HITS algorithm

- Each node i has a “hub” score h_i
- Each node i has an “authority” score a_i
- The hub score of a page is the sum of the authority scores of pages it links to
- The authority score of a page is the sum of hub scores of pages that link to it

Trust in networks

This “self-reinforcing” notion is the idea behind the HITS algorithm

Algorithm:

$$a_i^{(0)} = \frac{1}{\sqrt{n}} \quad h_i^{(0)} = \frac{1}{\sqrt{n}}$$

$A = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{bmatrix}$

iterate until convergence:

$$\forall_i a_i^{(t+1)} = \sum_{j \rightarrow i} h_j^{(t)}$$

pages that link to i

$$\forall_i h_i^{(t+1)} = \sum_{i \rightarrow j} a_j^{(t)}$$

pages that i links to

normalize:

$$\|a^{(t+1)}\|_2^2 = 1 \quad \|h^{(t+1)}\|_2^2 = 1$$

Trust in networks

This “self-reinforcing” notion is the idea behind the HITS algorithm

This can be re-written in terms of the adjacency matrix (A)

$$a_i^{(0)} = \frac{1}{\sqrt{n}} \quad h_i^{(0)} = \frac{1}{\sqrt{n}}$$

iterate until convergence:

$$\begin{array}{ll} a^{(t+1)} = A^T h^{(t)} & \\ h^{(t+1)} = A a^{(t)} & \end{array} \quad \begin{array}{l} \text{skipping} \\ \text{a step:} \end{array} \quad \begin{array}{l} a^{(t+2)} = (A^T A)^t a^{(t)} \\ h^{(t+2)} = (A A^T)^t h^{(t)} \end{array}$$

normalize:

$$\|a^{(t+1)}\|_2^2 = 1 \quad \|h^{(t+1)}\|_2^2 = 1$$

Trust in networks

This “self-reinforcing” notion is the idea behind the HITS algorithm

So at convergence we seek stationary points such that

$$A^T A a = c' \cdot a$$

$$A A^T h = c'' \cdot h$$

(constants don't matter since we're normalizing)

- This can only be true if the authority/hub scores are **eigenvectors** of $A^T A$ and $A A^T$
- In fact this will converge to the eigenvector with the largest eigenvalue (see: Perron-Frobenius theorem)

Trust in networks

The idea behind PageRank is very similar:

- Every page gets to “vote” on other pages
- Each page’s votes are proportional to that page’s importance
- If a page of importance x has n outgoing links, then each of its votes is worth x/n
- Similar to the previous algorithm, but with only a single term to be updated (the rank r_i of a page i)

$$\forall_i r_i^{(t+1)} = \sum_{j \rightarrow i} \frac{r_j^{(t)}}{|\Gamma(j)|}$$

rank of linking pages

of links from linking pages

Trust in networks

The idea behind PageRank is very similar:

Matrix formulation:

each **column** describes the out-links of one page, e.g.:

pages

$$M = \begin{pmatrix} \frac{1}{3} & 0 & \frac{1}{4} & 1 \\ 0 & 0 & \frac{1}{4} & 0 \\ \frac{1}{3} & \frac{1}{2} & \frac{1}{4} & 0 \\ \frac{1}{3} & \frac{1}{2} & \frac{1}{4} & 0 \end{pmatrix} \text{ pages}$$

this out-link gets 1/3 votes since this page has three out-links

column-stochastic matrix (columns add to 1)

Trust in networks

The idea behind PageRank is very similar:

Then the update equations become:

$$r^{(t+1)} = Mr^{(t)}$$

And as before the stationary point is given by the eigenvector of M with the highest eigenvalue

Summary

The level of “authoritativeness” of a node in a network should somehow be defined in terms of the pages that link to (it or the pages it links from), and *their* level of authoritativeness

- Both the HITS algorithm and PageRank are based on this type of “self-reinforcing” notion
 - We can then measure the centrality of nodes by some iterative update scheme which converges to a stationary point of this recursive definition
- In both cases, a solution was found by taking the principal eigenvector of some matrix encoding the link structure

Trust in networks

This week

- We've seen how to characterize networks by their degree distribution (degree distributions in many real-world networks follow power laws)
- We've seen some random graph models that try to mimic the degree distributions of real networks
- We've discussed the notion of "tie strength" in networks, and shown that edges are likely to form in "open" triads
 - We've seen that real-world networks often have small diameter, and exhibit "small-world" phenomena
- We've seen (very quickly) two algorithms for measuring the "trustworthiness" or "authoritativeness" of nodes in networks

Questions?

Further reading:

- Easley & Kleinberg, Chapter 14
- The “HITS” algorithm (aka “Hubs and Authorities”)
“Hubs, authorities, and communities” (Kleinberg,
1999)

http://cs.brown.edu/memex/ACM_HypertextTestbed/papers/10.html

CSE 258 – Lecture 13

Web Mining and Recommender Systems

Some midterm Qs

Midterm Qs

Section 1: Regression and Ranking (7 marks)

Unless specified otherwise the questions in this section are each worth **1 mark**.

