

Social and Information Networks

University of Toronto CSC303
Winter/Spring 2021

Week 6: February 22-26 (2021)

Mon. Feb 22: Announcements and Corrections

- For those of you who were wondering, yes, Dr. Stanley Milgram was the American sociologist who created both the Small World experiment we discussed, and the eponymous Milgram experiment on obedience. Aptly enough, it's (metaphorically) a small world ;)

Mon. Feb 22: Announcements and Corrections

- For those of you who were wondering, yes, Dr. Stanley Milgram was the American sociologist who created both the Small World experiment we discussed, and the eponymous Milgram experiment on obedience. Aptly enough, it's (metaphorically) a small world ;)
- The Midterm will be from March 12th to March 14th

Mon. Feb 22: Announcements and Corrections

- For those of you who were wondering, yes, Dr. Stanley Milgram was the American sociologist who created both the Small World experiment we discussed, and the eponymous Milgram experiment on obedience. Aptly enough, it's (metaphorically) a small world ;)
- The Midterm will be from March 12th to March 14th
 - ▶ Open book

Mon. Feb 22: Announcements and Corrections

- For those of you who were wondering, yes, Dr. Stanley Milgram was the American sociologist who created both the Small World experiment we discussed, and the eponymous Milgram experiment on obedience. Aptly enough, it's (metaphorically) a small world ;)
- The Midterm will be from March 12th to March 14th
 - ▶ Open book
 - ▶ Take-home

Mon. Feb 22: Announcements and Corrections

- For those of you who were wondering, yes, Dr. Stanley Milgram was the American sociologist who created both the Small World experiment we discussed, and the eponymous Milgram experiment on obedience. Aptly enough, it's (metaphorically) a small world ;)
- The Midterm will be from March 12th to March 14th
 - ▶ Open book
 - ▶ Take-home
 - ▶ No time limit other than submission by March 14th at 11:50PM, Toronto time

Mon. Feb 22: Announcements and Corrections

- For those of you who were wondering, yes, Dr. Stanley Milgram was the American sociologist who created both the Small World experiment we discussed, and the eponymous Milgram experiment on obedience. Aptly enough, it's (metaphorically) a small world ;)
- The Midterm will be from March 12th to March 14th
 - ▶ Open book
 - ▶ Take-home
 - ▶ No time limit other than submission by March 14th at 11:50PM, Toronto time
 - ▶ It will cover everything we've seen so far, plus power-laws and web-search (i.e., Ch 14 & 18 of E&K)

Mon. Feb 22: Announcements and Corrections

- For those of you who were wondering, yes, Dr. Stanley Milgram was the American sociologist who created both the Small World experiment we discussed, and the eponymous Milgram experiment on obedience. Aptly enough, it's (metaphorically) a small world ;)
- The Midterm will be from March 12th to March 14th
 - ▶ Open book
 - ▶ Take-home
 - ▶ No time limit other than submission by March 14th at 11:50PM, Toronto time
 - ▶ It will cover everything we've seen so far, plus power-laws and web-search (i.e., Ch 14 & 18 of E&K)
 - ▶ Further details to follow

Roadmap: where we have been and whats next

Chapter 20 started off with a discussion of the small worlds phenomena and an insightful understanding of how decentralized search can work.

Previously, we were led to the observation that geographical distance (or social distance) correlates with friendship such that $\text{Prob}[v \text{ is a friend of } u] \approx [\text{rank}_u(v)]^{-1}$.

Furthermore, there is a sense that long distance friendships are “rare”.

Roadmap: where we have been and whats next

Chapter 20 started off with a discussion of the small worlds phenomena and an insightful understanding of how decentralized search can work.

Previously, we were led to the observation that geographical distance (or social distance) correlates with friendship such that $\text{Prob}[v \text{ is a friend of } u] \approx [\text{rank}_u(v)]^{-1}$.

Furthermore, there is a sense that long distance friendships are “rare”.

We even saw a claim (by Oscar Sandberg) that decentralized search might implicitly be a partial explanation of network dynamics and structure

Roadmap: where we have been and whats next

Chapter 20 started off with a discussion of the small worlds phenomena and an insightful understanding of how decentralized search can work.

Previously, we were led to the observation that geographical distance (or social distance) correlates with friendship such that $\text{Prob}[v \text{ is a friend of } u] \approx [\text{rank}_u(v)]^{-1}$.

Furthermore, there is a sense that long distance friendships are “rare”.

We even saw a claim (by Oscar Sandberg) that decentralized search might implicitly be a partial explanation of network dynamics and structure

We have seen earlier (Chapters 3 and 4) how selection (i.e. homophily in the sense of “birds of a feather flock together”) causes friendship links. Chapter 5 also relates to how links can form or change to achieve structural balance.

Roadmap: where we have been and whats next

Chapter 20 started off with a discussion of the small worlds phenomena and an insightful understanding of how decentralized search can work.

Previously, we were led to the observation that geographical distance (or social distance) correlates with friendship such that $\text{Prob}[v \text{ is a friend of } u] \approx [\text{rank}_u(v)]^{-1}$.

Furthermore, there is a sense that long distance friendships are “rare”.

We even saw a claim (by Oscar Sandberg) that decentralized search might implicitly be a partial explanation of network dynamics and structure

We have seen earlier (Chapters 3 and 4) how selection (i.e. homophily in the sense of “birds of a feather flock together”) causes friendship links. Chapter 5 also relates to how links can form or change to achieve structural balance.

This week we will be building on these ideas.

Power law distributions

A *power law distribution* for discrete random variable X satisfies $\text{Prob}[X = k] \approx \frac{a}{k^c}$ for some constants a and c . (We often just focus on on the exponent c .)

Power law distributions

A *power law distribution* for discrete random variable X satisfies $\text{Prob}[X = k] \approx \frac{a}{k^c}$ for some constants a and c . (We often just focus on on the exponent c .)

Closely related (and sometimes used interchangeably) is Zipf's Law, which relates the frequency, f (i.e. count) of something with it's rank, r .

- $r = 1$ being the most frequent, $r = 2$ being the second most frequent, and so on.

Power law distributions

A *power law distribution* for discrete random variable X satisfies $\text{Prob}[X = k] \approx \frac{a}{k^c}$ for some constants a and c . (We often just focus on on the exponent c .)

Closely related (and sometimes used interchangeably) is Zipf's Law, which relates the frequency, f (i.e. count) of something with it's rank, r .

- $r = 1$ being the most frequent, $r = 2$ being the second most frequent, and so on.

A phenomena satisfies Zipf's Law when:

$$f \approx \frac{a}{r^c}$$

for some constants a and c

Why care about power laws?

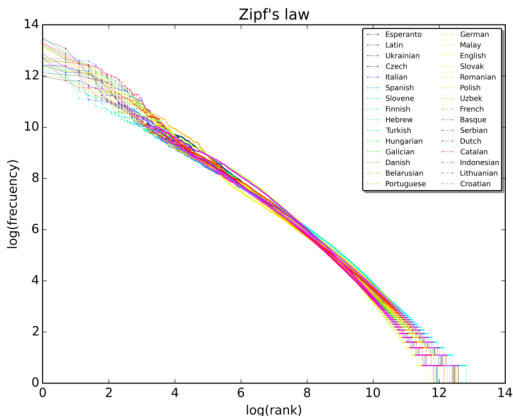
- Chapter 18 calls attention to the fact that power law distributions often occur in network and natural phenomena.
- We've already seen power laws emerge in the probability of friendship forming with respect to both distances, and ranks

Why care about power laws?

