

Five Lectures on Networks

Aaron Clauset

 @aaronclauset

Assistant Professor of Computer Science

University of Colorado Boulder

External Faculty, Santa Fe Institute

lecture 5: learning from network data and metadata



University of Colorado **Boulder**

Network Analysis and Modeling

Instructor: Aaron Clauset

This graduate-level course will examine modern techniques for analyzing and modeling the structure and dynamics of complex networks. The focus will be on statistical algorithms and methods, and both lectures and assignments will emphasize model interpretability and understanding the processes that generate real data. Applications will be drawn from computational biology and computational social science. No biological or social science training is required. (Note: this is not a scientific computing course, but there will be plenty of computing for science.)

Full lectures notes online (~150 pages in PDF)

<http://santafe.edu/~aarond/courses/5352/>

Software

[R](#)
[Python](#)
[Matlab](#)
[NetworkX \[python\]](#)
[graph-tool \[python, c++\]](#)
[GraphLab \[python, c++\]](#)

Standalone editors

[UCI-Net](#)
[NodeXL](#)
[Gephi](#)
[Pajek](#)
[Network Workbench](#)
[Cytoscape](#)
[yEd graph editor](#)
[Graphviz](#)

Data sets

[Mark Newman's network data sets](#)
[Stanford Network Analysis Project](#)
[Carnegie Mellon CASOS data sets](#)
[NCEAS food web data sets](#)
[UCI NET data sets](#)
[Pajek data sets](#)
[Linkgroup's list of network data sets](#)
[Barabasi lab data sets](#)
[Jake Hofman's online network data sets](#)
[Alex Arenas's data sets](#)

1. defining a network
2. describing a network
3. null models for networks
- 4. statistical inference**

what is structure?

what is structure?

- makes data different from noise
 - makes a network different from a random graph

what is structure?

- makes data different from noise
 - makes a network different from a random graph
- helps us compress the data
 - describe the network succinctly
 - capture most relevant patterns

what is structure?

- makes data different from noise
 - makes a network different from a random graph
- helps us compress the data
 - describe the network succinctly
 - capture most relevant patterns
- helps us generalize,
from data we've seen to data we haven't seen:
 - i. from one part of network to another
 - ii. from one network to others of same type
 - iii. from small scale to large scale (coarse-grained structure)
 - iv. from past to future (dynamics)

statistical inference

- imagine graph G is drawn from an ensemble or **generative model**: a probability distribution $\Pr(G | \theta)$ with parameters θ
- θ can be continuous or discrete; represents structure of graph

statistical inference

- imagine graph G is drawn from an ensemble or **generative model**: a probability distribution $\Pr(G | \theta)$ with parameters θ
- θ can be continuous or discrete; represents structure of graph
- inference (MLE): given G , find θ that maximizes $\Pr(G | \theta)$
- inference (Bayes): compute or sample from posterior distribution $\Pr(\theta | G)$

statistical inference

- imagine graph G is drawn from an ensemble or **generative model**: a probability distribution $\Pr(G | \theta)$ with parameters θ
 - θ can be continuous or discrete; represents structure of graph
 - inference (MLE): given G , find θ that maximizes $\Pr(G | \theta)$
 - inference (Bayes): compute or sample from posterior distribution $\Pr(\theta | G)$
-

- if θ is partly known, constrain inference and determine the rest
- if G is partly known, infer θ and use $\Pr(G | \theta)$ to generate the rest
- if model is good fit (application dependent), we can generate synthetic graphs structurally similar to G
- if part of G has low probability under model, flag as possible anomaly

statistical inference

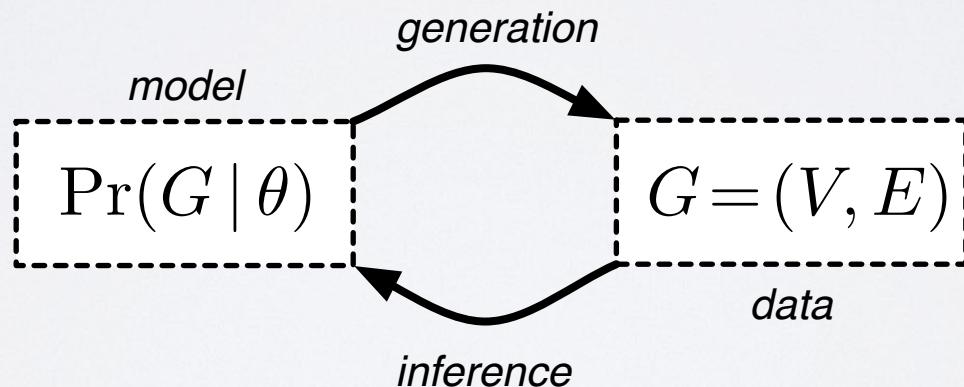
- imagine graph G model
 - θ can be continuous or discrete; represents structure of graph
 - inference (MLE): given G and θ that maximize $\Pr(G|\theta)$
 - inference (Bayes): compute sample from posterior distribution $\Pr(\theta|G)$
 - if θ is partly known, constrain inference and determine the rest
 - if G is partly known, infer θ and use $\Pr(G|\theta)$ to generate the rest
 - if model is good fit (application dependent), we can generate synthetic graphs structurally similar to G
 - if part of G has low probability under model, flag as possible anomaly
- statistical inference = principled approach to learning from data**
- combines tools from statistics, machine learning, information theory, and statistical physics**
- quantifies uncertainty**
- separates the model from the learning**

statistical inference: key ideas

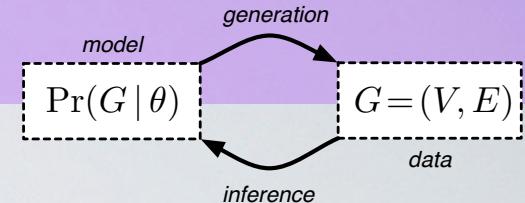
- interpretability
 - model parameters have meaning for scientific questions
- auxiliary information
 - node & edge attributes, temporal dynamics (beyond static binary graphs)
- scalability
 - fast algorithms for fitting models to big data (methods from physics, machine learning)
- model selection
 - which model is better? is this model bad? how many communities?
- partial or noisy data
 - extrapolation, interpolation, hidden data, missing data
- anomaly detection
 - low probability events under generative model

generative models for complex networks

- define a parametric probability distribution over networks $\Pr(G | \theta)$
- **generation** : given θ , draw G from this distribution
- **inference** : given G , choose θ that makes G likely



generative models for complex networks



general form

$$\Pr(G | \theta) = \prod_{i < j} \Pr(A_{ij} | \theta)$$

"attachment" function

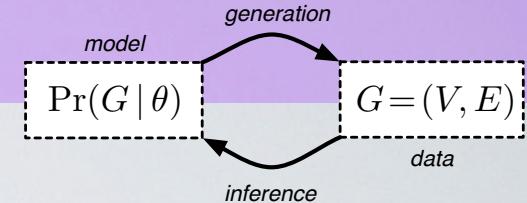
assumptions about "structure" go into $\Pr(A_{ij} | \theta)$

$$\text{consistency } \lim_{n \rightarrow \infty} \Pr(\hat{\theta} \neq \theta) = 0$$

requires that edges be conditionally independent [Shalizi, Rinaldo 2011]

two general classes of these models

generative models for complex networks



stochastic block models

k types of vertices, $\Pr(A_{ij} \mid M, z)$ depends only on types z_i, z_j
originally invented by sociologists [Holland, Laskey, Leinhardt 1983]

many, many flavors, including

binomial SBM [Holland, Laskey, Leinhardt 1983, Wang & Wong 1987]

mixed-membership SBM [Airoldi, Blei, Feinberg, Xing 2008]

hierarchical SBM [Clauset, Moore, Newman 2006, 2008, Peixoto 2014]

fractal SBM [Leskovec et al. 2005]

infinite relational model [Kemp et al. 2006]

simple assortative SBM [Hofman & Wiggins 2008]

degree-corrected SBM [Karrer & Newman 2011]

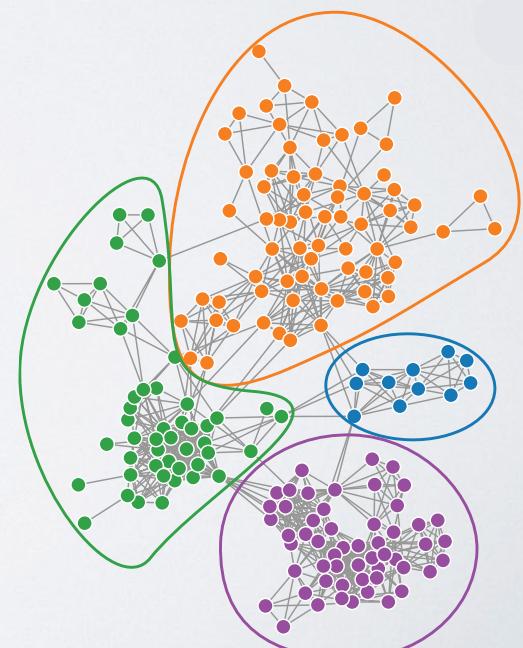
SBM + topic models [Ball, Karrer & Newman 2011]

SBM + vertex covariates [Mariadassou, Robin & Vacher 2010]

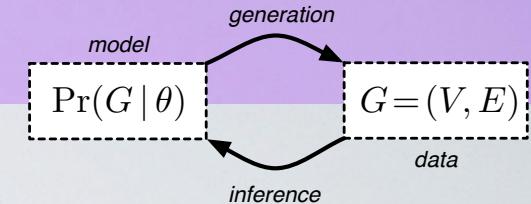
SBM + edge weights [Aicher, Jacobs & Clauset 2013, 2014]

bipartite SBM [Larremore, Clauset & Jacobs 2014]

and many others



generative models for complex networks



latent space models

nodes live in a latent space, $\Pr(A_{ij} \mid f(x_i, x_j))$ depends only on vertex-vertex proximity

many, many flavors, including

logistic function on vertex features [Hoff, Raftery, Handcock 2002]

social status / ranking [Ball, Newman 2013]

nonparametric metadata relations [Kim, Hughes, Sudderth 2012]

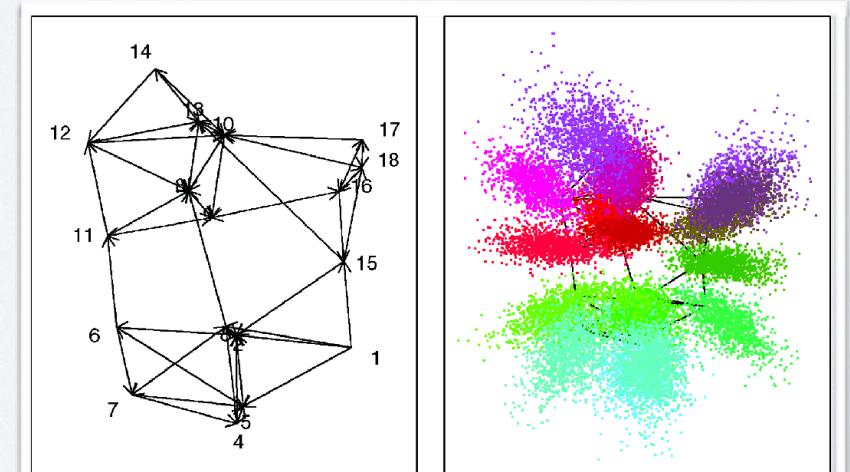
multiplicative attribute graphs [Kim, Leskovec 2010]

nonparametric latent feature model [Miller, Griffiths, Jordan 2009]

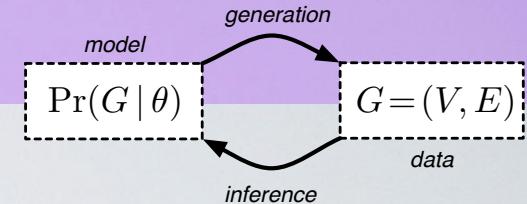
infinite multiple memberships [Morup, Schmidt, Hansen 2011]

ecological niche model [Williams, Anandanadesan, Purves 2010]

hyperbolic latent spaces [Boguna, Papadopoulos, Krioukov 2010]

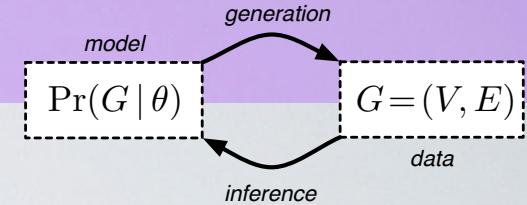


opportunities and challenges



- richly annotated data
 - edge weights, node attributes, time, etc.
 - = new classes of generative models
- generalize from $n = 1$ to ensemble
 - useful for modeling checking, simulating other processes, etc.
- many familiar techniques
 - frequentist and Bayesian frameworks
 - makes probabilistic statements about observations, models
 - predicting missing links \approx leave- k -out cross validation
 - approximate inference techniques (EM, VB, BP, etc.)
 - sampling techniques (MCMC, Gibbs, etc.)
- learn from partial or noisy data
 - extrapolation, interpolation, hidden data, missing data

opportunities and challenges



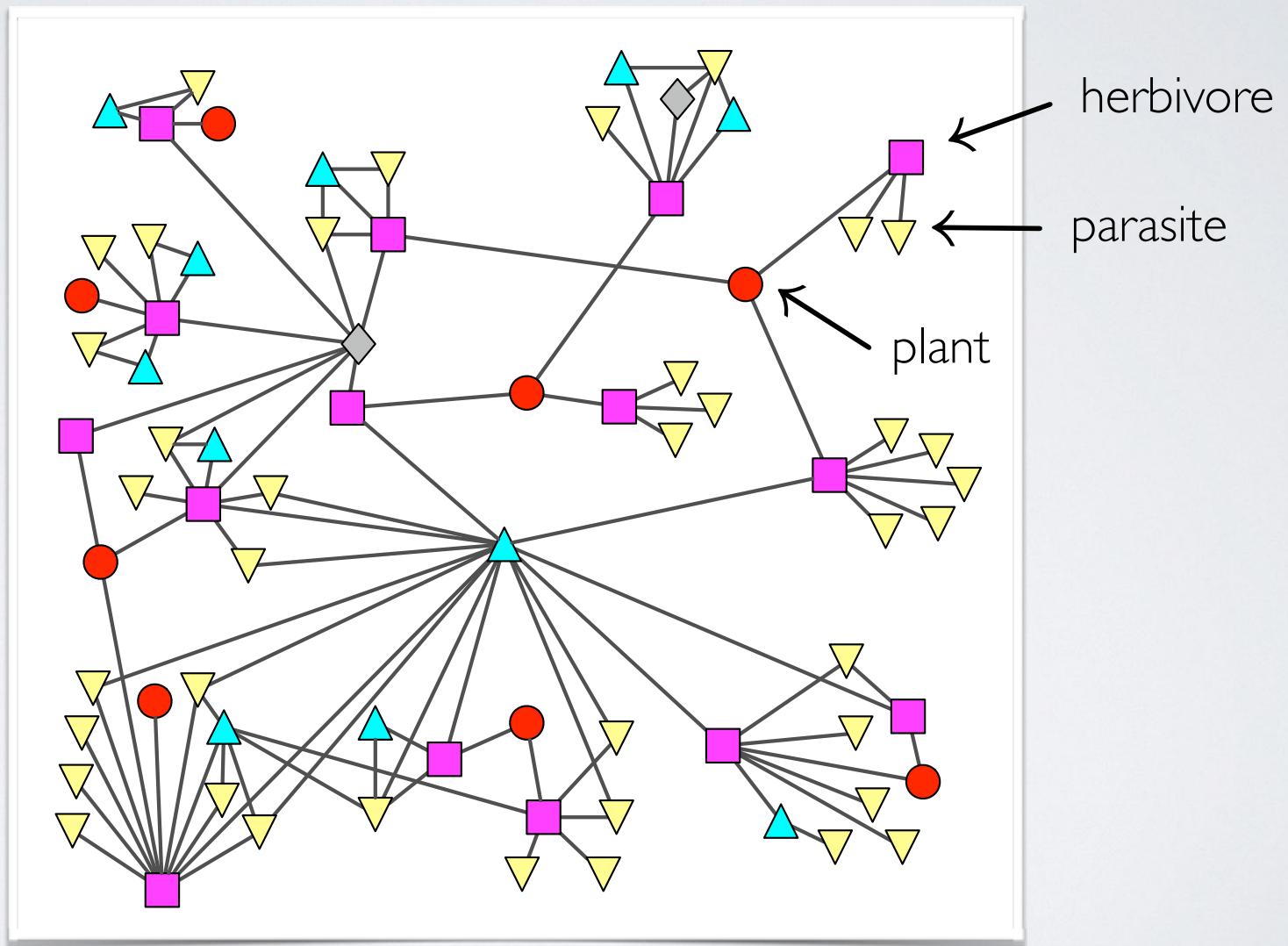
- only two classes of models
 - stochastic block models
 - latent space models
- bootstrap / resampling for network data
 - critical missing piece
 - depends on what is independent in the data
- model comparison
 - naive AIC, BIC, marginalization, LRT can be wrong for networks
 - what is goal of modeling: realistic representation or accurate prediction?
- model assessment / checking?
 - how do we know a model has done well? what do we check?
- what is v -fold cross-validation for networks?
 - Omit n^2/v edges? Omit n/v nodes? What?