The following is a list of Vin Diesel's recent films:

No.	Title	Year	IMDB score	MPAA rating	length in minutes	classifier score
1	XXX: The Return of Xander Cage	2017	5.6	PG-13	110	-0.5
2	Billy Lynn's Long Halftime Walk	2016	6.6	R	113	-2.1
3	The Last Witch Hunter	2015	6.3	PG-13	106	-3.2
4	Furious 7	2015	7.4	PG-13	137	4.8
5	Guardians of the Galaxy	2014	8.1	PG-13	121	2.2
6	Riddick	2013	6.4	R	119	-1.2
7	Fast & Furious 6	2013	7.2	PG-13	130	-0.8
8	Fast Five	2011	7.3	PG-13	131	1.2
9	Fast & Furious	2009	6.6	PG-13	107	0.1
10	The Fast and the Furious: Tokyo Drift	2006	6.0	PG-13	104	-0.3



1. Suppose you train a regressor of the following form to predict the IMDB score:

$$\text{IMDB score} \simeq \theta_0 + \theta_1[\text{'Fast' in title}] + \theta_2[\text{'R' rated}] + \theta_3[\text{length in minutes}]$$

relevant

"retrained"

What would be the feature representation of the first two movies?

1:
2:

Midterm Qs

2. After training the above regressor you obtain $\theta = (1.5, 0.05, -0.25, 0.05)$. What would you predict would be the IMDB score of Vin Diesel's next film, *The Fate of The Furious* (released 2017, PG-13, 140 minutes long). You can write down an expression rather than the exact value:

A:

Next, you train a Support Vector Machine to predict the binary outcome 'IMDB score ≥ 7.0 ' using the same features. Suppose the classifier produces the scores $(X_i \cdot \theta)$ shown in the right column of the table.

3. What is the accuracy, and the Balanced Error Rate of this classifier?

A:

4. What is the classifier's precision and recall?

A:

vs. $pr@k$
 $rec@k$

5. Perhaps you would like to add the 'year' variable to your classifier. Assuming a simple model with no regularizer, show that using the feature [year] is equivalent to using the feature [year - 2006].

A:

$$\begin{aligned} & \theta_0 + \theta_1 (\text{year} - 2006) \\ = & (\theta_0 - \theta_1 \cdot 2006) + \theta_1 \cdot \text{year} \\ = & \theta'_0 + \theta_1 \cdot \text{year} \end{aligned}$$

Midterm Qs

6. (Hard) Briefly explain why these two representations would *not* be equivalent when training a model with a regularizer (e.g. $\|\theta\|_2^2$) (2 marks).

A:

- representation changes param. values
- But, regularizer is sensitive to changes in relative param. scale

Midterm Qs

Section 2: Classification and Diagnostics (6 marks)

Each of the following questions is worth 2 marks.

Suppose you are trying to train a classifier to detect whether bicycle frames will catastrophically fail during five years of use (i.e., $y_i = 1$ if the i^{th} bicycle fails). You build a dataset containing features about 10,000 bicycle frames manufactured between 1978 and 2012 (weight, material, manufacture year, etc.) and whether or not they failed. You partition the dataset into a training and validation set using a 50%/50% split.

After trying several different classifiers on your data and measuring their error (percentage of incorrect classifications), you obtain some unexpected results. Briefly explain a possible reason for the results and suggest a possible solution.

7. You obtain low (1%) error on your training set but even *lower* error (0.5%) on your validation set

Diagnosis:	Non random splits
Solution:	random splits

8. You build an accurate classifier with only 1% error on your training set, and 1.5% error on your validation set. However, when you deploy the system, it fails to identify any instances of catastrophic failure

Diagnosis:	Predicts "F" always
Solution:	measure recall / F-score

Midterm Qs

9. Suppose that to fix the above issue, you want to adjust a logistic regression-based classifier so that it gives 100 times as much weight to false negatives (bicycles that fail but were predicted not to) as false positives. How would you adjust the objective to achieve this? Recall that the original objective for logistic regression is

$$\log \sum_{y_i=1} \log \sigma(X_i \cdot \theta) + \sum_{y_i=0} \log(1 - \sigma(X_i \cdot \theta))$$

A:

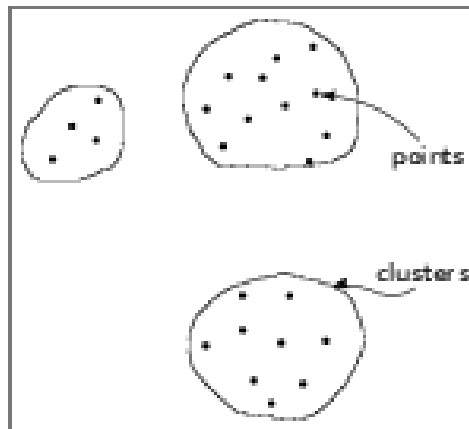
labeled T,
predicted F

Midterm Qs

Section 3: Clustering / Communities (8 marks)

The following questions are concerned with the *K-means* algorithm (see pseudocode at the end of the exam). Each question is worth **2 marks**.