- Chapter 18 calls attention to the fact that power law distributions often occur in network and natural phenomena.
- We've already seen power laws emerge in the probability of friendship forming with respect to both distances, and ranks
- Where else do they appear?

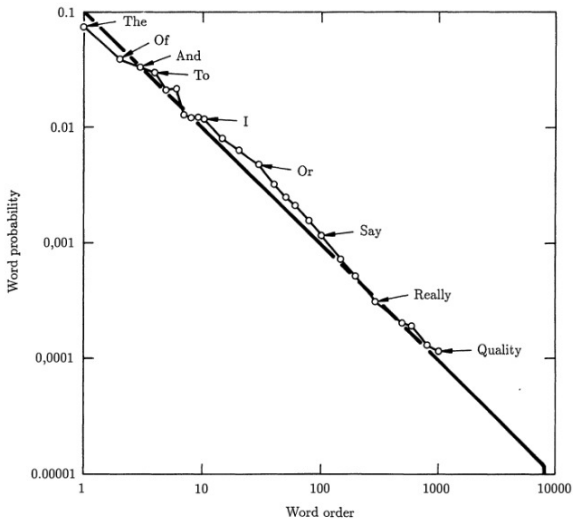
Zipf's law in text

- In any book, let w be a word within. If we calculate it's frequency f_w (i.e. number of occurrences in the text) and it's rank r_w (i.e. is the word the first, second, ... n th most common word in the book) then we find that $f_w \propto 1/r_w$ (or equivalently, $\log f_w \approx -\log r_w + C$)



[Image By SergioJimenez - Own work, CC BY-SA 4.0, link]

Power law in text



[Image from Fractals, Chaos, Power Laws: Minutes from an Infinite Paradise by Manfred Schroeder]

Zipf's law in population centres

- Similarly, if we consider cities t and let f_t be the population of a city, and r_t be the city's rank by population, we have $f_t \approx a/r_t$

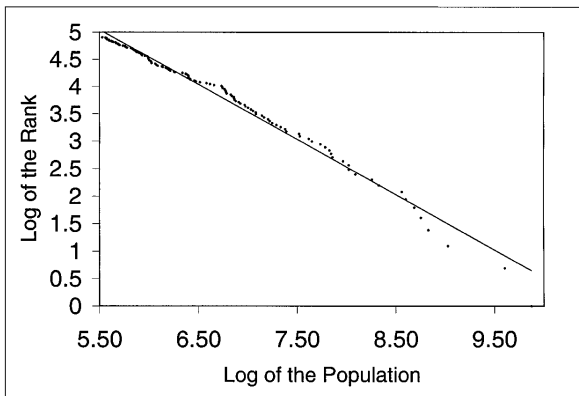


FIGURE I

Log Size versus Log Rank of the 135 largest U. S. Metropolitan Areas in 1991
Source: Statistical Abstract of the United States [1993].

[Image from Gabaix, 1999]

Power laws in websites and products

- Power laws also arise in the popularity of websites and commercial products

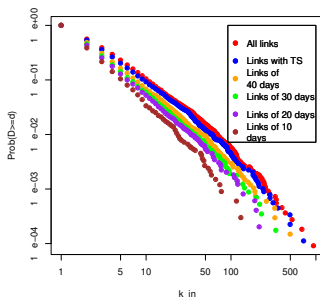


Fig. 6: Temporal changes in the in-degree distributions in TREC.

[Image from Shi et al.]

- Empirically, in the web network (i.e. an *information network*), the probability that a site will have k in-links is proportional to k^{-2} . (More precisely, proportional to $k^{-(2+\epsilon)}$ for some $\epsilon > 0$.)

Power law distributions

- Power law distributions in social and information networks often arise from coupled or correlated individual decisions.
- For example, the *popularity* of certain books or cities, occurrences of specific words in a natural language, etc. More specifically, we'll be considering the frequency of in-links to web sites.
- Events may be less rare than they appear at first glance

Power law distributions

- Power law distributions in social and information networks often arise from coupled or correlated individual decisions.
- For example, the *popularity* of certain books or cities, occurrences of specific words in a natural language, etc. More specifically, we'll be considering the frequency of in-links to web sites.
- Events may be less rare than they appear at first glance
- Key takeaway: extreme events (e.g., for a site to have very many in-links) is **not so rare** when compared with what would be predicted by independent decisions.

How rare is rare when compared with averages over independent actions?

- What if people chose where to live independent of the city? What would be (the distribution of) the population of cities?
- What if we all independently chose to read books not dependent on current events or what friends (or an online system) recommended? How rare would it be to have a huge best seller?
- What, if each web site chose their out-links independently and without some underlying dynamics to guide the process?

As is well understood, the **Central Limit Theorem** tells us that “a quantity that can be viewed as the sum (or average) of many small independent random effects will be well-approximated” by a *normal distribution*.

The normal distribution

The normal or Gaussian distribution has the following probability density function:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

As we know, normal distributions have a *bell shaped curve*.

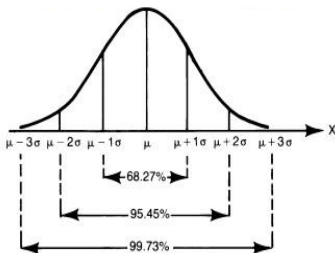


Figure 2

Percent	99.73%	99%	95.45%	95%	90%	80%	68.27%
No. of $\pm \sigma$'s	3.00	2.58	2.00	1.96	1.645	1.28	1.00

From: <http://www.answers.com/topic/normal-distribution>

So how rare is rare?

For a normal distribution, the probability that an outlier (i.e. an exceptional event) will occur decreases exponentially (with distance from the mean). In particular, if say in-links followed a normal distribution, then the probability that a given site would have k links would decrease *exponentially* in k . Very small or large “outliers” would be highly improbable.

So how rare is rare?

For a normal distribution, the probability that an outlier (i.e. an exceptional event) will occur decreases exponentially (with distance from the mean). In particular, if say in-links followed a normal distribution, then the probability that a given site would have k links would decrease *exponentially* in k . Very small or large “outliers” would be highly improbable.

So how rare is rare when for example we compare $\text{Prob}[k \text{ in-links}] \approx k^{-2}$ in comparison to $\text{Prob}[k \text{ in-links}] \approx 2^{-k}$?

So how rare is rare?

For a normal distribution, the probability that an outlier (i.e. an exceptional event) will occur decreases exponentially (with distance from the mean). In particular, if say in-links followed a normal distribution, then the probability that a given site would have k links would decrease *exponentially* in k . Very small or large “outliers” would be highly improbable.

So how rare is rare when for example we compare $\text{Prob}[k \text{ in-links}] \approx k^{-2}$ in comparison to $\text{Prob}[k \text{ in-links}] \approx 2^{-k}$?

For say $k = 30$, $2^{-30} \approx 1/10^9$ whereas $(30)^{-2} = 1/900$.
One in a billion vs better than 1 in a 1000.