generative models for complex networks

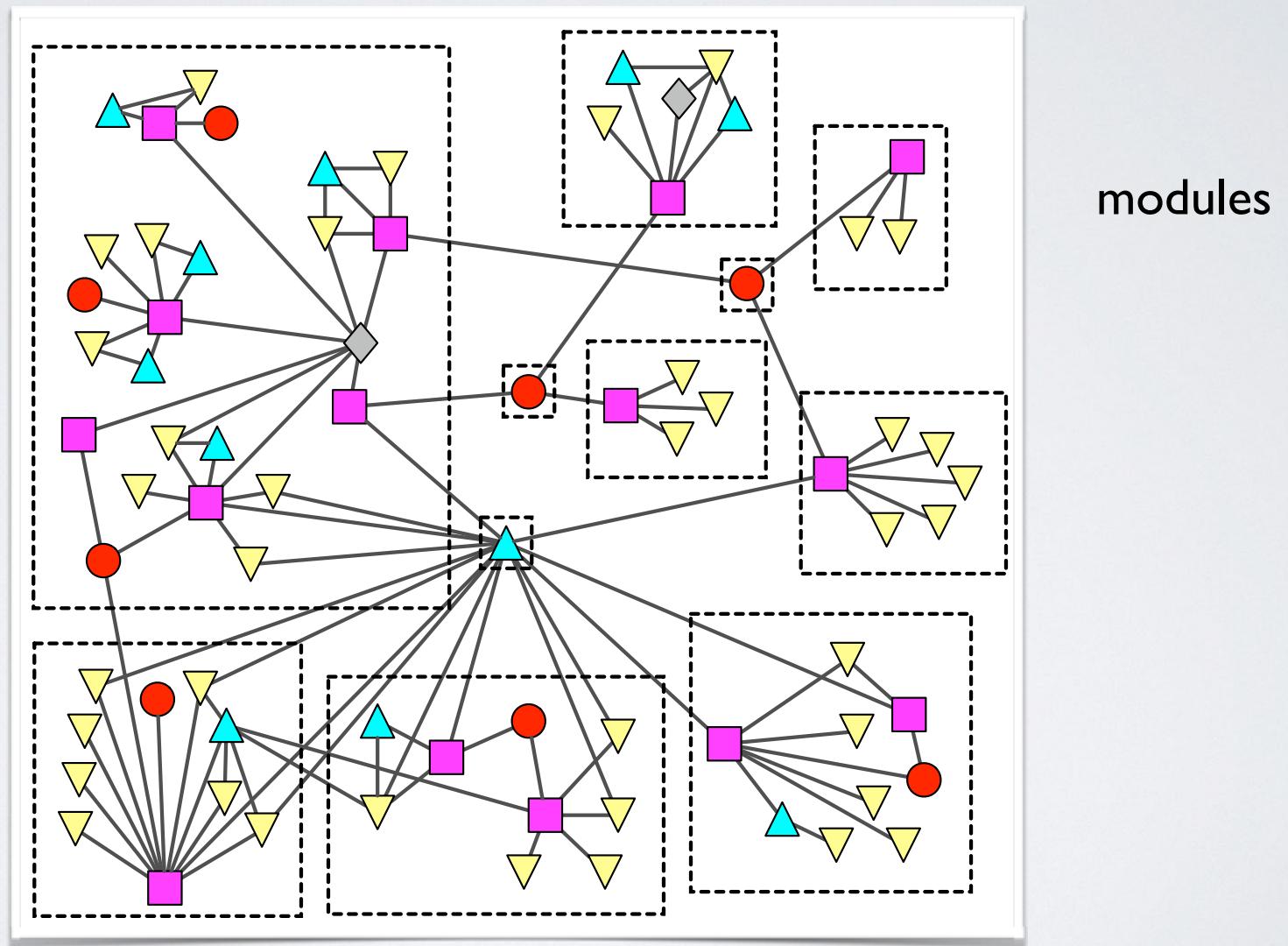
3 examples

- hierarchical random graphs (HRG)
- bipartite community structure (biSBM)
- weighted community structure (WSBM)

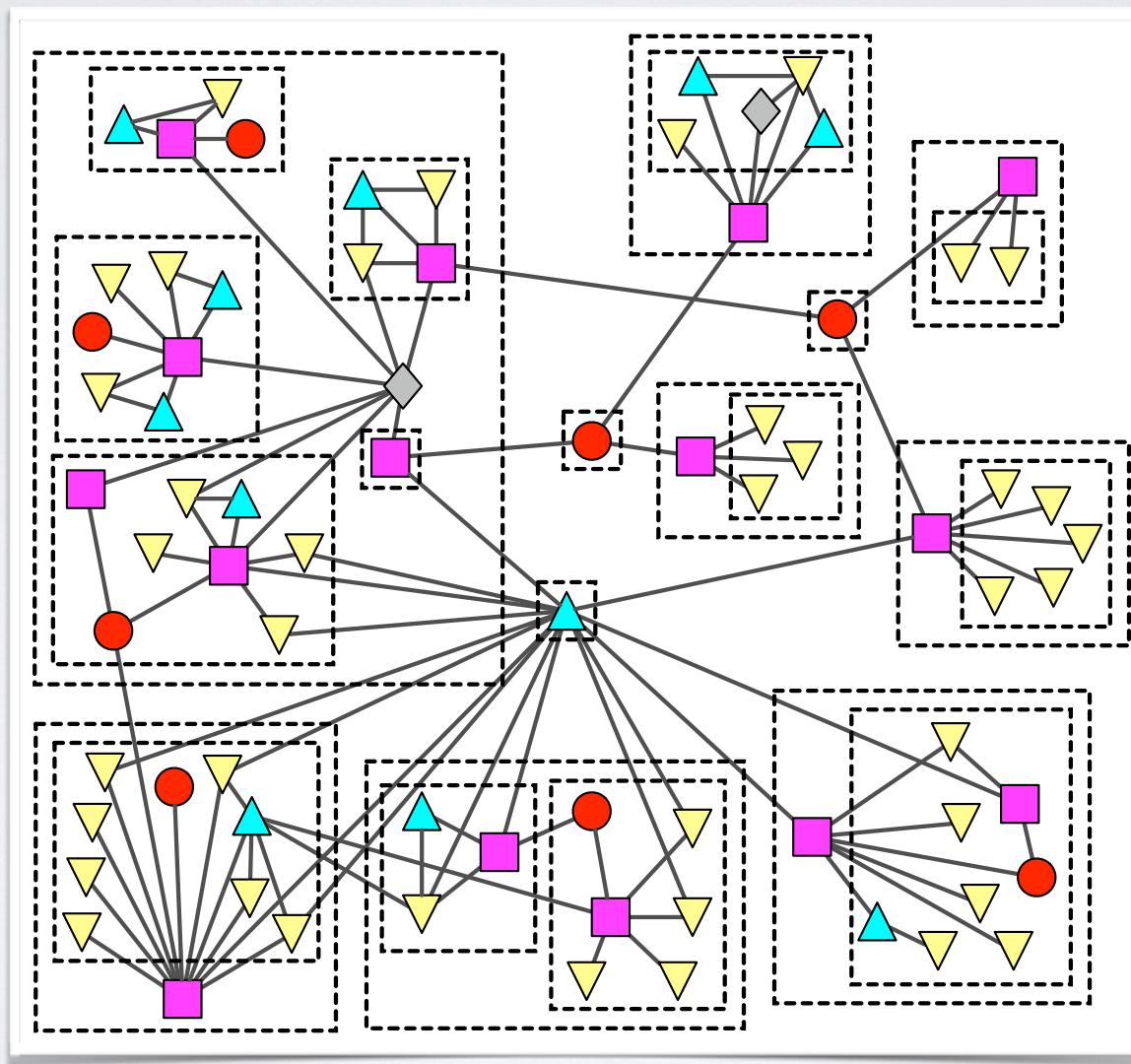
hierarchical structure



hierarchical structure



hierarchical structure

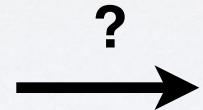
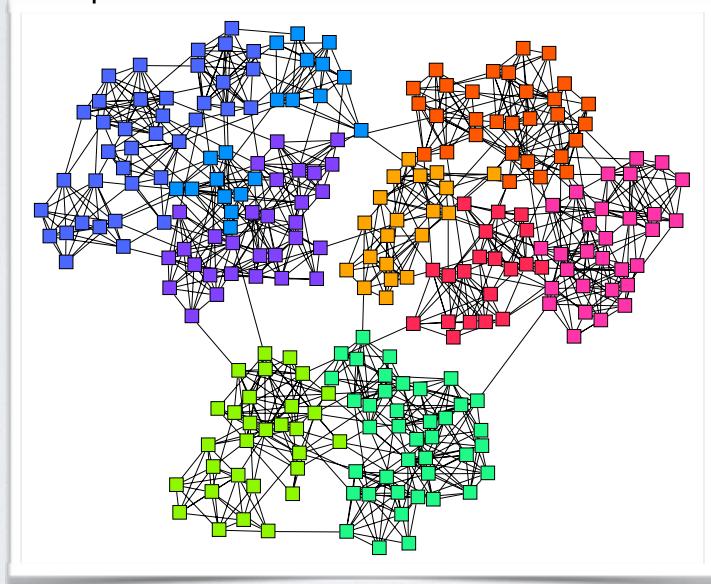


nested
modules

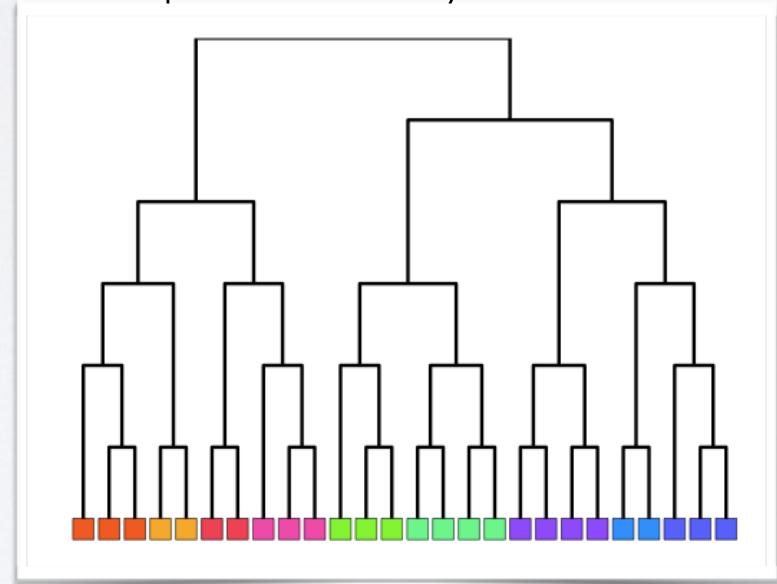
hierarchical structure

can we automatically extract hierarchies?

step 1: network data



step 3: hierarchy

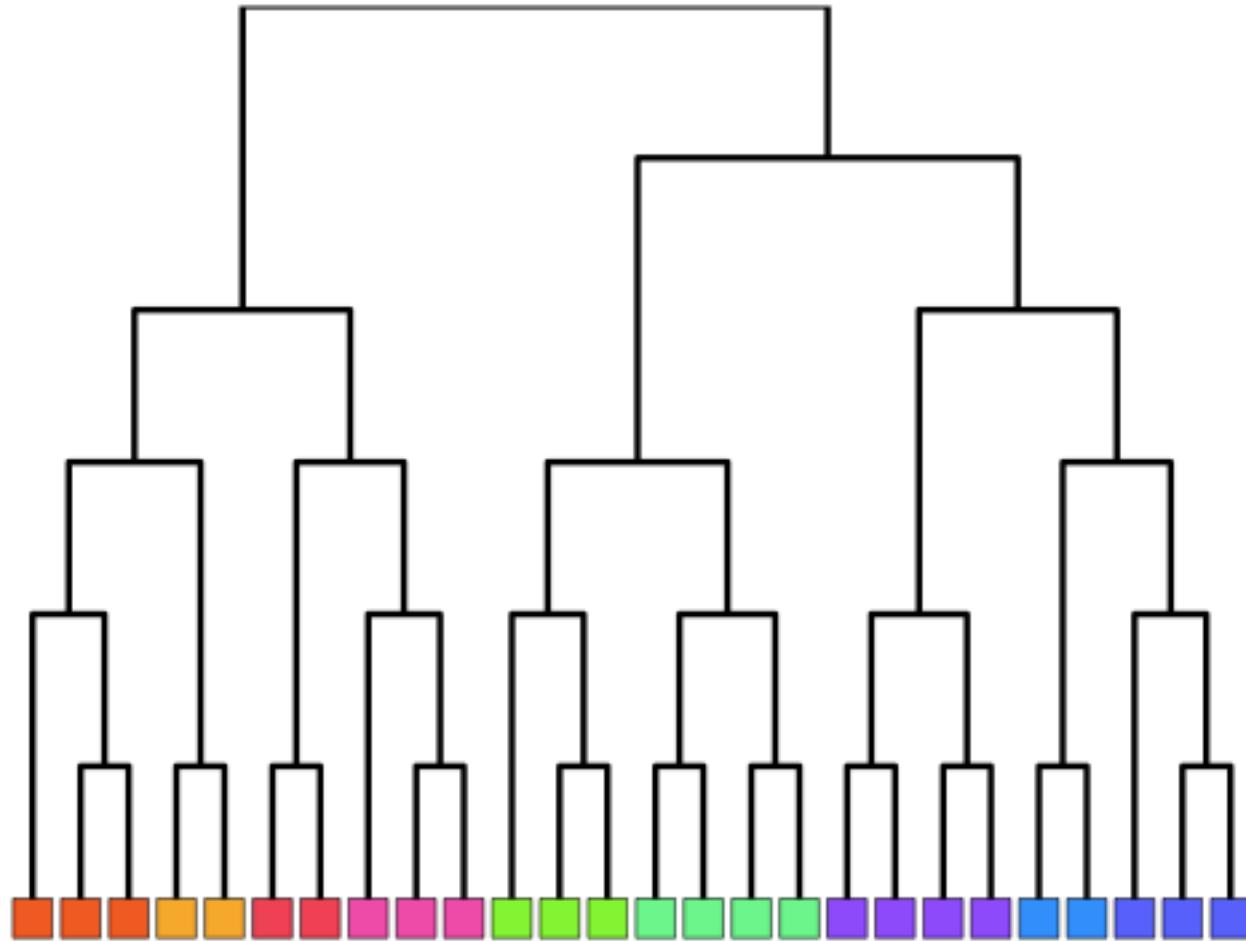


learning from network data

a direct approach

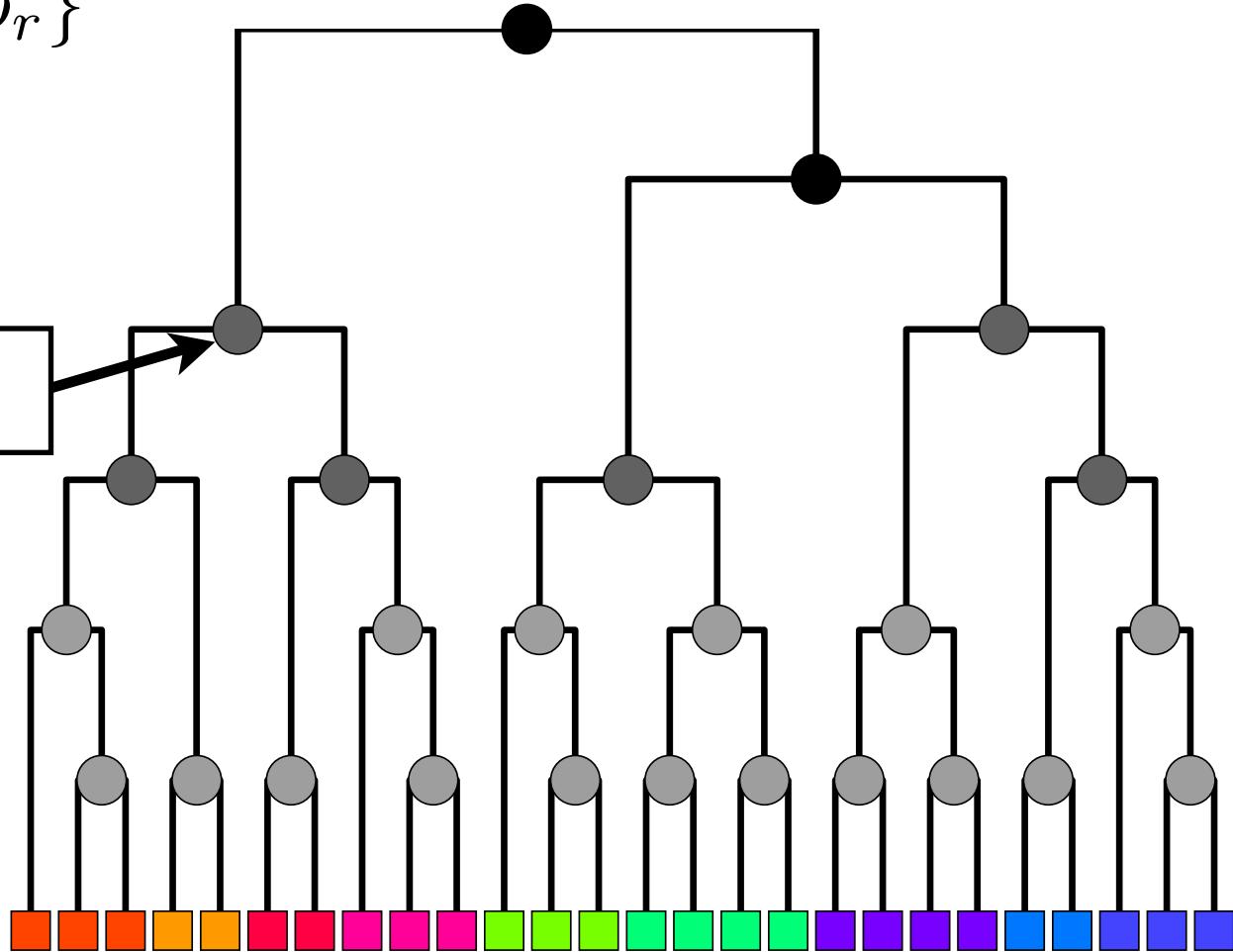
- write down (edge) generative model $\Pr(G \mid \theta)$
- sample models from posterior via MCMC $\rightarrow \{\theta_i\}$
- use ensemble $\{\theta_i\}$ to test fit, make predictions