When asked to draw examples, provide 2-d sets of points and clusters like the following:



10. Explain why the algorithm provided in the pseudocode will eventually converge (i.e., terminate).

A:

— each step reduces MSE
— finitely many steps

Midterm Qs

11. The K-Means algorithm will in general converge to a local optimum rather than a global one. Draw a simple 2-d example, containing a set of points and clusters, such that the solution is *not* optimal but for which the algorithm would not make further progress.



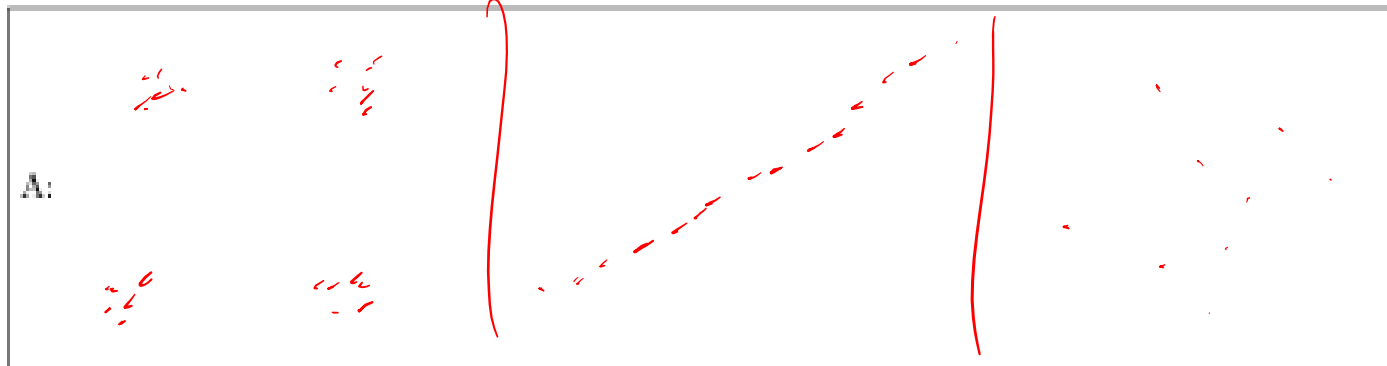
12. Suggest simple modifications to the k-means algorithm that might increase its chances of finding a good solution.

A:

- random restarts
- re-init. clusters that became empty

Midterm Qs

13. Compared to PCA, K-means will work well for different types of clusters. Give three examples of 2-d clustered data where (a) K-means will perform *better* than PCA (in terms of reconstruction error); (b) K-means will perform *worse* than PCA; and (c) K-means and PCA will both perform poorly.



Midterm Qs

Section 4: Recommender Systems (5 marks)

Unless specified otherwise the questions in this section are each worth 1 mark.

On a popular movie streaming website, a few users have watched the following recent movies:

Movie	Watched?				Rated?			
	Caroline	Mengting	Ruining	Zachary	Caroline	Mengting	Ruining	Zachary
<i>XXX: Return of Xander Cage</i>	1	1	0	1	5	1		2
<i>La La Land</i>	1	1	1	1	5	2	2	2
<i>John Wick 2</i>	1	1	1	0	4	2	1	
<i>Rogue One</i>	0	0	1	1			5	1
<i>Resident Evil</i>	1	0	0	1	4			1

14. Using 'watched' data: Which two users are *most similar* in terms of their Jaccard similarity (write down all pairs in case of a tie)?

A:

15. Which two *items* are most similar in terms of their Jaccard similarity?

A:

16. Which two users are most similar in terms of their *ratings*, based on their Pearson correlation (defined below)?

A:

CFZ

Midterm Qs

17. (Hard) Show that a latent factor model of the form

$$\text{rating}(u, i) = \alpha + \beta_u + \beta_i$$

is linear in its parameters ($\theta = (\alpha; \beta_u; \beta_i)$) but that a model of the form

$$\text{rating}(u, i) = \alpha + \beta_u + \beta_i + \gamma_u \cdot \gamma_i$$

is not linear in its parameters ($\theta = (\alpha; \beta_u; \beta_i; \gamma_u; \gamma_i)$) (recall the definition of linearity: $r_{\theta_1 + \theta_2}(u, i) = r_{\theta_1}(u, i) + r_{\theta_2}(u, i)$) (2 marks).

$$\begin{aligned} & (\alpha + \alpha') + (\beta_u + \beta_u') + (\beta_i + \beta_i') \\ &= (\alpha + \beta_u + \beta_i) + (\alpha' + \beta_u' + \beta_i') \end{aligned}$$

A:

$$\begin{aligned} & (\gamma_u + \gamma_u') \cdot (\gamma_i + \gamma_i') \\ & \neq (\gamma_u \cdot \gamma_i) + (\gamma_u' \cdot \gamma_i') \end{aligned}$$