So where are we going?

As we have mentioned before, one of the most fundamental questions for social networks concerns how they evolve. What is the interplay between selection and influence?

This is a difficult question. Perhaps the dynamics of information networks created by individuals can be better understood than the dynamics of friendships, political affiliations, opinion formation, etc.

We will see a network dynamic that leads to a power law distribution.

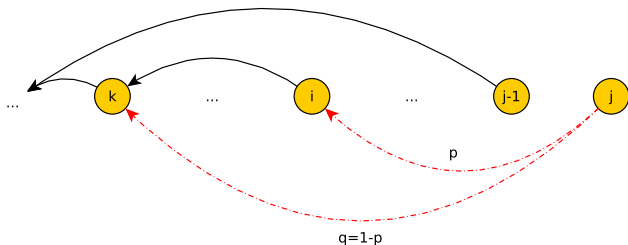
A power law distribution and network dynamics

- A *power law distribution* for a discrete random variable X satisfies $Prob[X = k] \approx \frac{a}{k^c}$ for some constants a and c
 - ▶ We often just focus on on the exponent c and say that $Prob[X = k] \propto k^{-c}$
- Having observed power law distributions emerge from events in social and information networks (e.g., website in-links), the big question is how this happens
 - ▶ We saw that it could *not* evolve from independent decisions that have averaged out; therefore it must arise from correlated decisions
- Kumar et al [2000] proposed a preferential attachment model that can explain the power law distribution
 - ▶ Recall, the observed distribution of in-links is that $Prob[\text{a site has } k \text{ in-links}] \propto k^{-(2+\epsilon)}$ for a small $\epsilon > 0$

A “rich get richer model” for in-links on the Web

Here is the model proposed in Kumar et al article:

- 1 Web pages are created sequentially, and named $1, 2, \dots, N$.
(Of course, N keeps growing but we are looking at the web at some point in time)
- 2 When the j^{th} page is created, we choose a page $i < j$ uniformly at random
 - ▶ With probability p , page j links to page i
 - ▶ Otherwise, page j links to the page (say $k < i$) to which i has a link



Aside

The full model by Kumar et al is more general in that multiple links from page j are created in this stochastic model. Chapter 18 simplifies the model and only creates one link. However, **this does not change the power law exponent**. As will be seen, the key parameter is p .

The linking model continued

There is an equivalent way to frame the indirect linking that highlights the “rich get richer” preferential attachment phenomena.

- [2'] With probability $q = 1 - p$, page j chooses a page $\ell < j$ with probability proportional to ℓ 's current number of in-links and creates a link to ℓ . Otherwise, $\ell < j$ is chosen uniformly at random.

This is, of course, the idea behind *popularity*. For example, the more people that are reading a current novel, the more likely that you might want to read it. And for various social and economic reasons why some large cities continue to grow.

The linking model continued

There is an equivalent way to frame the indirect linking that highlights the “rich get richer” preferential attachment phenomena.

- [2'] With probability $q = 1 - p$, page j chooses a page $\ell < j$ with probability proportional to ℓ 's current number of in-links and creates a link to ℓ . Otherwise, $\ell < j$ is chosen uniformly at random.

This is, of course, the idea behind *popularity*. For example, the more people that are reading a current novel, the more likely that you might want to read it. And for various social and economic reasons why some large cities continue to grow.

Note: As $p \rightarrow 0$ (and $q \rightarrow 1$), pages are more likely to copy the same previous pages and the more likely that the process is creating some popular pages.

The linking model continued

There is an equivalent way to frame the indirect linking that highlights the “rich get richer” preferential attachment phenomena.

- [2'] With probability $q = 1 - p$, page j chooses a page $\ell < j$ with probability proportional to ℓ 's current number of in-links and creates a link to ℓ . Otherwise, $\ell < j$ is chosen uniformly at random.

This is, of course, the idea behind *popularity*. For example, the more people that are reading a current novel, the more likely that you might want to read it. And for various social and economic reasons why some large cities continue to grow.

Note: As $p \rightarrow 0$ (and $q \rightarrow 1$), pages are more likely to copy the same previous pages and the more likely that the process is creating some popular pages.

Hedge: As the text states clearly, the goal of this model is not to capture all the reasons why people create links on the Web (or links in other networks) but rather to explain why it is reasonable to expect power laws to arise from such popularity phenomena.

An informal analysis for the simplified preferential attachment model proposed for Web in-links

A precise analysis of even the simple one link per page preferential attachment model is technical.

There is a heuristic argument that shows how the power law exponent is determined by the probability p (of the j^{th} page linking uniformly at random)

An informal analysis for the simplified preferential attachment model proposed for Web in-links

A precise analysis of even the simple one link per page preferential attachment model is technical.

There is a heuristic argument that shows how the power law exponent is determined by the probability p (of the j^{th} page linking uniformly at random)

While we often discretize continuous processes, it is often advantageous to model a sequence of discrete events as a continuous process

An informal analysis for the simplified preferential attachment model proposed for Web in-links

A precise analysis of even the simple one link per page preferential attachment model is technical.

There is a heuristic argument that shows how the power law exponent is determined by the probability p (of the j^{th} page linking uniformly at random)

While we often discretize continuous processes, it is often advantageous to model a sequence of discrete events as a continuous process

Specifically, we'll consider a continuous deterministic variable $x_\ell(t)$, that approximates the discrete random variable $X_\ell(t)$, the number of in-links to a page ℓ at time $t \geq 0$.

The deterministic continuous model of the random discrete process

- Assuming that page ℓ is added at time ℓ , then $X_\ell(\ell) = x_\ell(\ell) = 0$.
 - ▶ the soonest a link to ℓ can be added is at $\ell + 1$
- In the discrete model, for $t \geq \ell$, the probability that the number of links to a page ℓ increases at time $t + 1$ is:

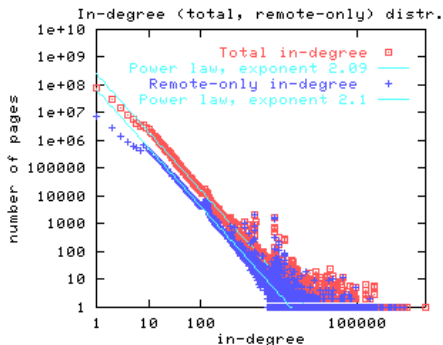
$$\frac{p}{t} + \frac{q \cdot X_\ell(t)}{t}$$

- For $t \geq \ell$, and the corresponding continuous model obeys the differential equation:

$$\frac{dx_\ell}{dt} = \frac{p}{t} + \frac{q \cdot x_\ell(t)}{t}$$

- From some basic calculus (see Ch 18.7) this leads to a power law distribution proportional to k^{-c} with $c = 1 + 1/q$

The deterministic continuous model of the random discrete process



[Fig 18-2 in E&K]

- As $p \rightarrow 0$, the exponent $c = 1 + 1/q$ limits to the observed exponent $c = 2 + \epsilon$ for the observed in-link power law distribution
- As $p \rightarrow 1$, the exponent limits to ∞ making a large number of in-links very unlikely.