\mathcal{D}

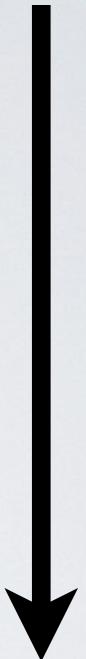


$\mathcal{D}, \{p_r\}$

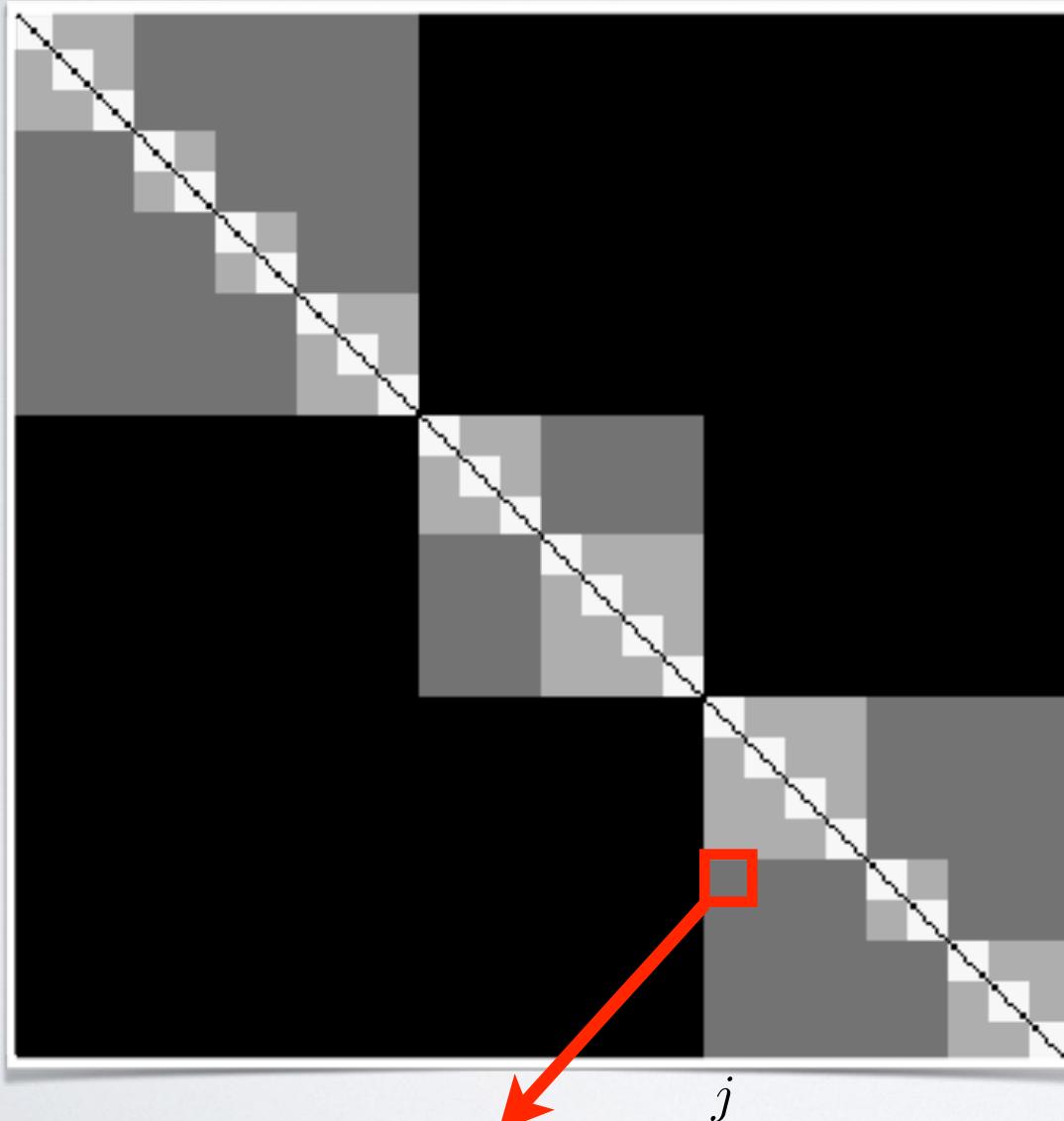
probability p_r



assortative modules



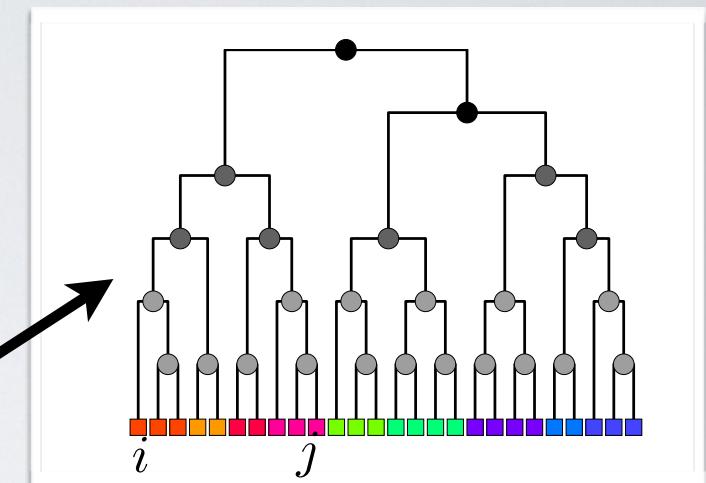
“inhomogeneous” random graph



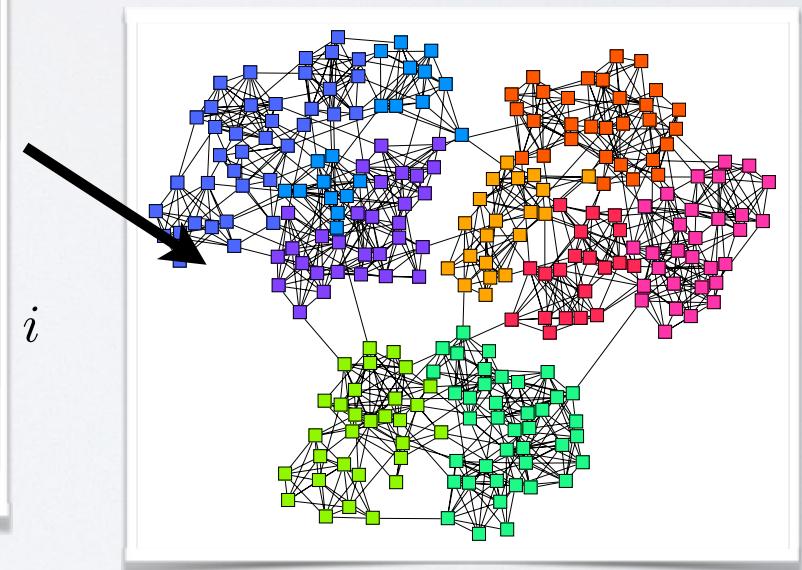
$$\Pr(i, j \text{ connected}) = p_r$$

$$= p_{(\text{lowest common ancestor of } i, j)}$$

model



instance



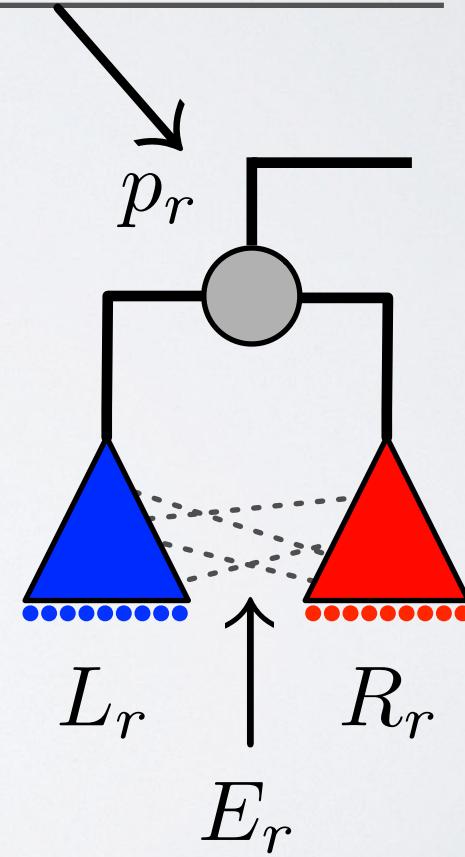
likelihood function

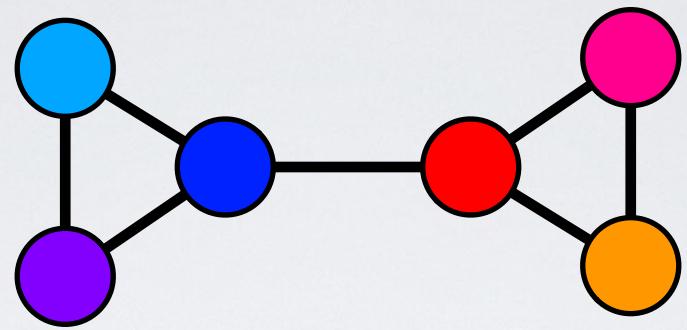
$$\mathcal{L}(\mathcal{D}, \{p_r\}) = \prod_r p_r^{E_r} (1 - p_r)^{L_r R_r - E_r}$$

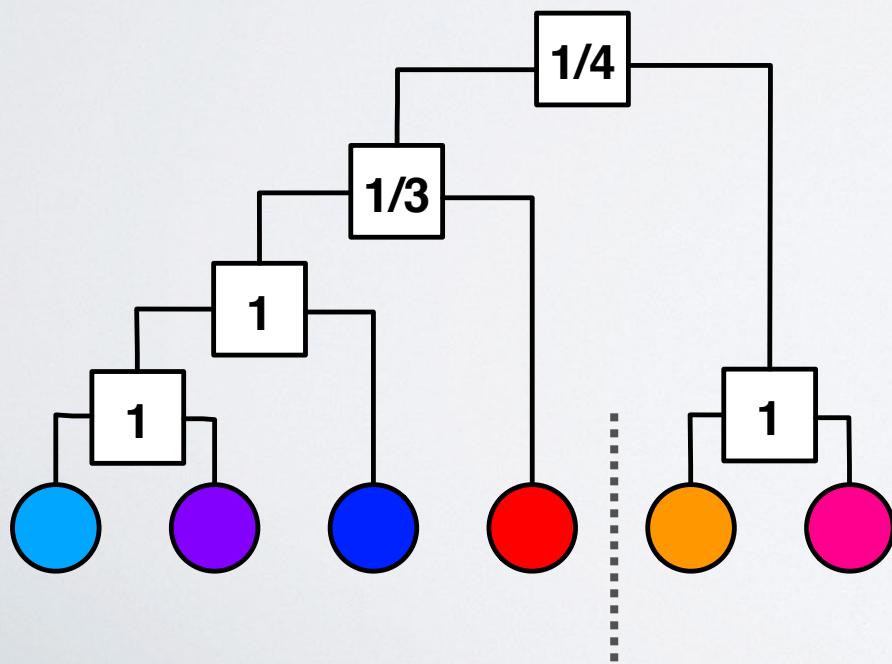
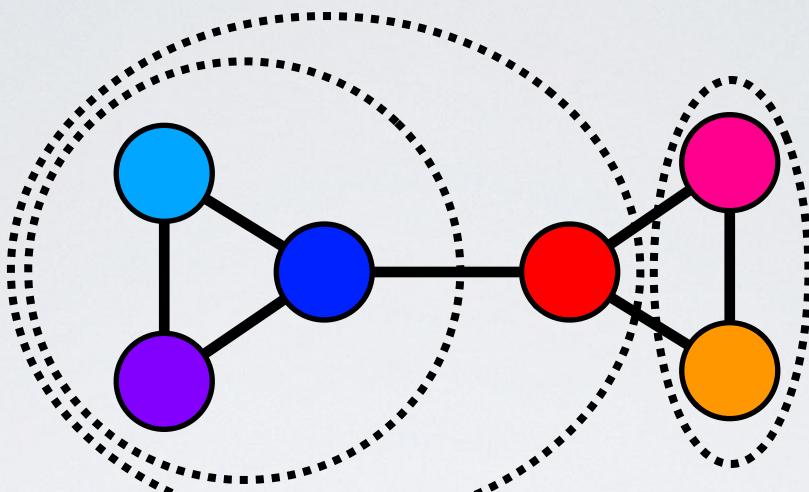
L_r = number nodes in left subtree

R_r = number nodes in right subtree

E_r = number edges with r as lowest common ancestor



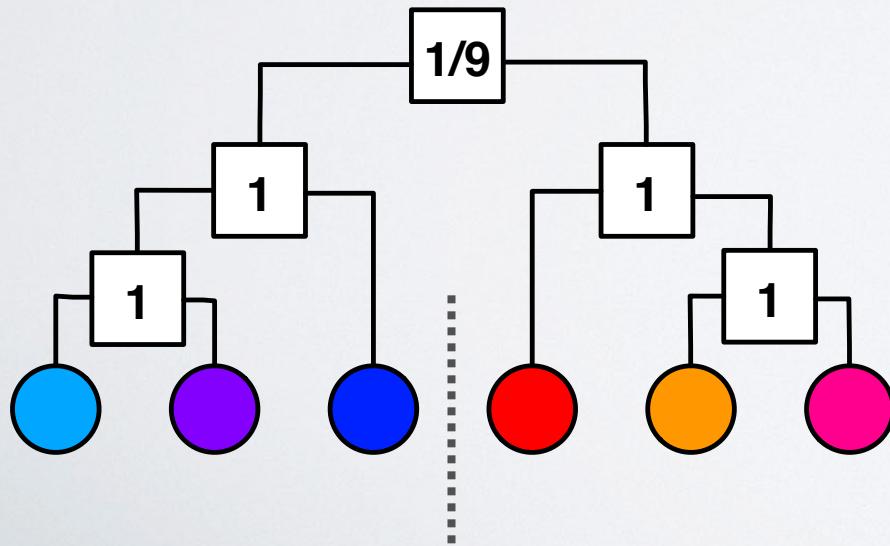
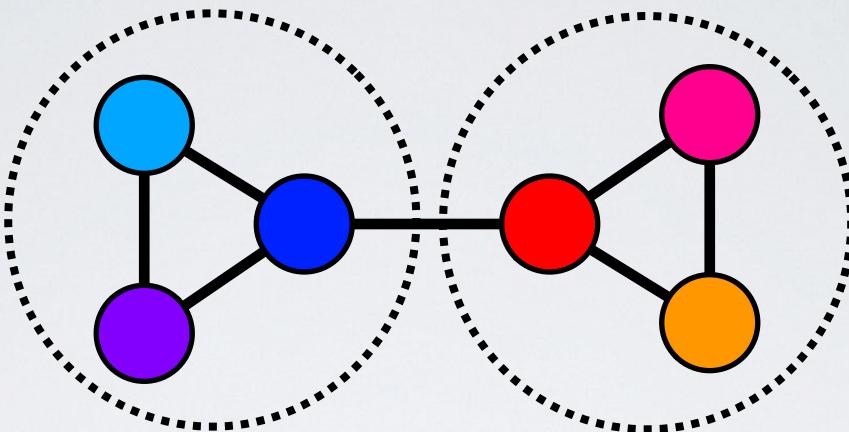




$$\mathcal{L}(\mathcal{D}, \{p_r\}) = \prod_r p_r^{E_r} (1 - p_r)^{L_r R_r - E_r}$$

$$\mathcal{L} = \left[\left(\frac{1}{3} \right)^1 \left(\frac{2}{3} \right)^2 \right] \cdot \left[\left(\frac{1}{4} \right)^2 \left(\frac{3}{4} \right)^6 \right]$$

$$\mathcal{L} = 0.0016$$



$$\mathcal{L} = \left[\left(\frac{1}{9} \right)^1 \left(\frac{8}{9} \right)^8 \right]$$

$$\mathcal{L} = 0.0433$$

$$\mathcal{L}(\mathcal{D}, \{p_r\}) = \prod_r p_r^{E_r} (1 - p_r)^{L_r R_r - E_r}$$

generalizing from a single example

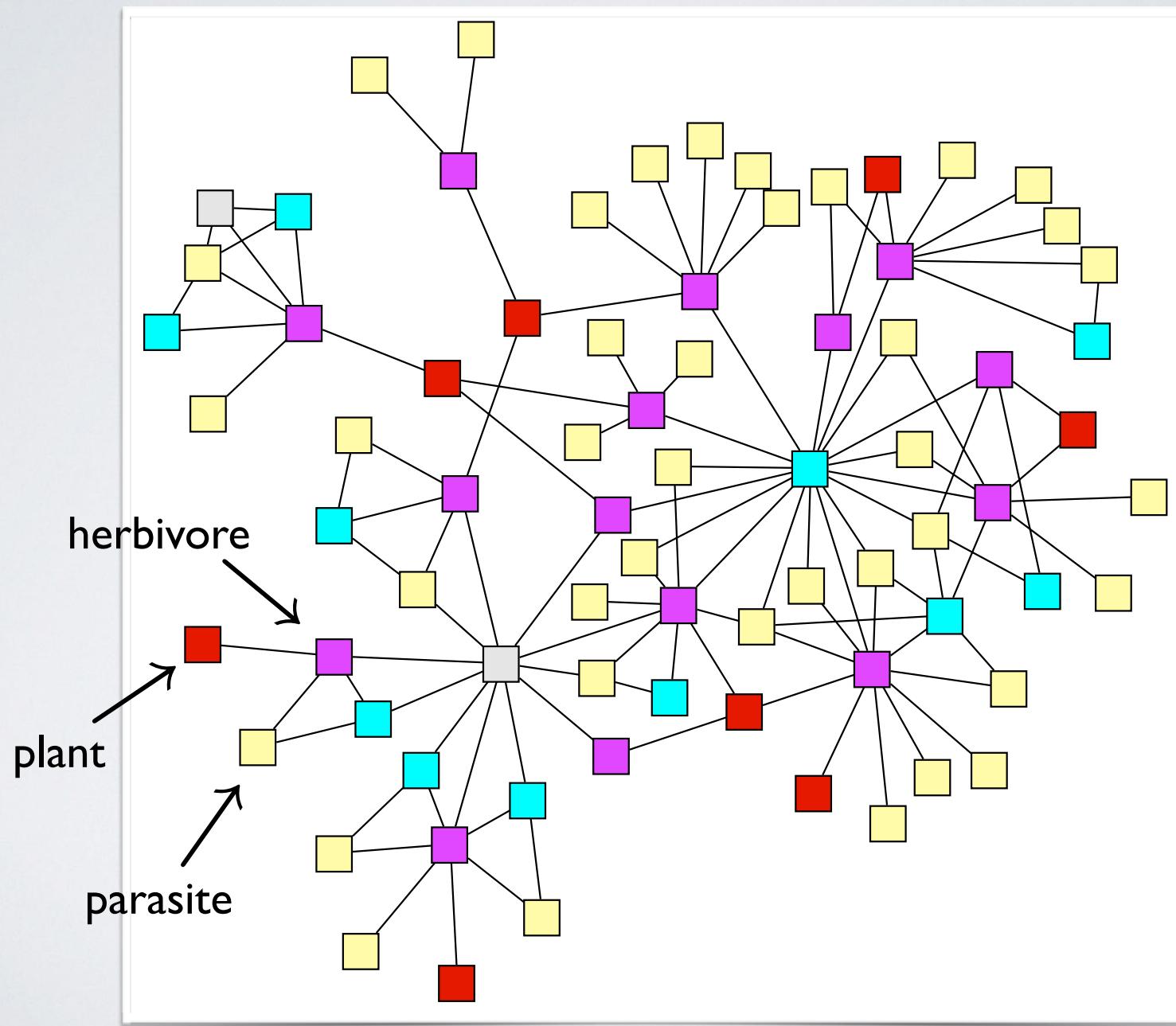
generalizing from a single example

- given graph G
- run MCMC to equilibrium
- for each sampled \mathcal{D} , draw a new graph G' from ensemble

checking our models:

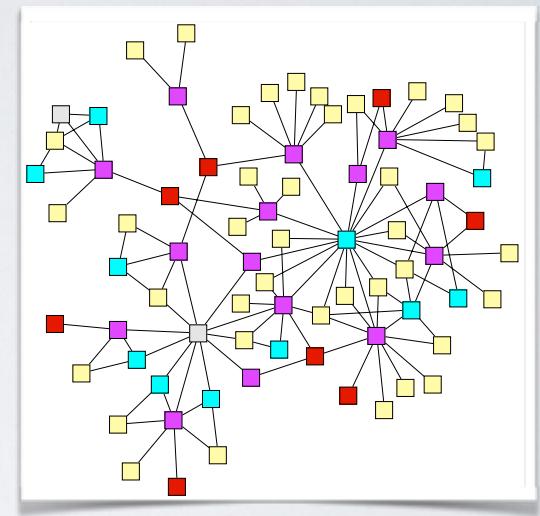
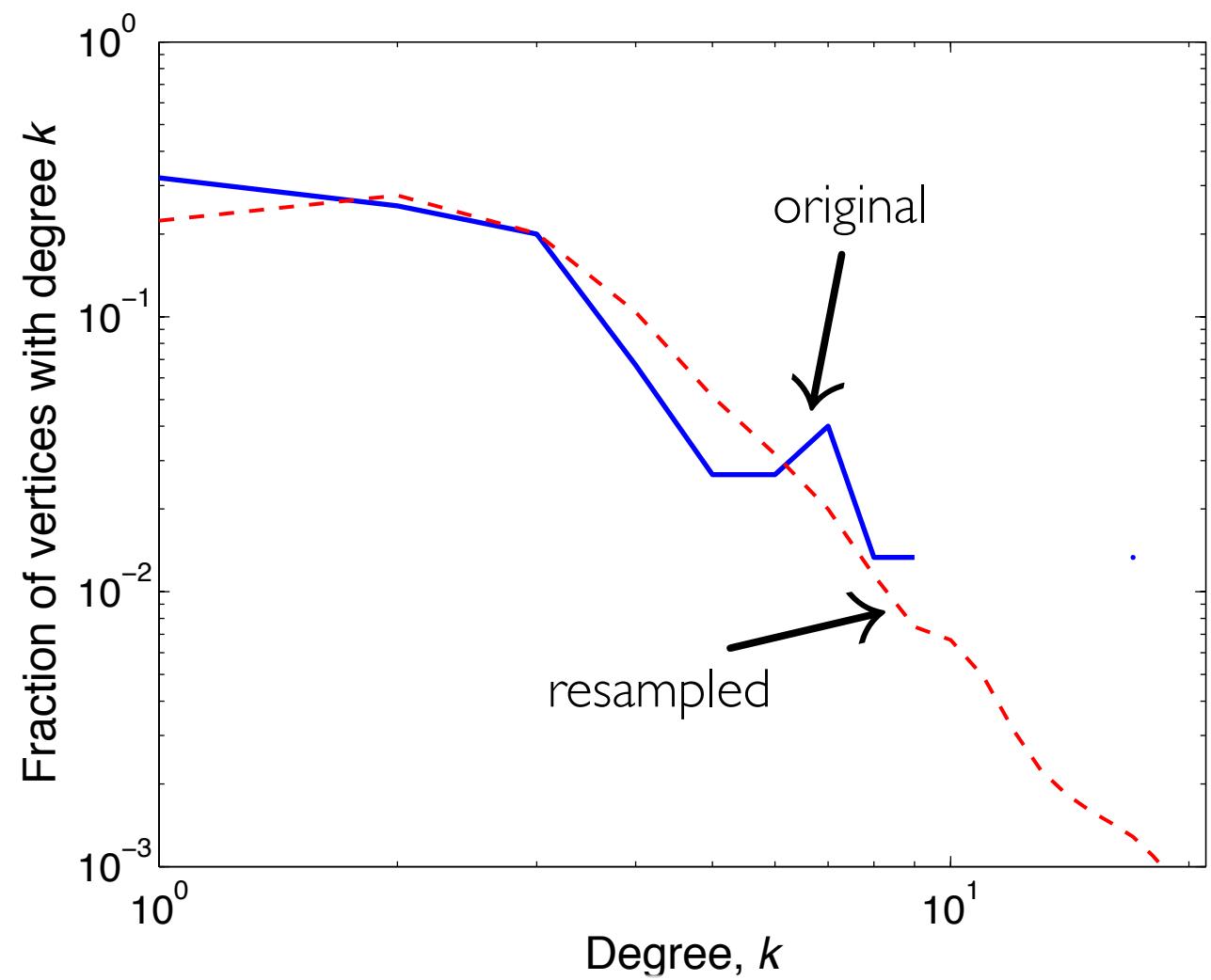
do resampled graphs look like original?

generalizing from a single example



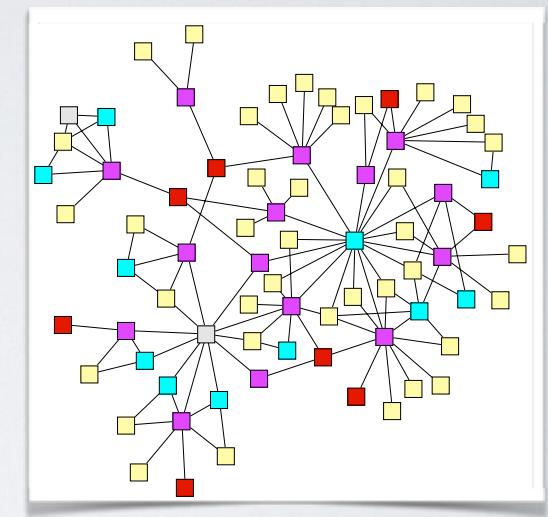
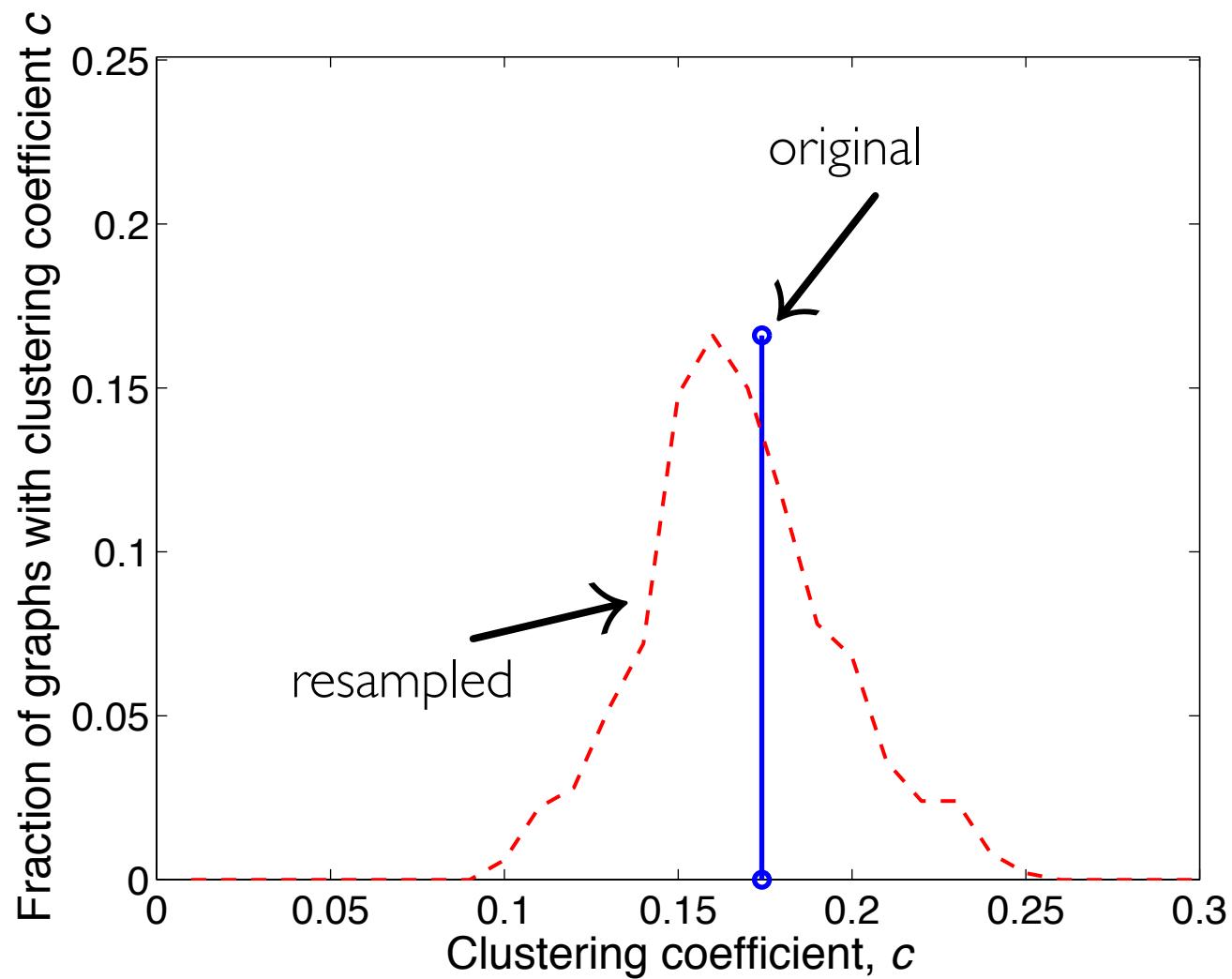
generalizing from a single example

degree distribution



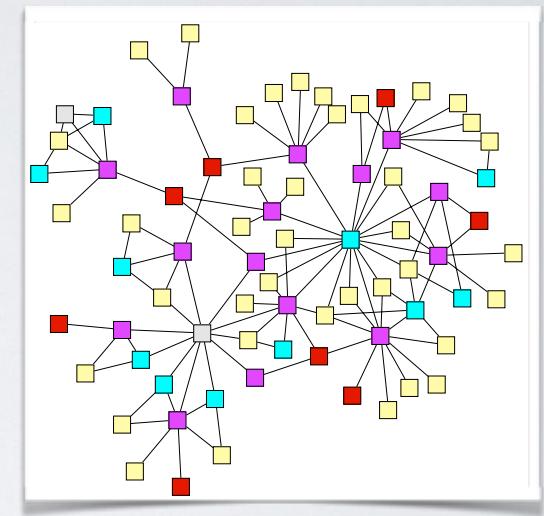
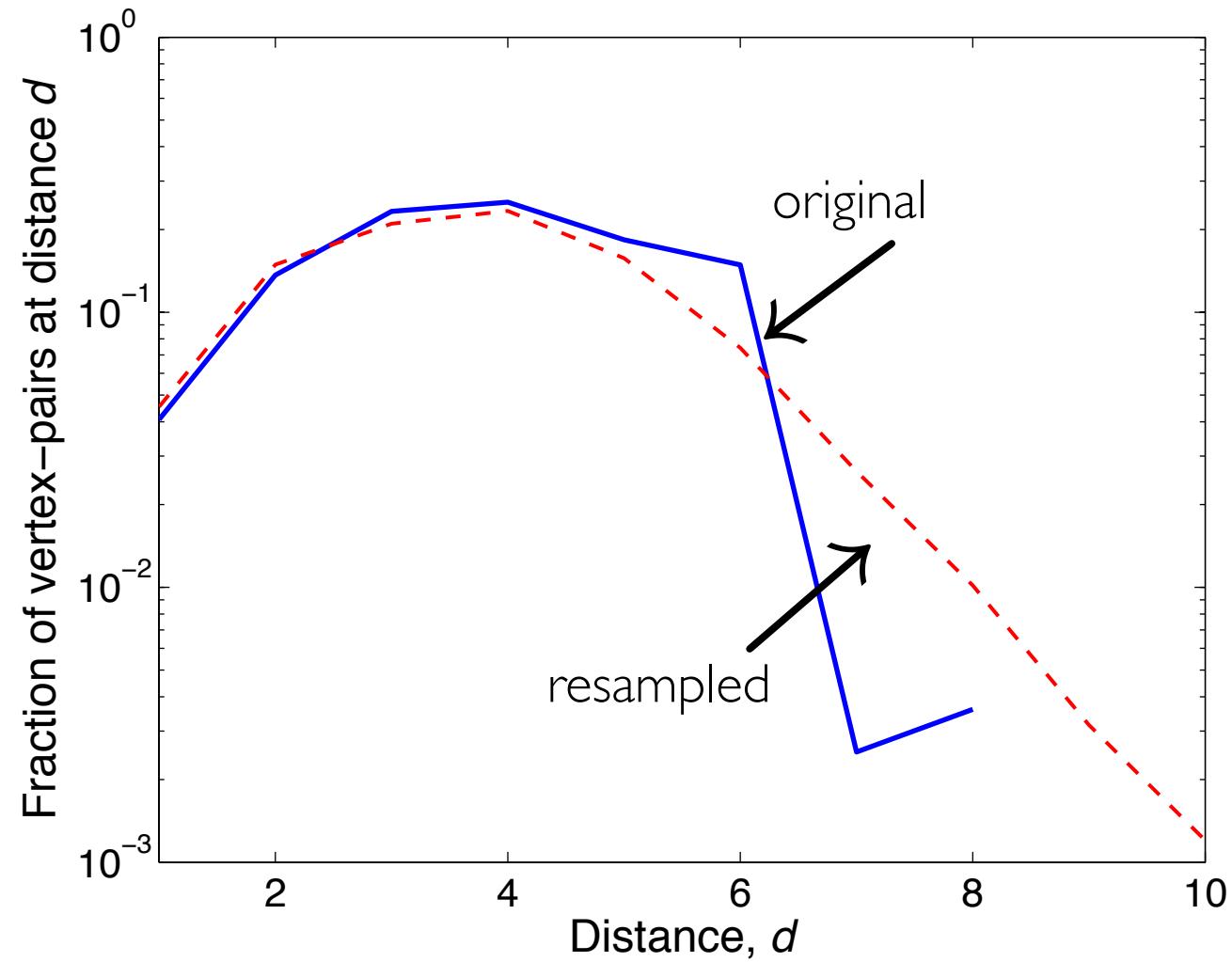
generalizing from a single example

density of triangles



generalizing from a single example

geodesic distances



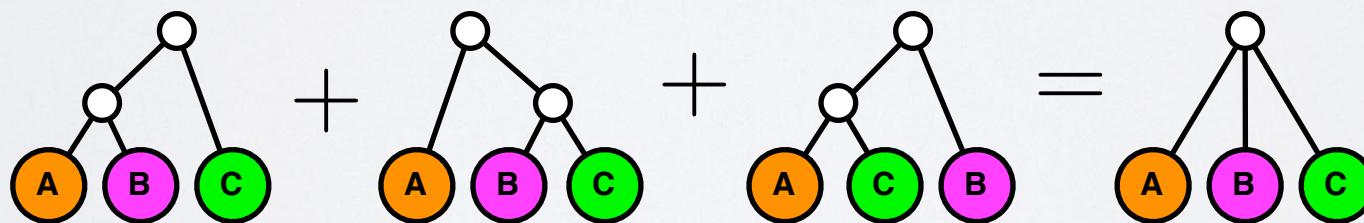
consensus hierarchies

MCMC produces a sample of hierarchies $\{\mathcal{D}_i\}$

how can we combine these trees into a **consensus tree**?

borrow technique from phylogenetic tree reconstruction:

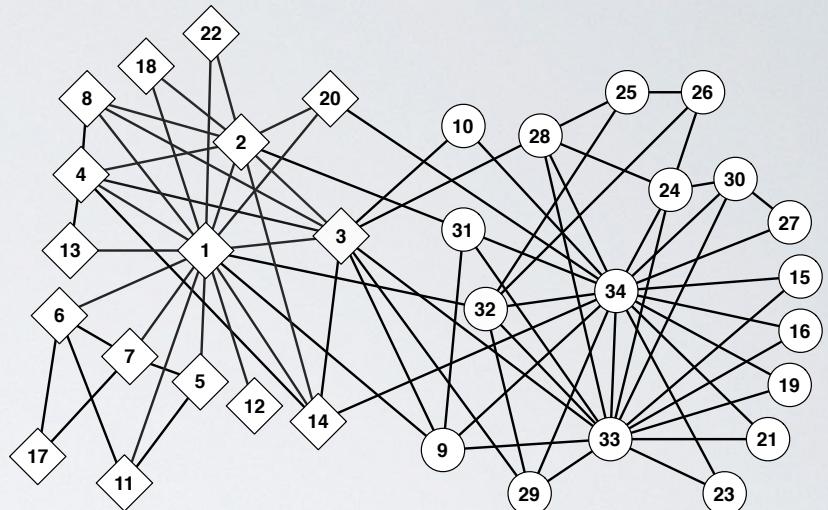
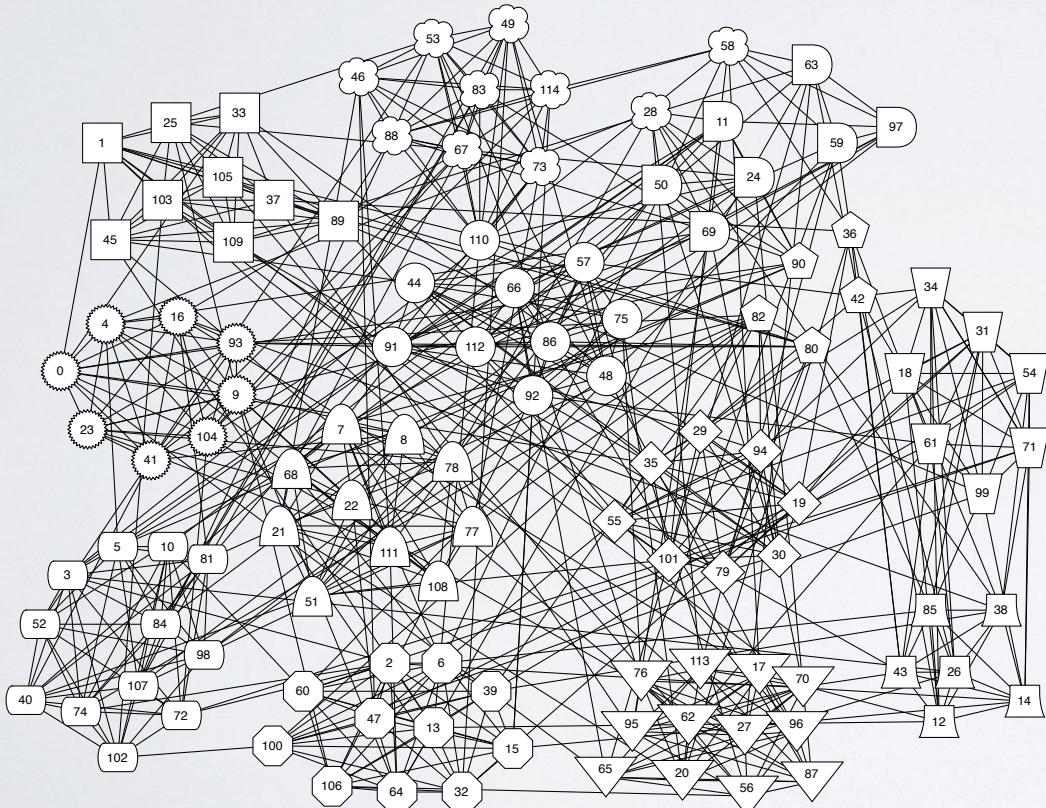
- extract set of hierarchical relationships \mathcal{D}_c contained in majority of sampled hierarchies $\{\mathcal{D}_i\}$
- for instance:



consensus hierarchies

NCAA Schedule 2000

$n = 115$ $m = 613$



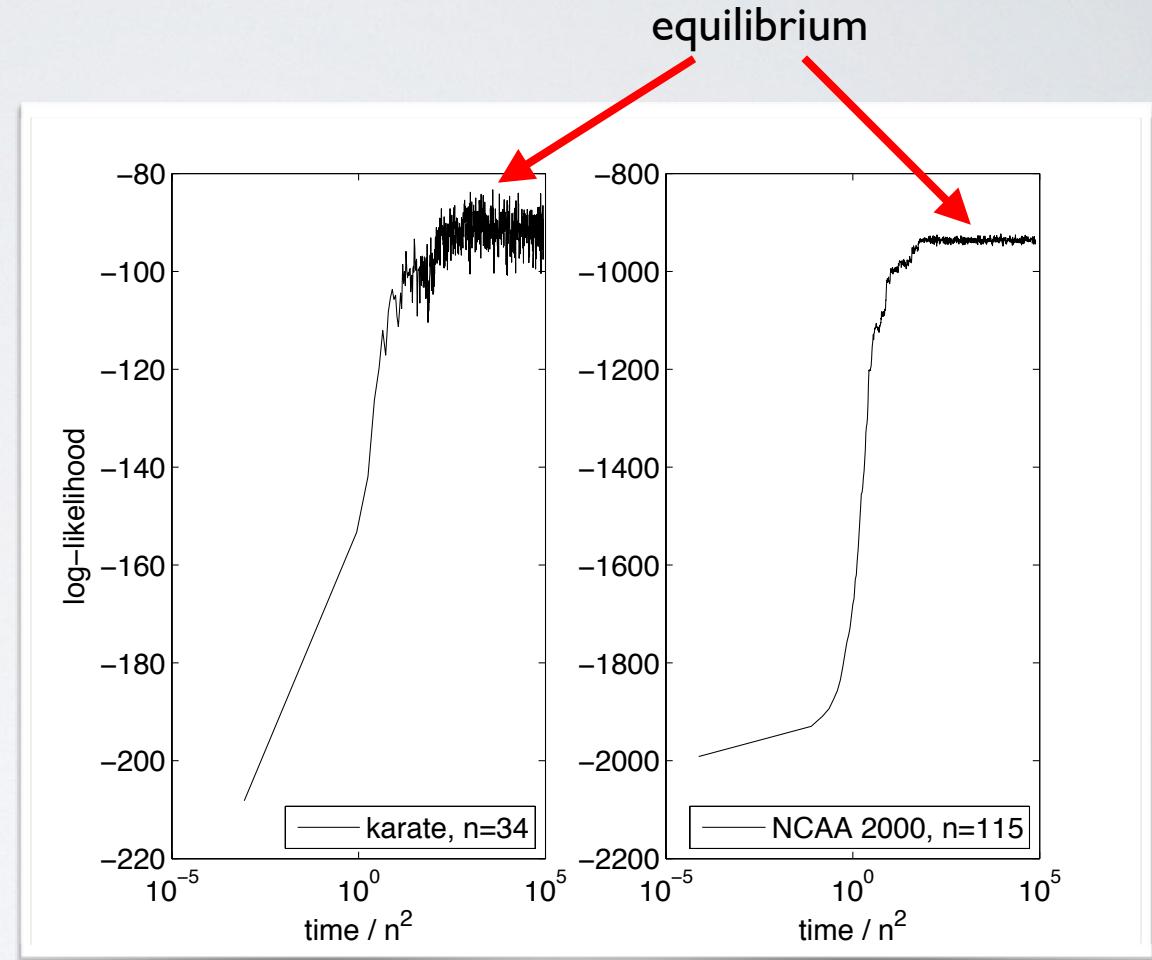
Zachary's Karate Club

$n = 34$ $m = 78$

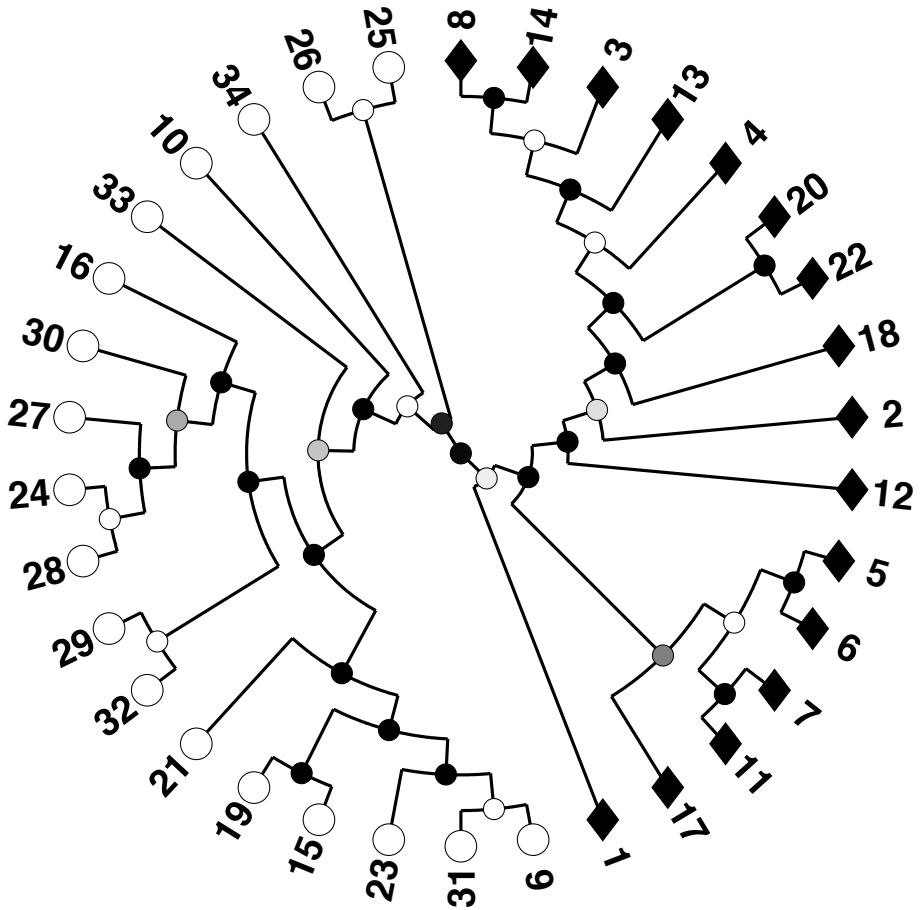
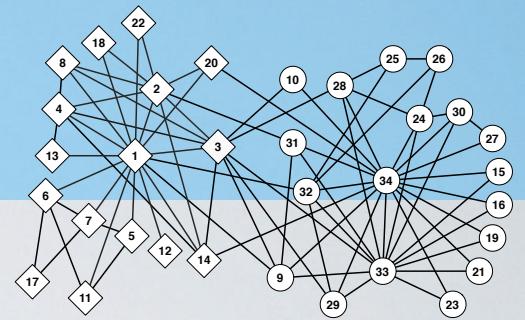
consensus hierarchies

MCMC mixes relatively quickly

equilibrium in $\approx n^2$ steps

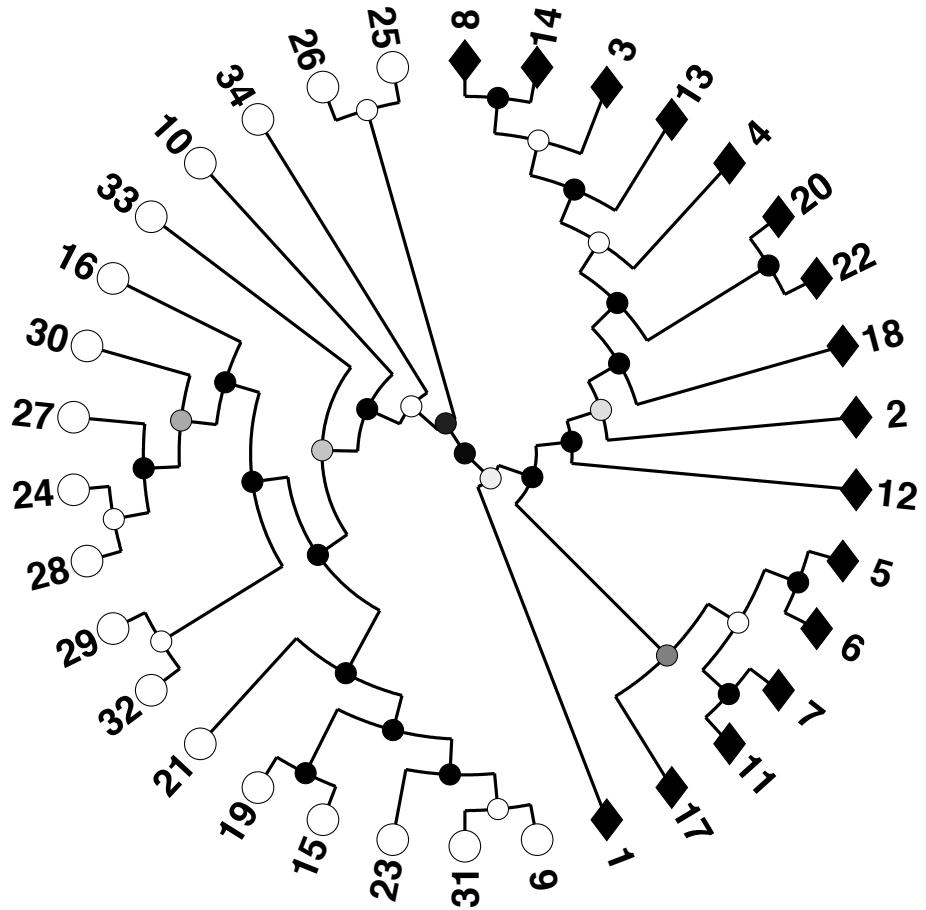
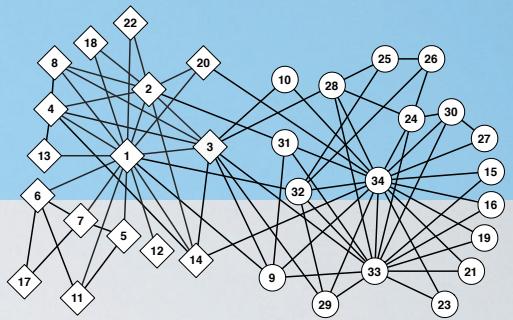