Aside: Open Questions

- It's worth noting that although the preferential attachment model suggests that popularity phenomena lead to power laws, it still cannot explain all the examples we saw
- $c = 1 + 1/q \in [2, \infty)$, yet in text and cities the observed exponent is 1

Aside: Open Questions

- It's worth noting that although the preferential attachment model suggests that popularity phenomena lead to power laws, it still cannot explain all the examples we saw
- $c = 1 + 1/q \in [2, \infty)$, yet in text and cities the observed exponent is 1

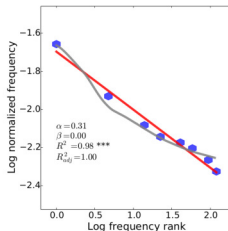


Figure 8: An approximate power law distribution of novel alien names used by subjects in making up a story.

[Figure from Piantadosi (2014)]

- Piantadosi (2014)'s prompt was "An alien space ship crashes in the Nevada desert. Eight creatures emerge, a Wug, a Plit, a Blicket, a Flark, a Warit, a Jupe, a Ralex, and a Timon. In at least 2000 words, describe what happens next"; observed is $c = 0.31$

Fri. Feb 26: Announcements and Corrections

- A1 marks are coming soon™, we will cover the solutions in next week's tutorial

Fri. Feb 26: Announcements and Corrections

- A1 marks are coming soonTM, we will cover the solutions in next week's tutorial
- Critical Review groups & paper selections due in 1 week (March 6)

Fri. Feb 26: Announcements and Corrections

- A1 marks are coming soonTM, we will cover the solutions in next week's tutorial
- Critical Review groups & paper selections due in 1 week (March 6)
 - ▶ Some paper selection advice is on the discussion board

Fri. Feb 26: Announcements and Corrections

- A1 marks are coming soonTM, we will cover the solutions in next week's tutorial
- Critical Review groups & paper selections due in 1 week (March 6)
 - ▶ Some paper selection advice is on the discussion board
- A2 has been released early, see course website for the .pdf and Quercus for the .tex

Fri. Feb 26: Announcements and Corrections

- A1 marks are coming soon™, we will cover the solutions in next week's tutorial
- Critical Review groups & paper selections due in 1 week (March 6)
 - ▶ Some paper selection advice is on the discussion board
- A2 has been released early, see course website for the .pdf and Quercus for the .tex
 - ▶ It's due on March 26th, we haven't yet covered all of the material but you can start on Q1-Q3

Fri. Feb 26: Announcements and Corrections

- A1 marks are coming soon™, we will cover the solutions in next week's tutorial
- Critical Review groups & paper selections due in 1 week (March 6)
 - ▶ Some paper selection advice is on the discussion board
- A2 has been released early, see course website for the .pdf and Quercus for the .tex
 - ▶ It's due on March 26th, we haven't yet covered all of the material but you can start on Q1-Q3
- The Midterm is in 2 weeks, from March 12th to March 14th

Fri. Feb 26: Announcements and Corrections

- A1 marks are coming soon™, we will cover the solutions in next week's tutorial
- Critical Review groups & paper selections due in 1 week (March 6)
 - ▶ Some paper selection advice is on the discussion board
- A2 has been released early, see course website for the .pdf and Quercus for the .tex
 - ▶ It's due on March 26th, we haven't yet covered all of the material but you can start on Q1-Q3
- The Midterm is in 2 weeks, from March 12th to March 14th
 - ▶ Submission via MarkUs (should be visible)
 - ▶ No late submissions! Grace tokens will NOT work here.
 - ▶ PDF will be released around midnight of March 11th on Quercus
 - ▶ An announcement will be made on Quercus, and on the course website
 - ▶ 10 questions (effectively 9)

Be sure to include your name and student number with your test. All tests are to be submitted on Markus. You CANNOT use grace tokens to extend the deadline. Like the assignments, any legible PDF will be accepted. However, if the TAs or myself can't read your solutions then you will get a zero. For this reason, I strongly suggest typesetting your solutions (e.g., in \LaTeX).

You will receive 20% of the points for any (sub)problem for which you write "I do not know how to answer this question." You will receive 10% if you leave a question blank. If instead you submit irrelevant or erroneous answers, you will receive 0 points. You may receive partial credit for the work that is clearly on the right track.

The only resources you are allowed to consult are material from class (e.g., lecture slides, E&K textbook, tutorial recordings, your own notes etc...). You are not allowed to collaborate with other students.

If you have a question, then please email me (see the address on Quercus) and I will reply to you. I will try to answer questions as I receive them between 9AM and 5PM. Outside of these times I'll do my best, but I can't make any guarantees. I will put any clarifications that I think are of value onto the message board. **During the test, do not post or answer to posts on the message board.** Any corrections will be announced via Quercus, in addition to the message board.

You have 3 days to work on the test (March 12, 13, & 14). It *should* only take you 2-4 hours to finish the midterm. This is our first year running a take-home open book midterm so we may have misjudged the difficulty – please don't panic if you complete the midterm faster or slower.

Please read all questions before beginning, and contact me if you have any questions ASAP. Make sure to check the course forum to see if I've posted a clarification there.

Sensitivity to unpredictable initial stages in network dynamics

As we are all are familiar, it is never clear why say some “pop” singers become so popular while other (perhaps of equal talent) never “make it”. Clearly, the initial stages of a dynamic process are critical and that is why advertising, promotions, etc. are so important.

How can we better understand the impact of the randomness in the initial stages of a dynamic process? What if we could replay history many times? We would, of course, expect the resulting distribution to be the same. But would the same books, the same movies, the same pop stars, the same web pages, etc continue to be the most popular?

Sensitivity to unpredictable initial stages in network dynamics

As we are all are familiar, it is never clear why say some “pop” singers become so popular while other (perhaps of equal talent) never “make it”. Clearly, the initial stages of a dynamic process are critical and that is why advertising, promotions, etc. are so important.

How can we better understand the impact of the randomness in the initial stages of a dynamic process? What if we could replay history many times? We would, of course, expect the resulting distribution to be the same. But would the same books, the same movies, the same pop stars, the same web pages, etc continue to be the most popular?

Our intuition (and experience) suggests that there is often considerable “luck” in exactly who or what becomes popular, On the flip side, we also believe that “quality” is also important.

But how do we “rewind history”?

An experiment to “rewind history”



Although we can't rewind history, Salganik et al perform an interesting experiment (in fact, two experiments at different times with different participants) to observe the impact of the initial random stages in a dynamic process. (the article is available on the course website).

The Salganik et al experiment

Here is their experiment:

- They created 9 copies of a music streaming site with 48 “obscure” (as determined by some experts) songs of varying “quality”
- Approximately 7200 young participants were recruited to listen to the music, knowing only the name of the band and the song.
- In each of the copies, participants sequentially listened to some music selections, rated the music and then were given the opportunity to download copies of songs they liked.
- In 8 copies of the site (each with 10% of the participants), they were also given the number of times each song had been previously downloaded.
- In the 9th version, this previous history of downloads was not provided to the remaining 20% of the participants. The average of the ratings (from 1 = “I hated it” to 5 = “I loved it”) in this “no influence” version determined the “true” song “quality”.
- The experiment was then repeated, with the 8 site copies displaying songs sorted by downloads instead of randomly

The findings in the Salganik experiment

The experiment was designed to measure the extent that social influence leads to different outcomes in the “success” (i.e. the number of downloads) of a particular song.