consensus hierarchies

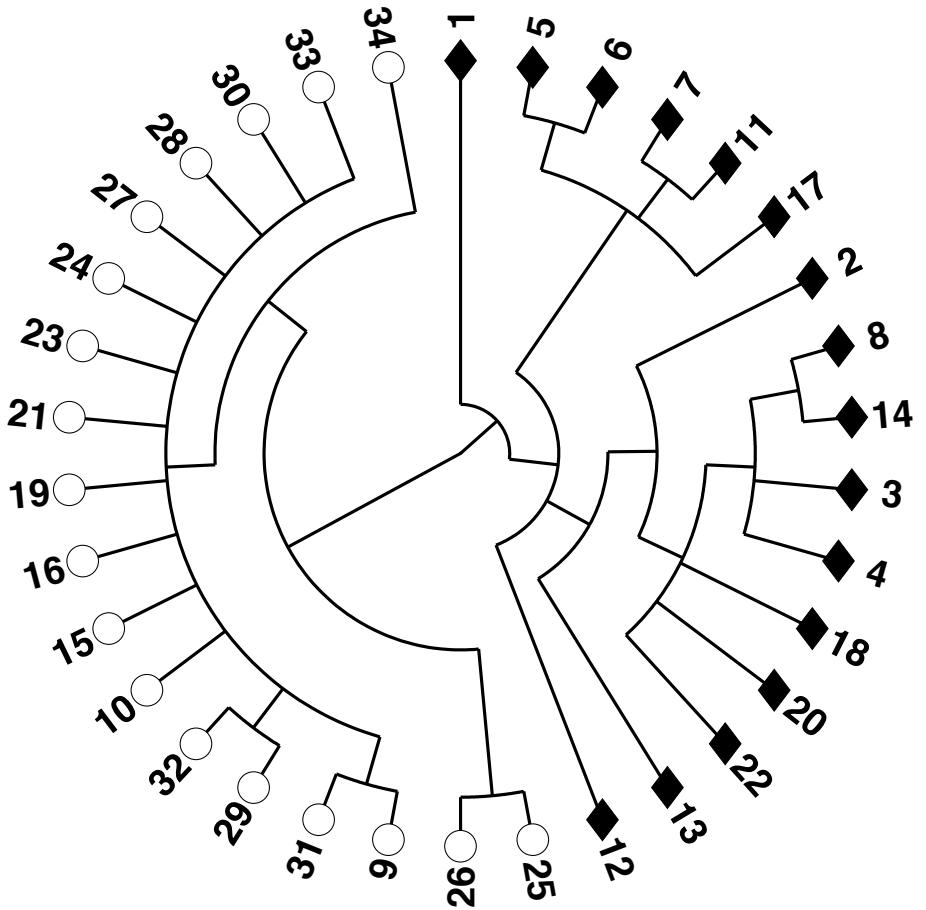


point estimate

consensus hierarchies

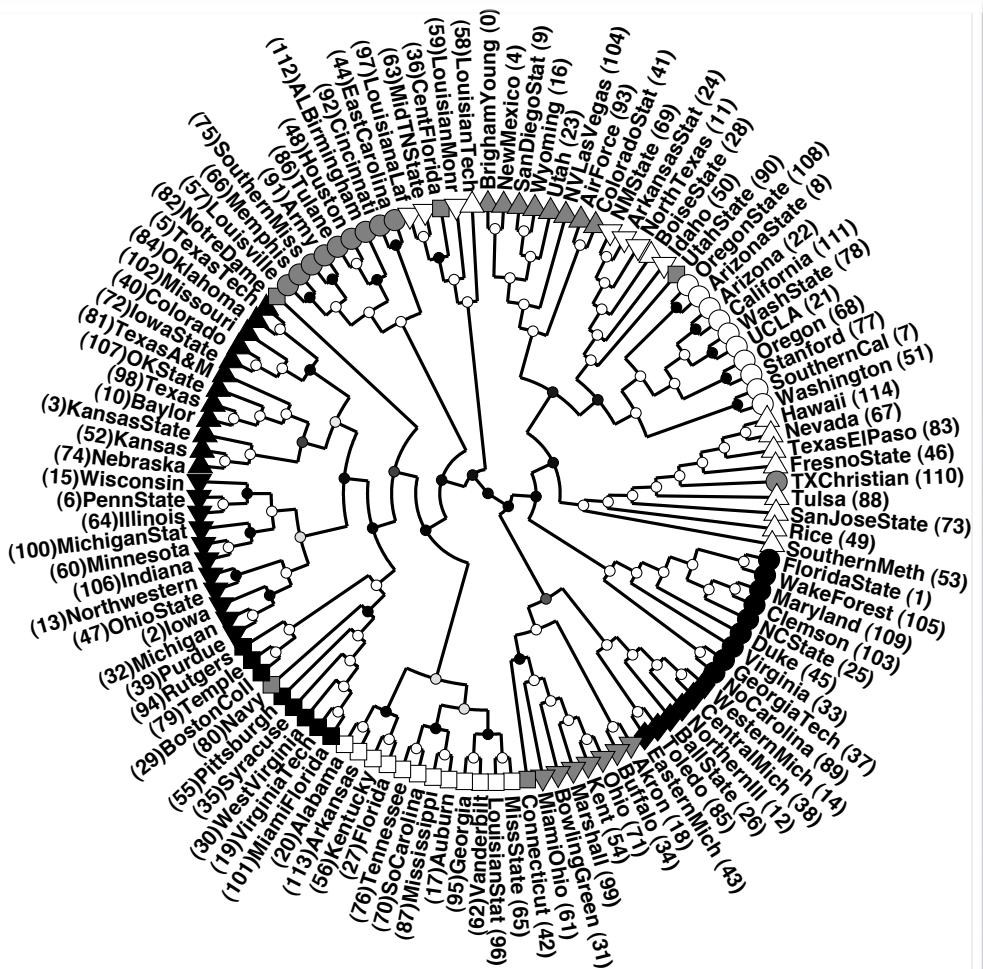


point estimate

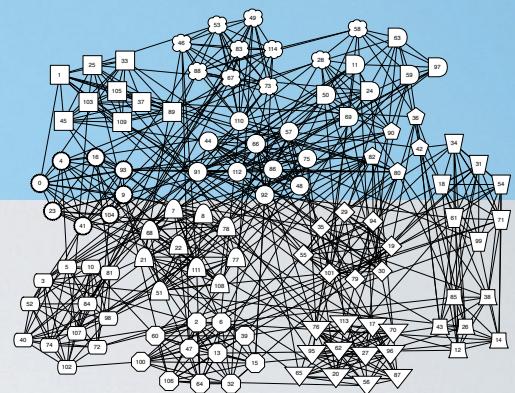


consensus hierarchy

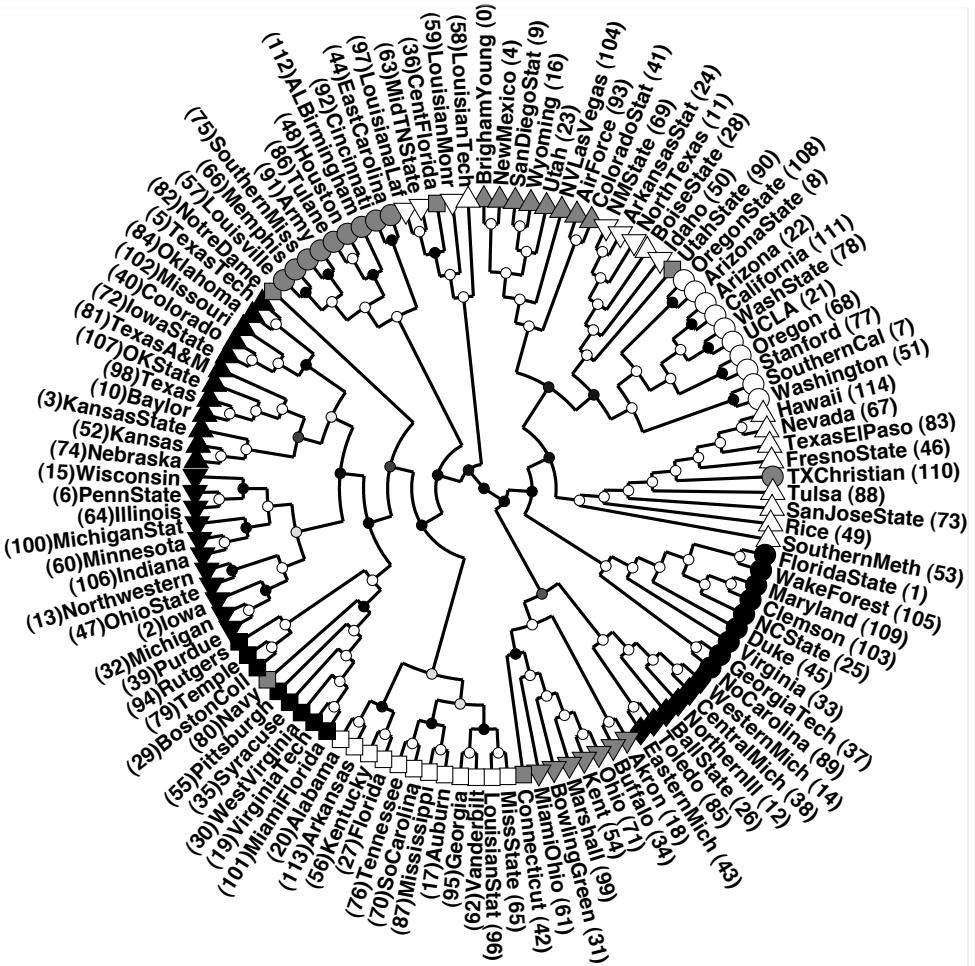
consensus hierarchies



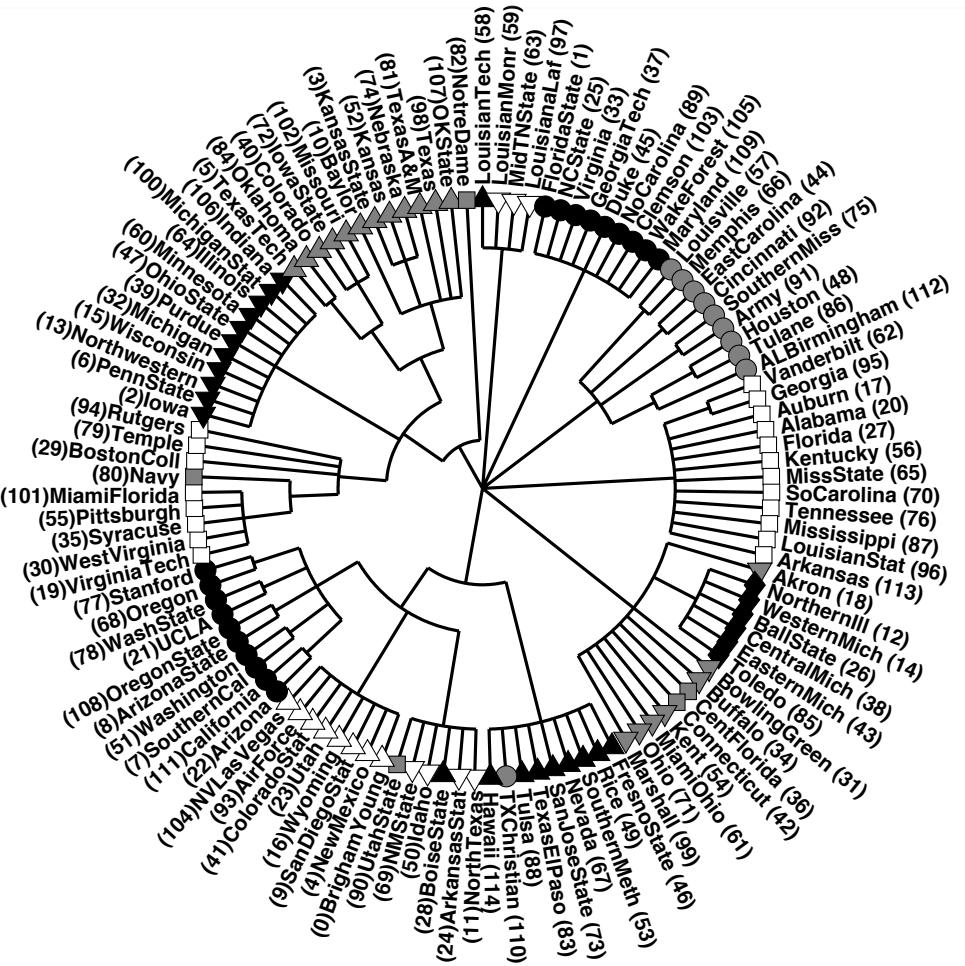
point estimate



consensus hierarchies

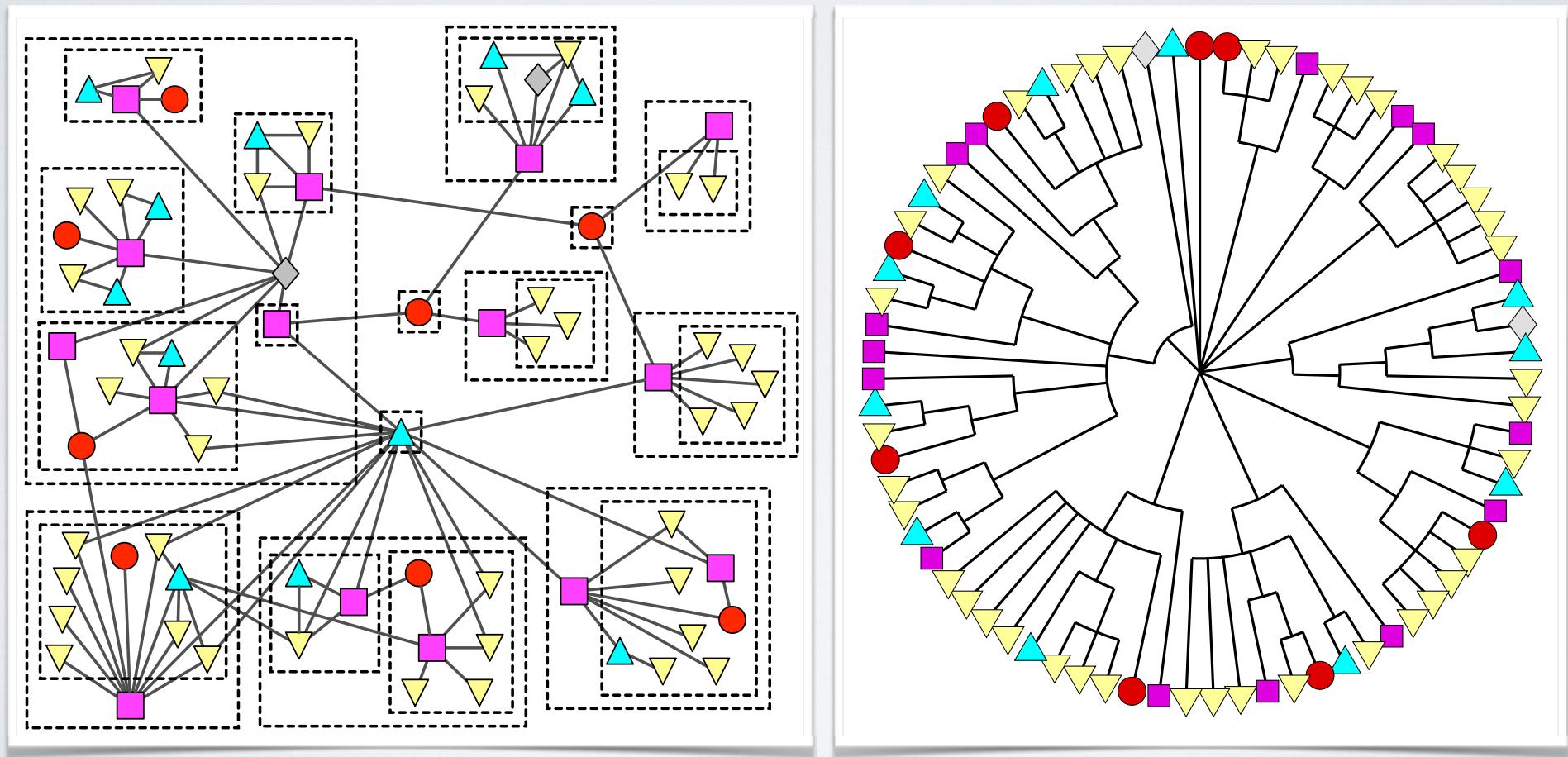


point estimate



consensus hierarchy

consensus hierarchies



predicting missing links

many networks partially known, noisy

- social nets, foodwebs, protein interactions, etc.

use generative model to predict missing links

[can do the same for spurious links]

other approaches

- Liben-Nowell & Kleinberg (2003)
- Goldberg & Roth (2003)
- Szilagy et al. (2005)
- and now many others

predicting missing links

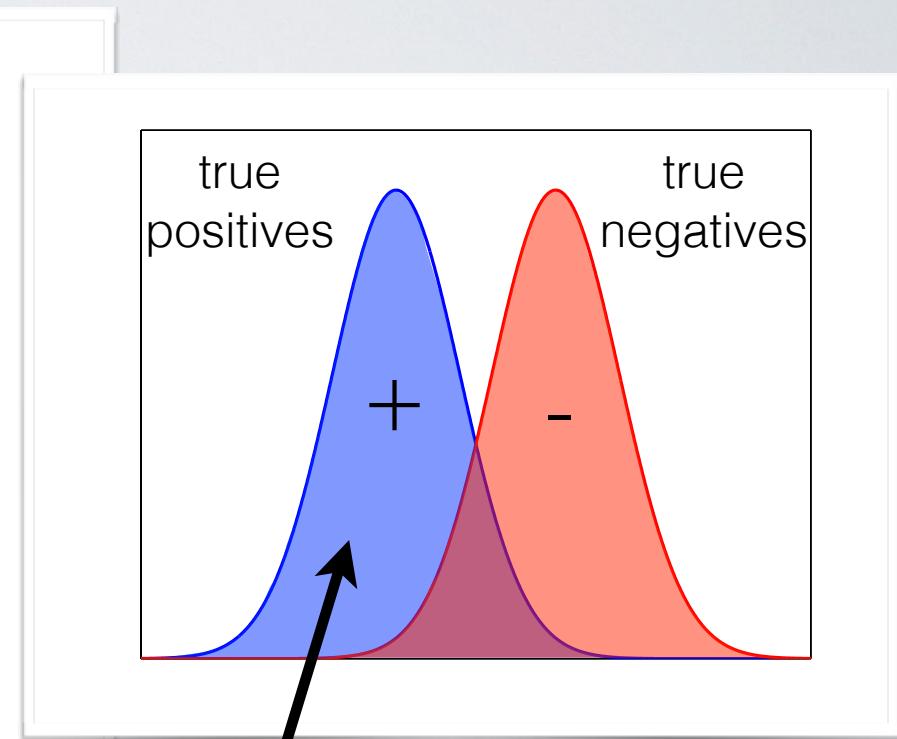
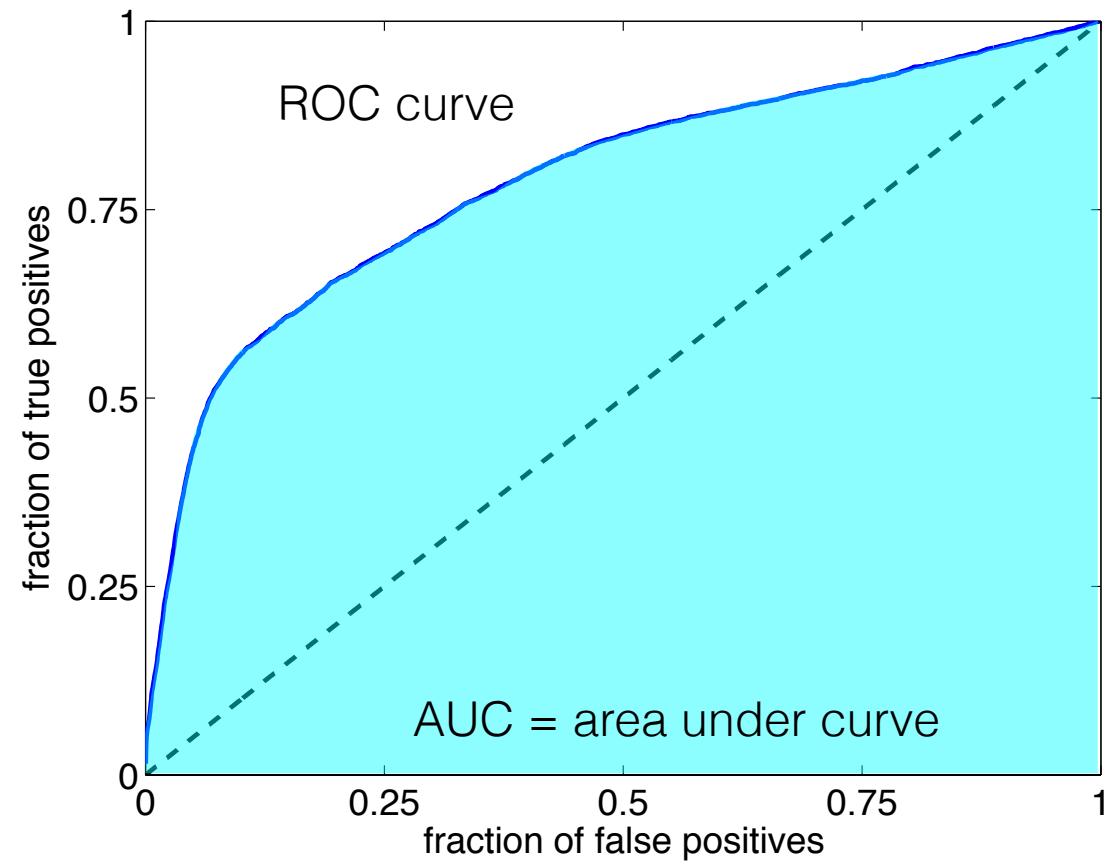
- given incomplete graph G
- run MCMC to equilibrium
- then, over sampled \mathcal{D} , compute average $\langle p_r \rangle$ for links $(i, j) \notin G$
- predict links with high $\langle p_r \rangle$ values are missing