Simply stated, the results show that:

- Increasing the strength of social influence (by sorting songs by downloads) increased both the inequality (i.e. degrees of popularity) and unpredictability (i.e., relation to quality) of success.
- However, quality was also a factor: the best rated songs rarely did poorly and the the worst songs rarely did well.

As I said, this is an interesting study and one where the authors carefully try to eliminate sources of bias. The article is worth reading.

As the text points out in section 18.6, how recommendation systems are designed can impact how people make choices, leading to increased “rich get richer” phenomena, or alternatively exposing people to less popular items.

Visualizing the long tail of a power law distribution

Once we accept a power law nature of popularity, it is instructive to consider the consequences for a given industry. Namely, the nature of the sales curve that would be dictated by a power law distribution.

The shape of the long tail in a power distribution raises the question as to how many sales can be obtained from less popular (e.g. niche items).

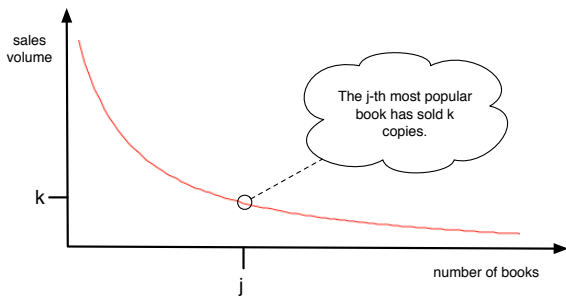


Figure: [Fig 18-4 in E&K] text; how many copies of the j^{th} most popular items have been sold.

Search and ranking on the Web

Our next topic is to understand how the popularity of a web page is determined and how that impacts its rank in the responses to a query.

But first, how do search engines find and rank responses to a query?

Search and ranking on the Web

Our next topic is to understand how the popularity of a web page is determined and how that impacts its rank in the responses to a query.

But first, how do search engines find and rank responses to a query?

The specific algorithms used by search engines such as Bing and Google is a trade secret. To some extent this has to be kept secret as there is always a “war” between a search engine and companies that create web sites to enhance the ranking of a site.

However, we do have a basic idea as to how these search engines rank sites given a query. In fact, at the most elementary level, the main idea is an old one.

Aside: In the 1960s and 70s, there was a basic argument as to whether online search and ranking was a more or less normal algorithmic search and optimization problem or one that required “intelligence” (i.e. the ability to understand natural language). **Who won this argument?**

Search and ranking of Web documents; the role of link popularity

The most basic approach is to treat a document as a bag of words and then use “normalized” word counts (and pairs, triplets of words) to identify and rank documents relating to the query. This became enhanced by more sophisticated contextual aspects of word occurrences, etc and today machine learning algorithms are also used in classifying a search query.

But early in the development of popular search engines, a popularity aspect was added where the ranking of a document also depended on the link structure and the popularity of a Web page in the Web network (or at least in that part that seems relevant to the query).

Two algorithms were independently proposed for determining the popularity of a Web page, namely Hubs and Authorities developed at IBM, and Page Rank, developed and integrated into Google’s search engine.

Link analysis and page popularity

Neither Hubs and Authorities nor Page Rank use link in-degree as the popularity measure but link analysis is (or at least was) used to determine page popularity. Currently, it seems clear that popularity also depends on recent behaviour of users to related queries.

We will not try to infer more precisely how say Google (or any search engine) precisely determines the ranking of a document in response to a query. In particular, we do not know how much page ranking depends on content vs link analysis. But we do know that this ranking is essential in determining how often a page will be downloaded. The quality of the ranking algorithm leads to user activity and thus the resulting advertising.

We will begin with the Hubs and Authorities ranking algorithm and then the Page Rank algorithm.

Hubs and Authorities

- A simple way to utilize links to rank web pages would be to think of each link from A to B as an endorsement or vote by A for B .
- **Question:** Assuming it's tractable, then what's wrong with just counting the number of in-links?

Hubs and Authorities

- A simple way to utilize links to rank web pages would be to think of each link from A to B as an endorsement or vote by A for B .
- **Question:** Assuming it's tractable, then what's wrong with just counting the number of in-links?
- If we use the number (or weight) of endorsements to determine rank, then one would have to adjust such scores coming from say the same domain name.

Hubs and Authorities

- A simple way to utilize links to rank web pages would be to think of each link from A to B as an endorsement or vote by A for B .
- **Question:** Assuming it's tractable, then what's wrong with just counting the number of in-links?
- If we use the number (or weight) of endorsements to determine rank, then one would have to adjust such scores coming from say the same domain name.
- Even after adjusting for such "vote fixing", if Wayne Gretzky has a web site with a link suggesting where to buy hockey equipment, that should be more meaningful than my recommendation about hockey equipment.

Hubs and Authorities

- A simple way to utilize links to rank web pages would be to think of each link from A to B as an endorsement or vote by A for B .
- **Question:** Assuming it's tractable, then what's wrong with just counting the number of in-links?
- If we use the number (or weight) of endorsements to determine rank, then one would have to adjust such scores coming from say the same domain name.
- Even after adjusting for such “vote fixing”, if Wayne Gretzky has a web site with a link suggesting where to buy hockey equipment, that should be more meaningful than my recommendation about hockey equipment.
Spoiler alert: I don't play or watch hockey.

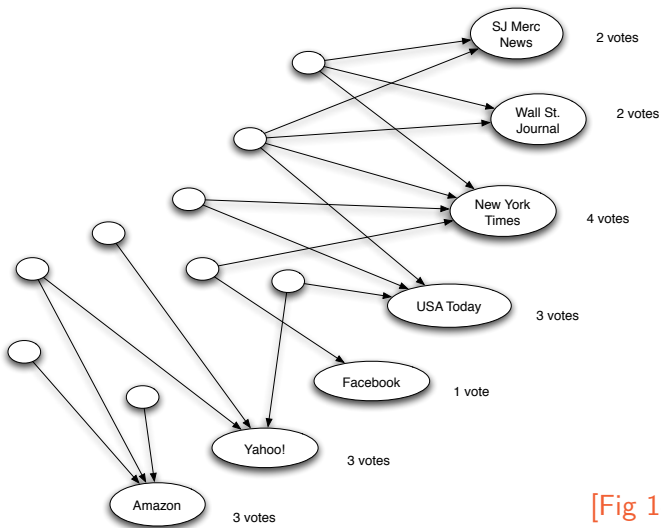
Reinforcement of Hubs and Authorities.

- This then becomes the motivation (and seemingly circular reasoning) behind hubs and authorities.
- The **best “authorities”** on a subject (places to buy equipment) are being **endorsed by the best hubs** (people who know where to buy equipment).
- Similarly, the **best hubs** are those sites that **recommend the best authorities**. Conceptually the link structure induces a bipartite graph, however the same web page can be both a hub and an authority.
- **Comment:** The word **“authority”** is not generally an accurate way to describe high ranking documents. These might better be referred to (barring other information) as the **most relied upon sites**. This is also different from “the most popular” sites which might better be measured in terms of the number of clicks being received. **Hubs** then are the **most reliable endorsers**.