compute AUC via leave-k-out (edges) cross-validation

$AUC = 1/2 \implies$ no better than chance

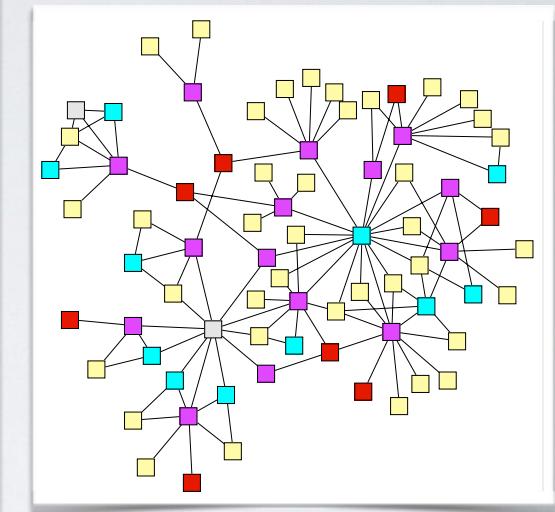
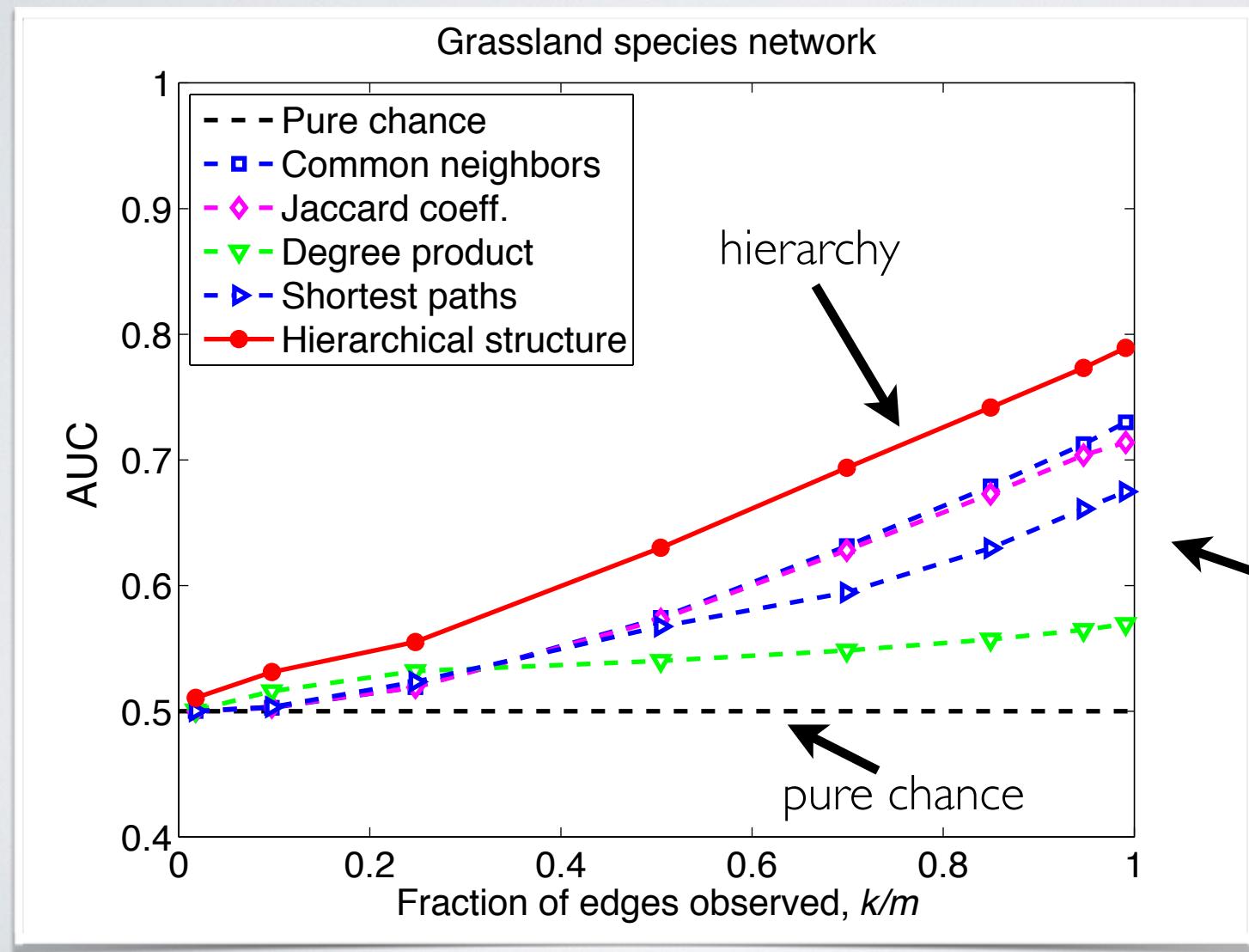
predicting missing links

scoring the predictions



$AUC = \Pr(\text{ distinguish } + \text{ from } -)$

predicting missing links



simple predictors

hierarchical random graphs

- **multi-scale structural inference**
mixtures of assortative, disassortative groups
- **inference is expensive (MCMC)**
but likelihood function extremely rugged
- **hierarchies can explain low-level network patterns**
degrees, triangles, distances, etc.
- **link-prediction is a hard validation**

generative models for complex networks

3 examples

- hierarchical random graphs (HRG)
 - an application to temporal networks
- bipartite community structure (biSBM)
- weighted community structure (WSBM)

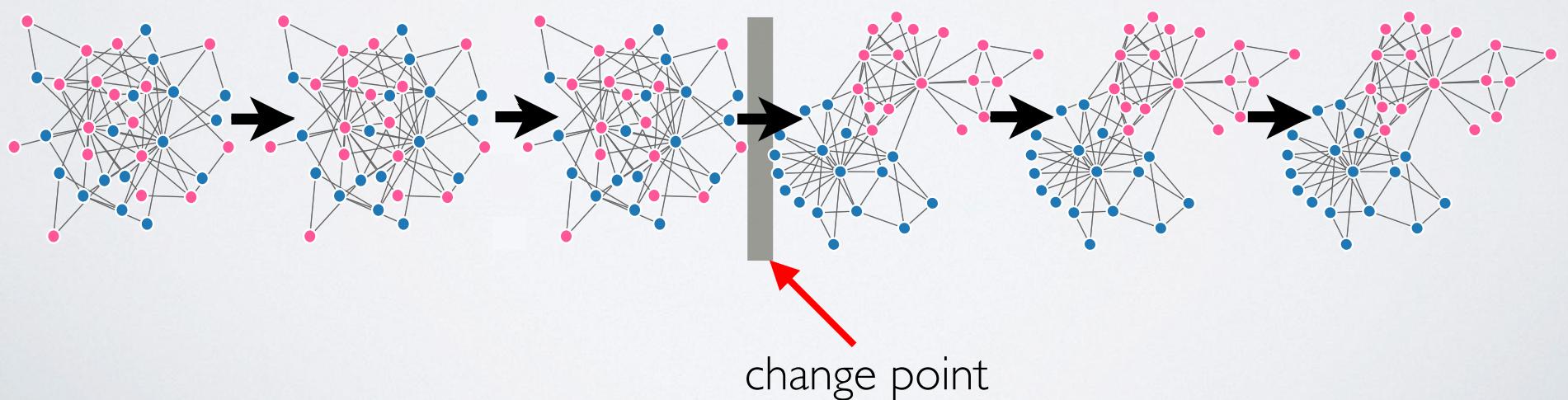
detecting network change points



temporal networks

Dr. Leto Peel
(Colorado)

- noise vs. structural "change points"
- can we detect when large-scale changes occur?



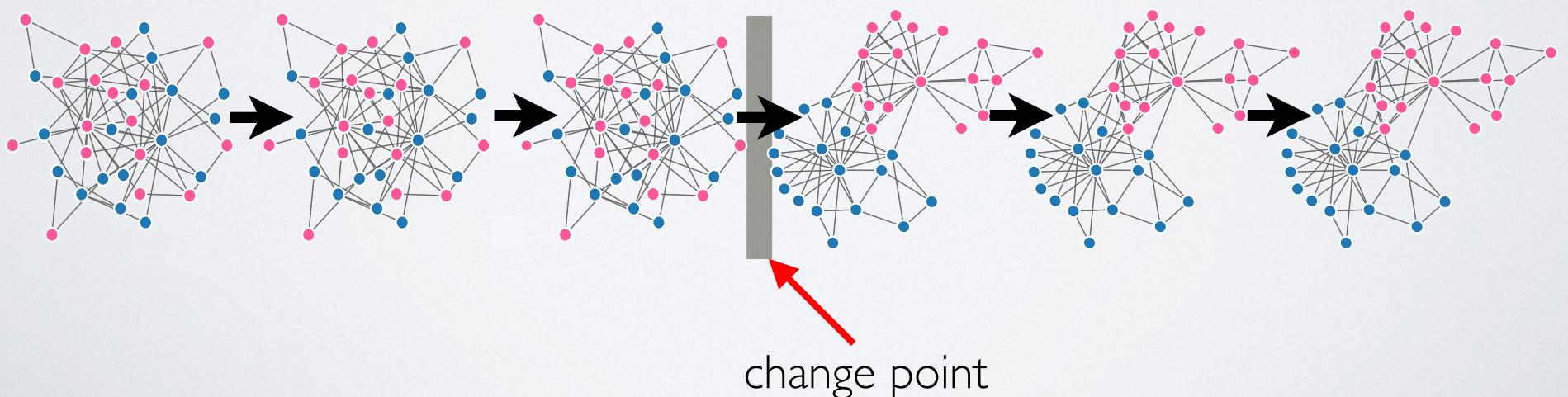
detecting network change points



Dr. Leto Peel
(Colorado)

temporal networks

- noise vs. structural "change points"
- can we detect when large-scale changes occur?
- apply generative model $\Pr(G | \theta)$ to sliding window of snapshots $\{G_t, G_{t+1}, \dots, G_{t+\tau}\}$
- detecting network change-points becomes detecting changes in model parameters



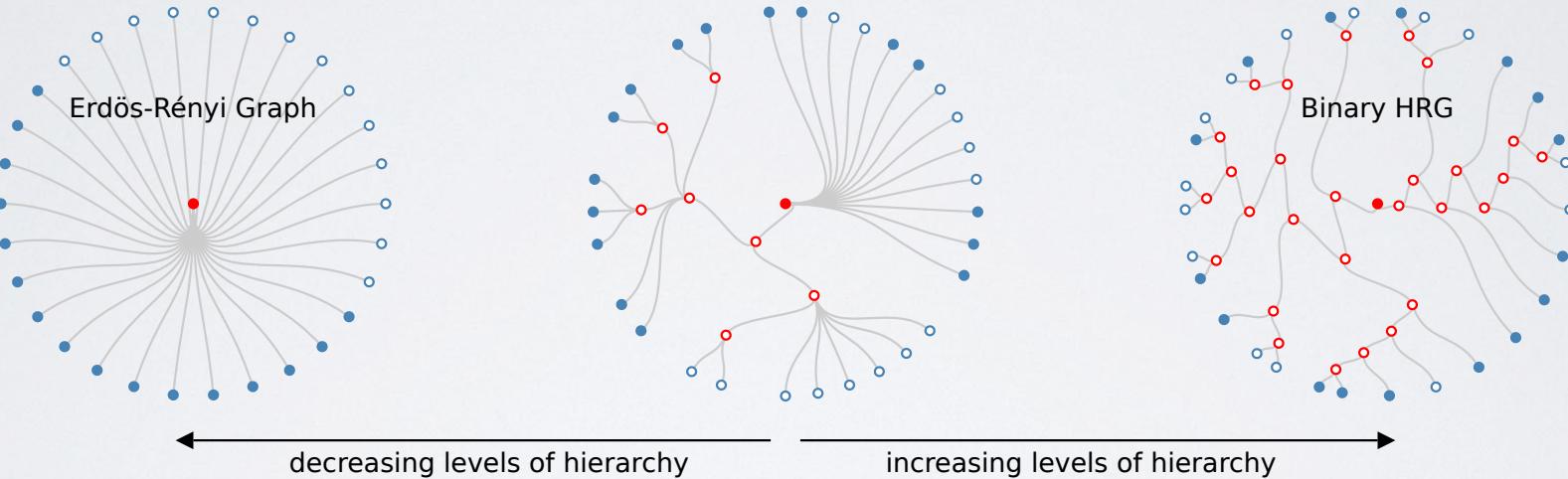
detecting network change points



Dr. Leto Peel
(Colorado)

generalized HRG

- like consensus HRG
- add Bayesian priors on the $\{p_r\}$



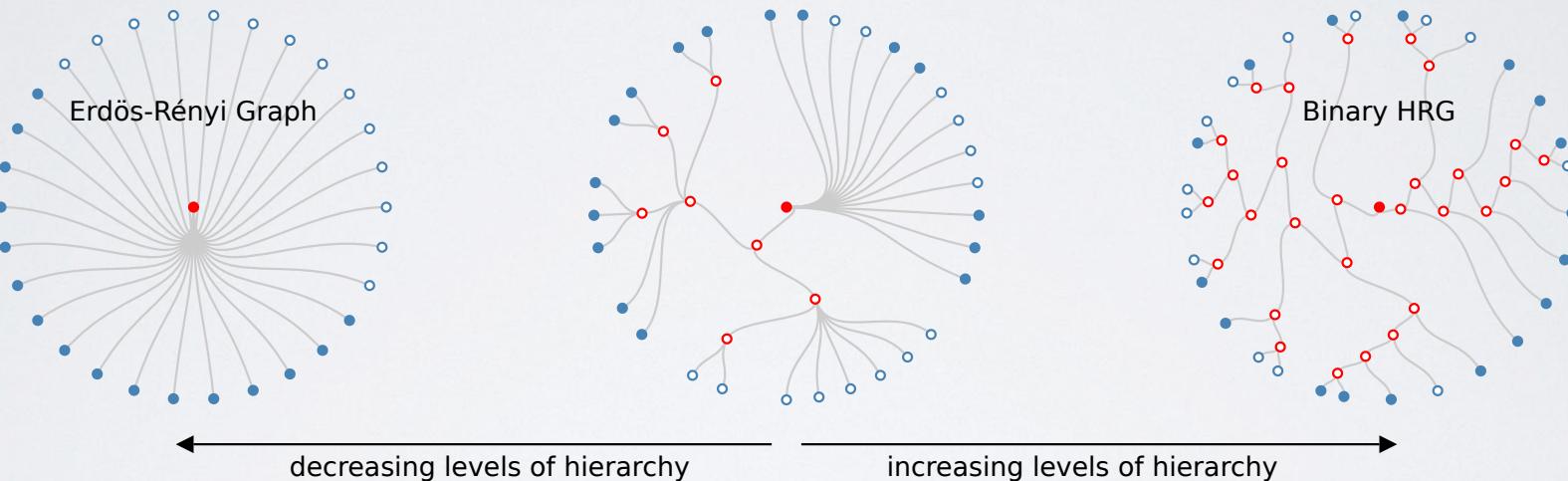
detecting network change points



Dr. Leto Peel
(Colorado)