Hubs & Authorities procedure

- The hubs and authorities procedure is as follows:
 - ▶ Initialize each node' hub value to some positive number (perhaps depending on usage or content)
 - ▶ For sufficiently large k , perform the following k times
 - ★ Apply authority update rule to each page, p (i.e., set the authority value to the sum of the hub values of the nodes endorsing the page p)
 - ★ Apply hub update rule to each page, p (i.e., set the hub value to the sum of the authority values of the nodes endorsed by the page p)
 - ★ Normalize so that sum of A and H weights = 1.

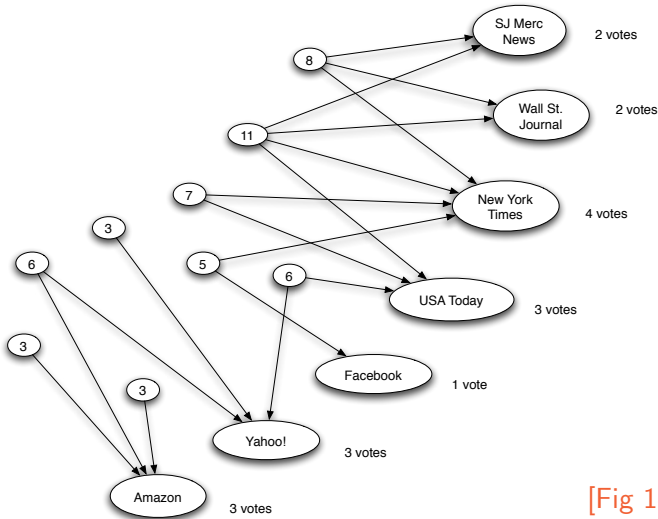
- The result of applying the **authority update rule** with all hub values initially 1: for each page p , $\text{auth}(p)$ is the sum of hub values (initially just the number) of hubs pointing to p .



[Fig 14.1, E&K]

Figure: Counting in-links to pages for the query “newspapers.”

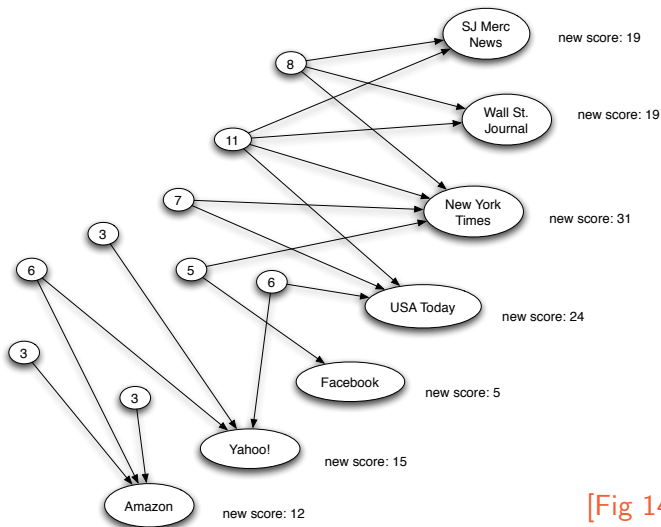
- Then to recalibrate hub values, we use the **hub update rule**: for each page p , $\text{hub}(p)$ is the sum of values of all authorities that p points to.



[Fig 14.2, E&K]

Figure: Finding good lists for the query “newspapers”: each page’s value as a list is written as a number inside it.

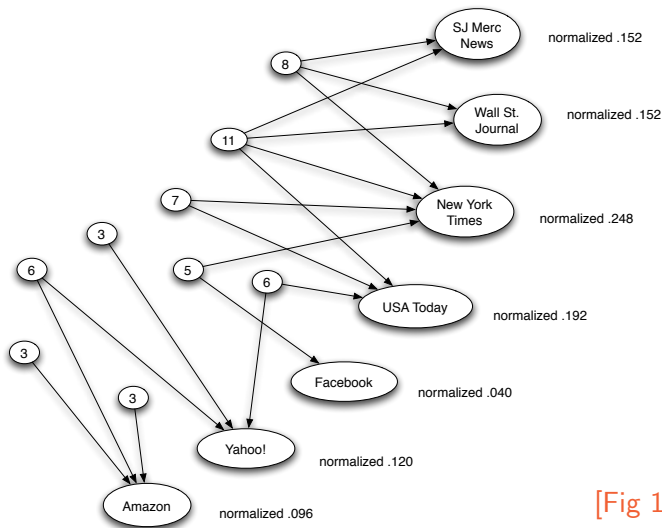
- Applying the authority update rule again we get figure 14.3.



[Fig 14.3, E&K]

Figure: Re-weighting votes for the query “newspapers”: each of the labeled pages new score is equal to the sum of the values of all lists that point to it.

- Since we only care about the relative values of these numbers, both authority and hub scores can be normalized to sum to 1 (to allow convergence and avoid dealing with large numbers).

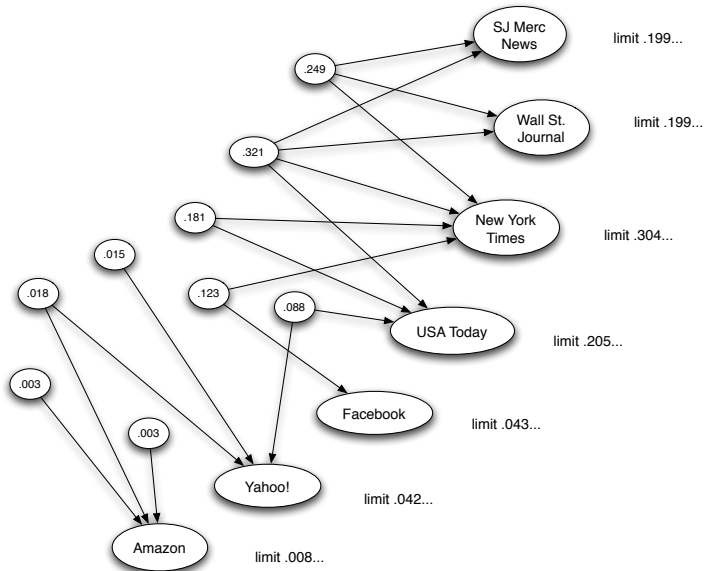


[Fig 14.4, E&K]

Figure: Re-weighting votes after normalizing for the query “newspapers”. 39 / 58

Keep repeating a good idea

- Now having recalibrated and normalized both the authority and hub scores, we can continue this process to continue to refine these scores.
- That is, the **hubs and authorities procedure** is as follows:
 - ▶ Initialize all hub values to some positive vector (perhaps depending on usage or content)
 - ▶ For sufficiently large k , perform the following k times
 - ★ Apply **authority update rule** to each page
 - ★ Apply **hub update rule** to each page
 - ★ Normalize so that sum of A and H weights = 1.
- Using linear algebra, it can be shown (in Section 14.6) that these A and H normalized values will **converge to a limit** as $k \rightarrow \infty$ (which can be approximated by some sufficiently large k)! The limiting value is an equilibrium.
- Hubs and Authorities can be extended to work for weighted edges (e.g. weighting links in anchor text, or near a section heading, etc.)



[Fig 14.5, E&K]

Figure: Limiting hub and authority values for the query "newspapers".

Page Rank

- The motivation behind page rank is a somewhat different view of how authority is conferred.
 - ▶ Endorsement of authority is conveyed by other authorities
 - ▶ That is, no hub concept
 - ▶ This is how peer review works in the academic and scholarly world.

- Authorities themselves convey authority on those they link to. This naturally leads to a formulation in terms of two equivalent views of page rank:
 - 1 Authorities directly conveying authority (without hubs)
 - 2 Authority values resulting from long term behaviour of a random walk on a graph.

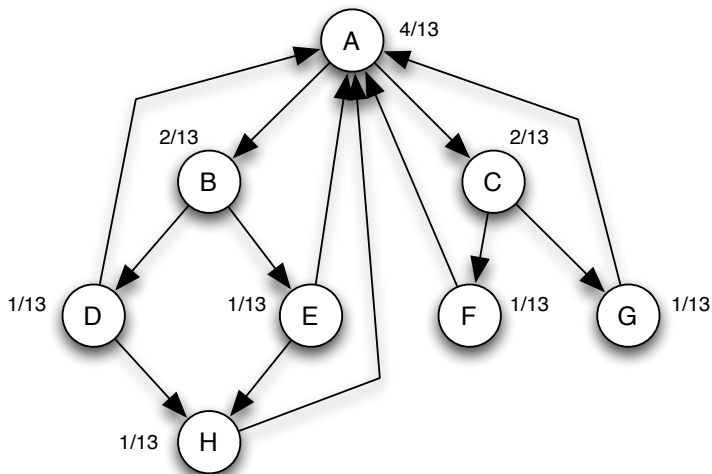
How does Page rank spread authority?

- Suppose at any point of time we have relevant authority scores.
 - ▶ A page **spreads its authority equally amongst all of its out links**.
 - ▶ If a page has no outlinks then all authority stays there.
- This **redistributes the authority scores**. (We are not creating or losing any authority, we are just redistributing it.)
- We can initially start with every relevant page having authority $1/n$ where there are n pages. Then we **repeat this process k times** for some sufficiently large k .
- With the exception of some “degenerate cases” (e.g. the process is periodic) it can be proven (again using linear algebra) that this process has a limiting behavior as $k \rightarrow \infty$.
- The resulting limit values will form an equilibrium.
- If the network is strongly connected then there is a unique equilibrium,

Remark

In many cases this won't reflect the desired authority. Namely, if the network has any sinks which it will surely have, then all of the authority will pass to such sinks.

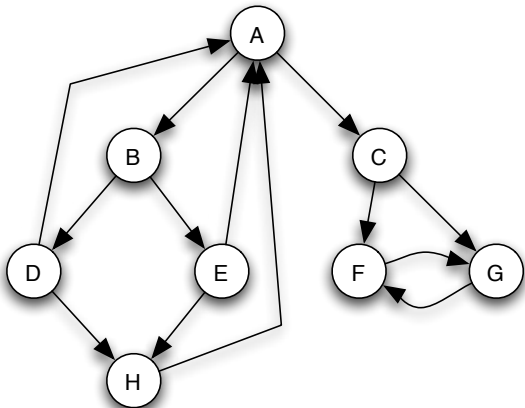
Page rank equilibrium for a network



[Fig 14.7, E&K]

Figure: Equilibrium PageRank values for the network of eight Web page.

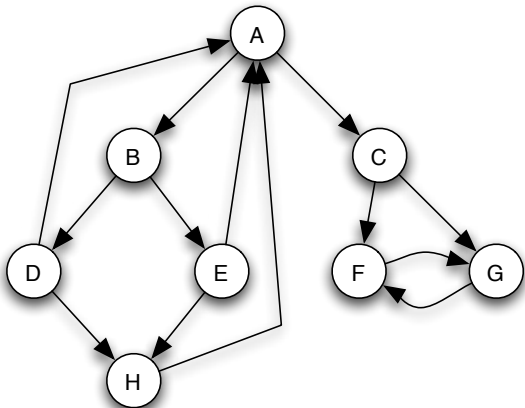
Where has all the authority gone when we redirect (F, A) and (G, A) edges?



[Fig 14.8, E&K]

The same collection of eight pages, but F and G have changed their links to point to each other instead of to A .

Where has all the authority gone when we redirect (F, A) and (G, A) edges?



[Fig 14.8, E&K]

The same collection of eight pages, but F and G have changed their links to point to each other instead of to A . Without “**scaling**”, all the PageRank would go to F and G .

How does PageRank spread authority?

- If the network is strongly connected then there is a unique equilibrium, however it may be undesirable

Definition (Sink)

Strict (typically small) subset of the nodes with no outgoing edges that are reachable from all nodes in the network

Remark

In many cases PageRank won't reflect the desired authority. Namely, if the network has any sinks which it will surely have, then all of the authority will pass to such sinks.

Scaled page rank

- The way around this sink hole of authority is to have a **scaled version of page rank** where
 - ▶ only a fraction s of the authority of a page is distributed to its out links
 - ▶ the remaining $(1 - s)$ fraction is distributed equally amongst all relevant pages.
- For any value of $s < 1$ (which effectively makes the graph strongly connected), we get **convergence to a unique set of scores** for each page and that is its page rank (for that particular value of s). It is reported that Google used $0.8 \leq s \leq 0.9$.
- (See the footnote on page 410 of E&K as to why in the previous example, nodes F and G will still get most of the authority but that for realistically large networks, the process works well.
Hint: “bow-tie” structure)

Some additional remarks

- The limiting scores for both the authority and hubs approach and the page rank approach are **equilibrium points for an appropriate algebraic process**.
- That is, if we actually were in the limiting state, we would be in the equilibrium state. In practice, we can **stop the process when the change in each iteration is sufficiently small**.
- We can **weight the network edges** (say according to some concept of link importance) and apply the same authority and hubs or page rank approach **distributing authority in proportion to these weights**.

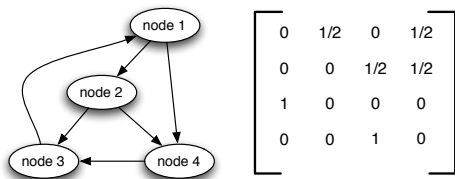
Advanced material (section 14.6): Handwaving argument why these processes converge

We have already suggested that both the page rank and hubs and authorities processes can be understood in terms of an algebraic process, namely, a linear transformation.

- Suppose we are considering a web network of n pages. We can represent the hub, authority or page rank values at any time k of the process by an n -dimensional (column) vector, denoted (respectively) by $\mathbf{h}^{(k)}$, $\mathbf{a}^{(k)}$, $\mathbf{r}^{(k)}$.
- Here we are using boldface $\mathbf{v} = \langle v_1, \dots, v_n \rangle$ to represent a vector whose components are the v_j so that (for example), $r_j^{(k)}$ represents the page rank of the j^{th} web page after k steps of the page rank process.
- Let \mathbf{v} be any of the hub, authority or page rank vectors. In each case it is not difficult to see that the process can be viewed as a linear transformation $\mathbf{v}^{(k+1)} = M\mathbf{v}^{(k)}$ for some appropriate $n \times n$ matrix M whose entries are non negative real numbers.

Advanced material continued: page rank convergence

- Section 14.6 tells us how to define the appropriate matrices and gives the conditions that will guarantee the convergence of the process; that is, when there exists $\mathbf{v}^{(*)} = \lim_{k \rightarrow \infty} \mathbf{v}^{(k)}$ and when this limit vector $\mathbf{v}^{(*)}$ is unique and independent of the starting vector $\mathbf{v}^{(0)}$.
- Figure 14.3 of the text illustrates a simple directed graph and the matrix N that defines the unscaled page rank update process. That is, $\langle r_1^{k+1}, \dots, r_n^{k+1} \rangle = N^{tr} \langle r_1^k, \dots, r_n^k \rangle$ where N^{tr} is the transpose of matrix N .

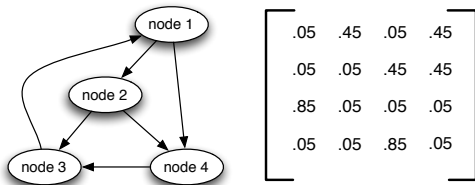


[Fig 14.13, E&K]

Figure: A toy web graph and the associated matrix N describing the unscaled update process.

Page rank analysis for the scaled update

Similarly Figure 14.4 illustrates the same graph and the matrix \tilde{N} that defines the scaled page rank update process with scaling factor $s = .8$.



[Fig 14.14, E&K]

Figure: The same toy web graph and the associated matrix \tilde{N} describing the scaled update process with $s = 0.8$.

- $\tilde{N} = sN + \frac{(1-s)}{n} \mathbf{1}^T \mathbf{1}$
- It follows that $\mathbf{r}^k = (\tilde{N}^{tr})^k \mathbf{r}^0$
- If the process is converging then it would be converging to some \mathbf{r}^* satisfying $\mathbf{r}^* = N^{tr} \mathbf{r}^*$

Now comes the necessary linear algebra

So far we have mainly used matrices as a convenient way to represent the process. But to understand convergence we need to mention some more essential aspects of linear algebra.

- Let $M_{n \times n}$ be a full rank matrix. Recall that the matrix-vector multiplication $M\mathbf{v}$ can rotate and expand/shrink the vector \mathbf{v} .
- Since it is hard to “visualize” an n -dimensional vector space, we can simply think about the meaning of such a linear transformation in 2-space or 3-space.
- A vector \mathbf{v} is an **eigenvector** of M with associated **eigenvalue** λ if $M\mathbf{v} = \lambda\mathbf{v}$. It follows that \mathbf{v} is also an eigenvector of M^k with eigenvalue λ^k .
- When $\lambda = 1$, the eigenvector then becomes an equilibrium of the process!

More linear algebra

- For some matrices there is a set of n eigenvectors with (not necessarily distinct) associated eigenvalues $\lambda_1, \dots, \lambda_n$; these eigenvectors span the n -dimensional Euclidean space so that any vector can be expressed as a linear combination of the eigenvectors.
- An important result from linear algebra (Perron's Theorem) states that any matrix which has all positive entries has a unique eigenvector \mathbf{y} corresponding to the largest positive eigenvalue λ_1 and furthermore $\lambda_1 > |\lambda_i|$ for $i > 1$.
- Since $\lambda_1 > |\lambda_i|$ for $i > 1$, and since every vector is a linear combination of the eigenvectors, it follows that as $k \rightarrow \infty$, the transformation M^k is being dominated by the largest eigenvalue acting on its associated eigenvector.
- For the scaled matrix \tilde{N}^{tr} , all entries are positive and the largest eigenvalue is 1. It follows that as $k \rightarrow \infty$, $(\tilde{N}^{tr})^k \mathbf{v}$ will converge to the eigenvector \mathbf{y} associated with the largest eigenvalue 1.

Aside: Random Walk

- Remember that we said that PageRank was equivalent to a random walk on the graph?

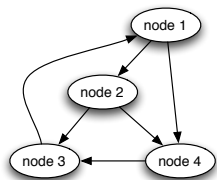
Aside: Random Walk

- Remember that we said that PageRank was equivalent to a random walk on the graph?
 - ▶ We can view this as a Markov chain (a type of probabilistic finite state machine that's represented as a graph, where each timestep we follow an edge based on the corresponding probabilities)

Aside: Random Walk

- Remember that we said that PageRank was equivalent to a random walk on the graph?
 - ▶ We can view this as a Markov chain (a type of probabilistic finite state machine that's represented as a graph, where each timestep we follow an edge based on the corresponding probabilities)
- Similarly, scaled page rank can be viewed as the same Markov chain but with added low probability edges between every pair of nodes

Aside: Random Walk

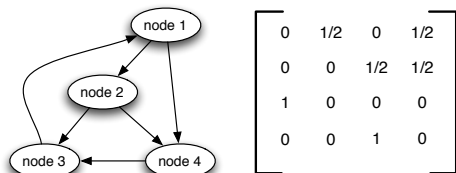


$$\begin{bmatrix} 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 1/2 & 1/2 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

[Fig 14.13, E&K]

- Specifically, the *transition matrices* of these Markov chains are N and \tilde{N} respectively (N_{ij} is the probability of transitioning from state i to state j)

Aside: Random Walk



[Fig 14.13, E&K]

- Specifically, the *transition matrices* of these Markov chains are N and \tilde{N} respectively (N_{ij} is the probability of transitioning from state i to state j)
- Under the assumption that the chain is *irreducible* (we can reach state j from all states i and vice versa), and *aperiodic* (there is no state i such that if you leave i , you can only return on timesteps that are multiples of some $p > 1$) then there is a unique *stationary distribution* π that the chain converges to
- This is the distribution that we found

Similar analysis for hubs and authorities

- If M is the adjacency matrix of the web graph, then the process can be described by $\mathbf{h} = M\mathbf{a}$ and $\mathbf{a} = M^{tr}\mathbf{h}$.
- Then
 - 1 $\mathbf{a}^{(1)} = M^{tr}\mathbf{h}^{(0)}$
 - 2 $\mathbf{h}^{(1)} = M\mathbf{a}^{(1)} = MM^{tr}\mathbf{h}^{(0)}$
- It follows that
 - 1 $\mathbf{a}^{(k)} = (M^{tr}M)^{k-1}M^{tr}\mathbf{h}^{(0)}$
 - 2 $\mathbf{h}^{(k)} = (MM^{tr})^k\mathbf{h}^{(0)}$

Hubs and authorities analysis continued

- The matrices (MM^{tr}) and $(M^{tr}M)$ are **symmetric** and have non-negative entries.
- Any $n \times n$ symmetric matrix S with non negative entries has an orthonormal set of n eigenvectors all of whose associated eigenvalues are real. By normalizing the scores, we can assume that the largest eigenvalue $\lambda_1 = 1$.
- If the largest eigenvalue is unique (which is what would happen in a real web graph), then the same analysis for page rank applies (assuming that the starting hub scores are all positive).

Returning to the issue of influence

In some sense or another we are often talking about social influence in this course. Even in web ranking (Ch 14), we can view hubs as influencing which Web pages will be ranked highly.

In chapter 18, we observed two sequential processes where previous individual decisions had a significant impact:

- 1) The evolution of links on the Web, and
- 2) The evolution of opinions in evaluating music.

The music evaluation experiment is closer to reality in the sense that it explicitly integrates a measure of quality (a simplification of selection?) into the decision making process.