generalized HRG

- like consensus HRG
- add Bayesian priors on the $\{p_r\}$



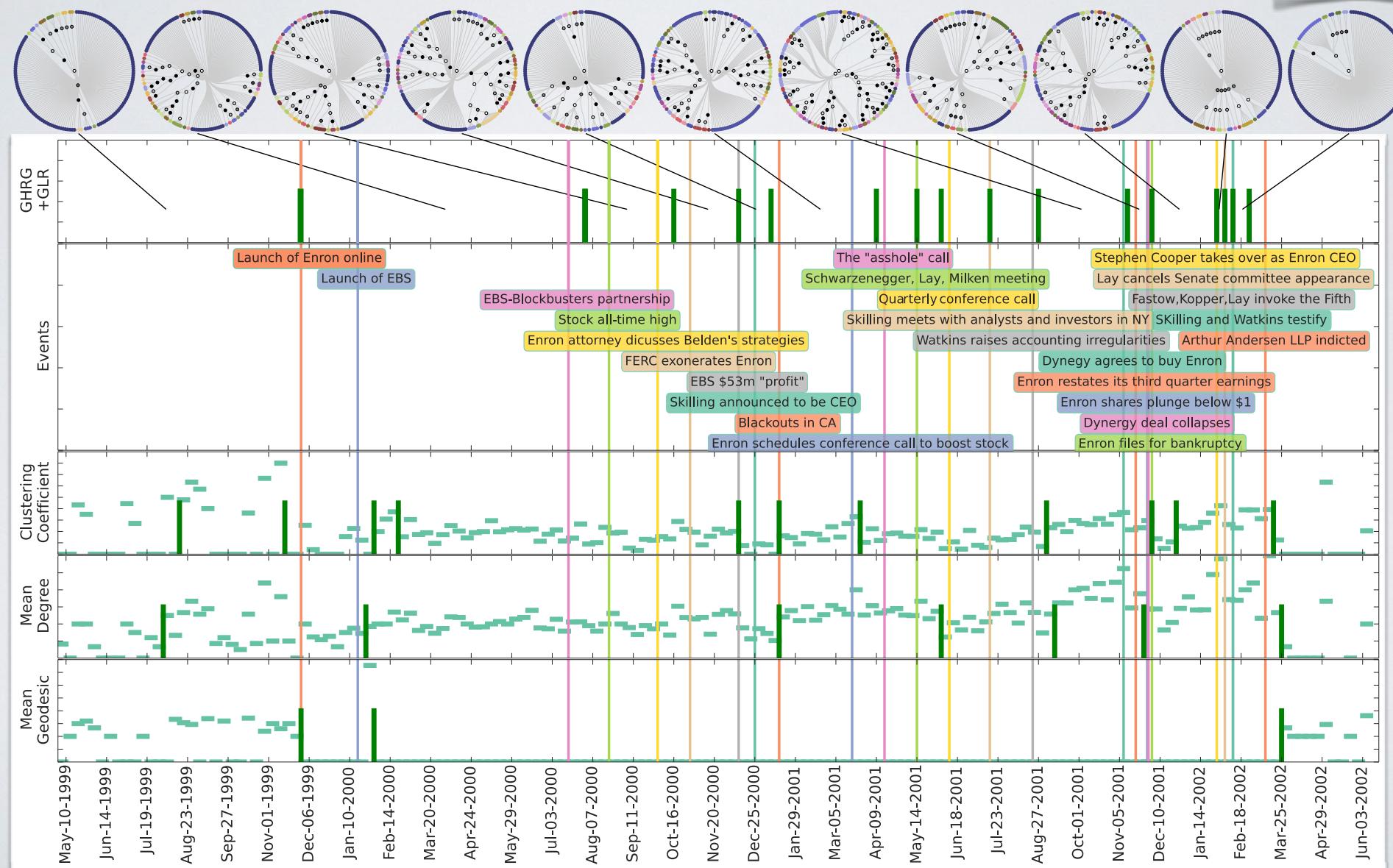
- use generalized likelihood ratio test for a change-point within our window

$$\frac{\underbrace{\mathcal{L}(G_{\leq t^*} \mid \theta_{\leq t^*})}_{\text{likelihood up to } t^*} \times \underbrace{\mathcal{L}(G_{>t^*} \mid \theta_{>t^*})}_{\text{likelihood after } t^*}}{\underbrace{\mathcal{L}(G_{\text{all}} \mid \theta_{\text{all}})}_{\text{likelihood of no change point}}}$$

detecting network change points



Enron email corpus + external events



generative models for complex networks

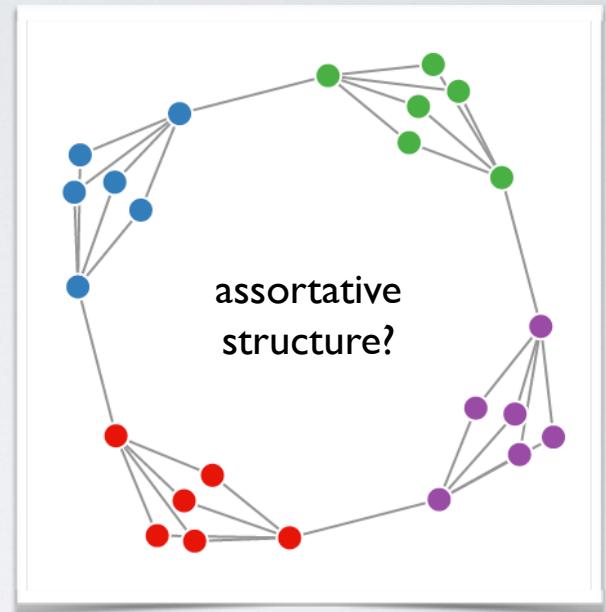
3 examples

- hierarchical random graphs (HRG)
- bipartite community structure (biSBM)
- weighted community structure (WSBM)

bipartite networks

many networks are bipartite

- scientists and papers (co-authorship networks)
- actors and movies (co-appearance networks)
- words and documents (topic modeling)
- plants and pollinators
- genes and genomes
- etc.

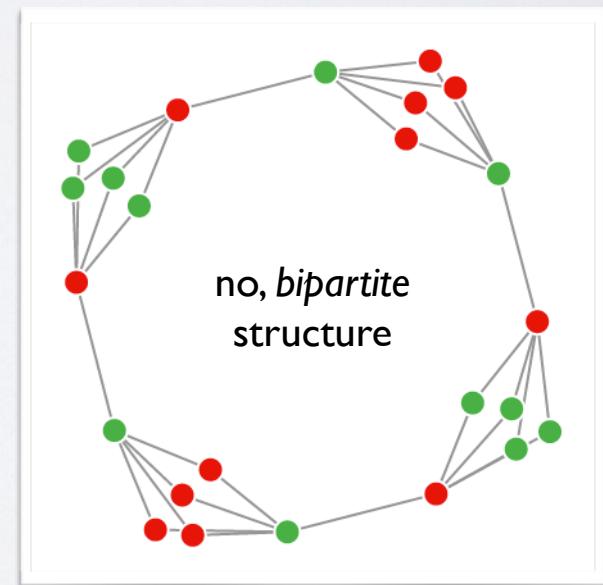
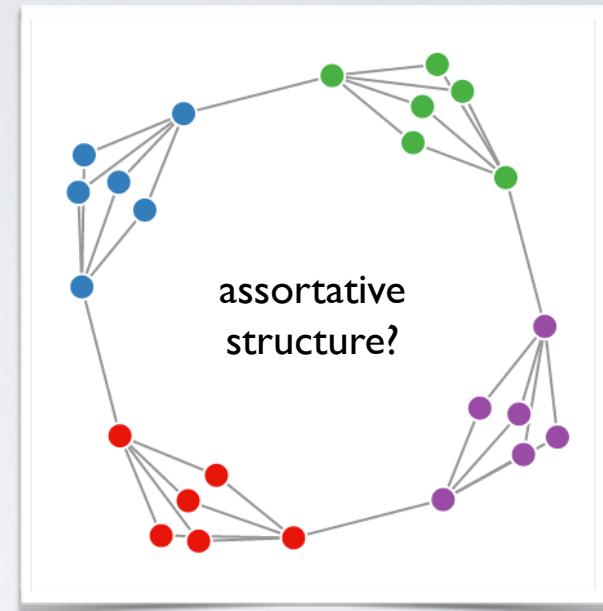


bipartite networks

many networks are bipartite

- scientists and papers (co-authorship networks)
- actors and movies (co-appearance networks)
- words and documents (topic modeling)
- plants and pollinators
- genes and genomes
- etc.

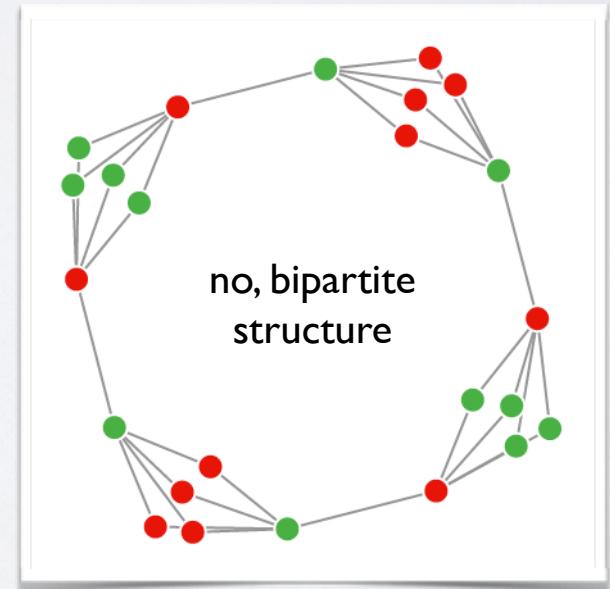
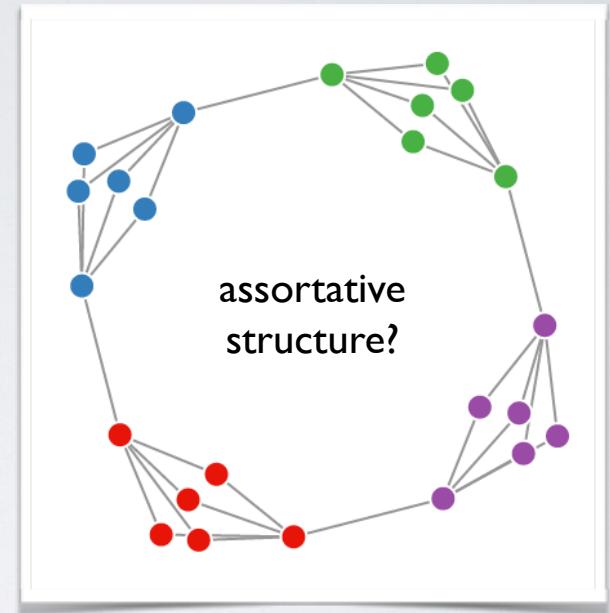
most analyses focus on one-mode projections
which discard information



bipartite networks

bipartite stochastic block model (biSBM)

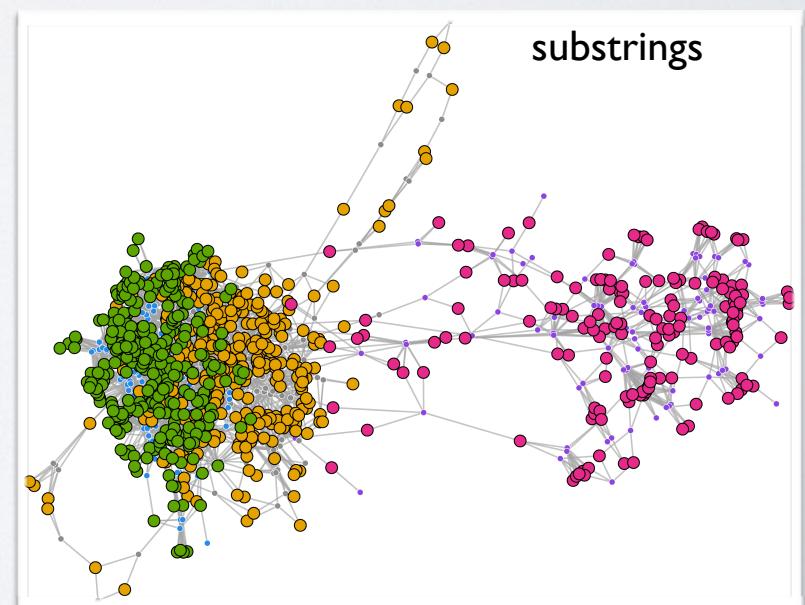
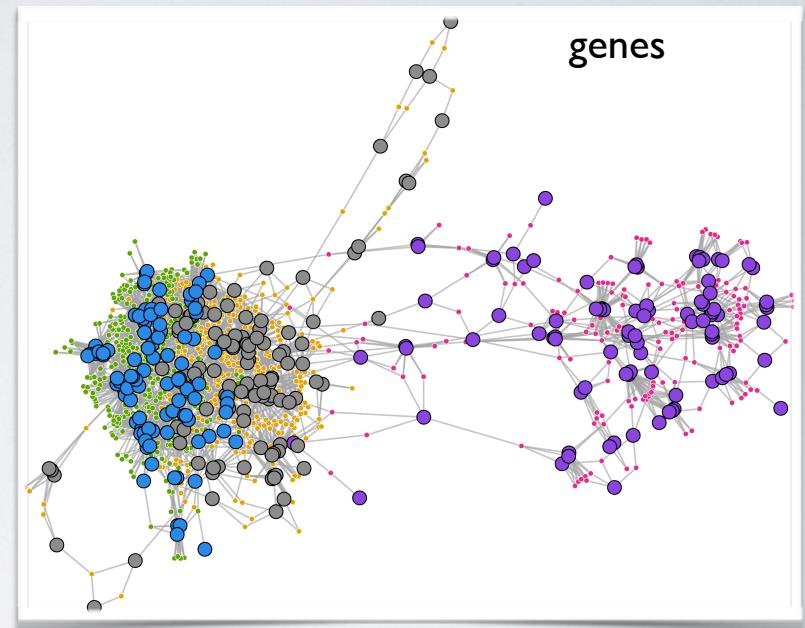
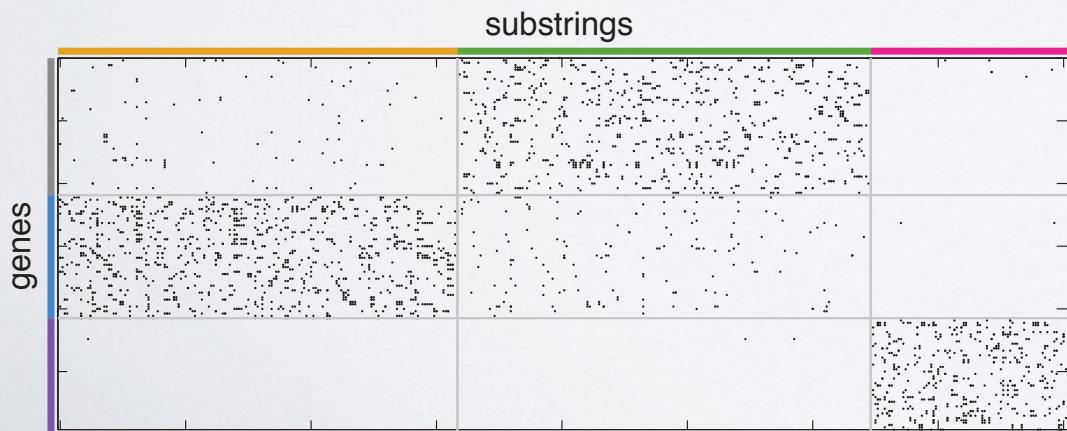
- exactly the SBM, but model knows network is bipartite
- if $\text{type}(z_i) = \text{type}(z_j)$
then require $M_{z_i, z_j} = 0$
- inference proceeds as before



bipartite networks

bipartite stochastic block model (biSBM)

- SBM can *learn* bipartite structure, but biSBM much more efficient, accurate
- biSBM always find pure-type communities
- more accurate than modeling one-mode projections (even weighted projections)
- finds communities in both modes

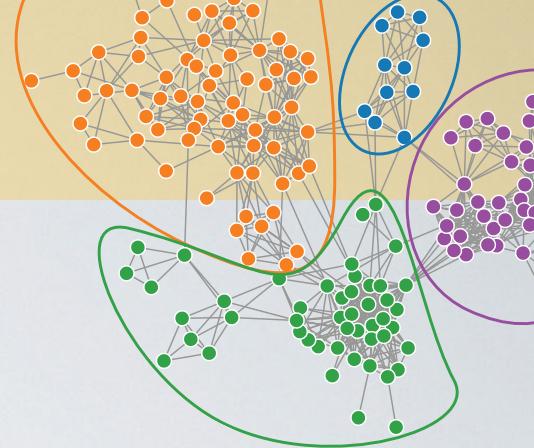


generative models for complex networks

3 examples

- hierarchical random graphs (HRG)
- bipartite community structure (biSBM)
- weighted community structure (WSBM)

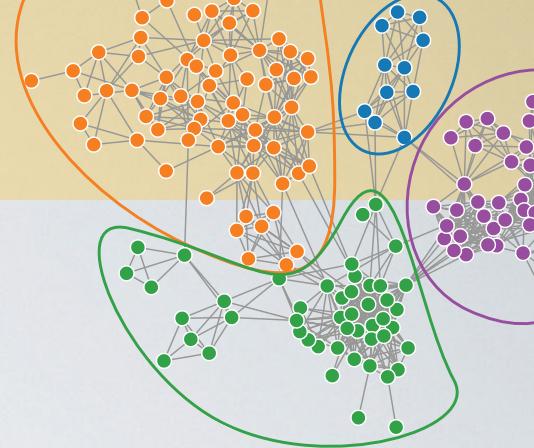
weighted networks



most interactions are weighted

- interaction frequency, strength, character, outcome, etc.
- thresholding discards information, can obscure underlying structure

weighted networks



most interactions are weighted

- interaction frequency, strength, character, outcome, etc.
- thresholding discards information, can obscure underlying structure

weighted SBM:

$$\ln \Pr(G | M, z, \theta, f) = \alpha \ln \Pr(G | M, z) + (1 - \alpha) \ln \Pr(G | \theta, z, f)$$

infer z, M, θ

edge-existence
[binomial distribution]

$$M_{z_i, z_j}$$

edge-weights
[exponential-family distribution]

$$\theta_{z_i, z_j}$$

Poisson, Normal, Gamma,
Exponential, Pareto, etc.

weighted networks

NFL 2009 season

- 32 teams, 2 “divisions”, 4 “subdivisions”
- *edge existence*: who plays whom
- *edge weight*: mean score difference

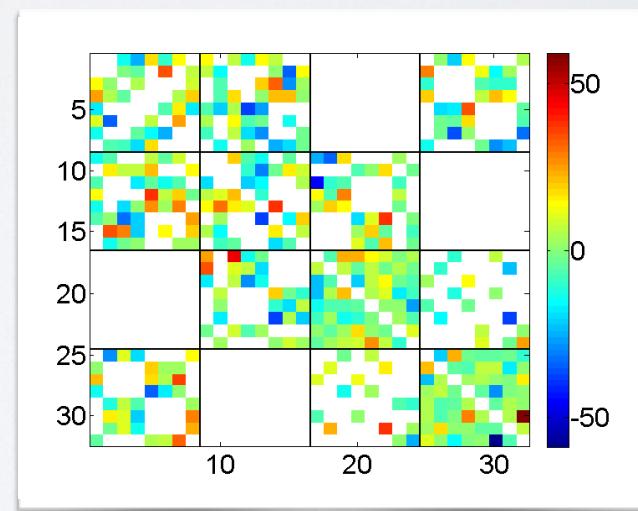
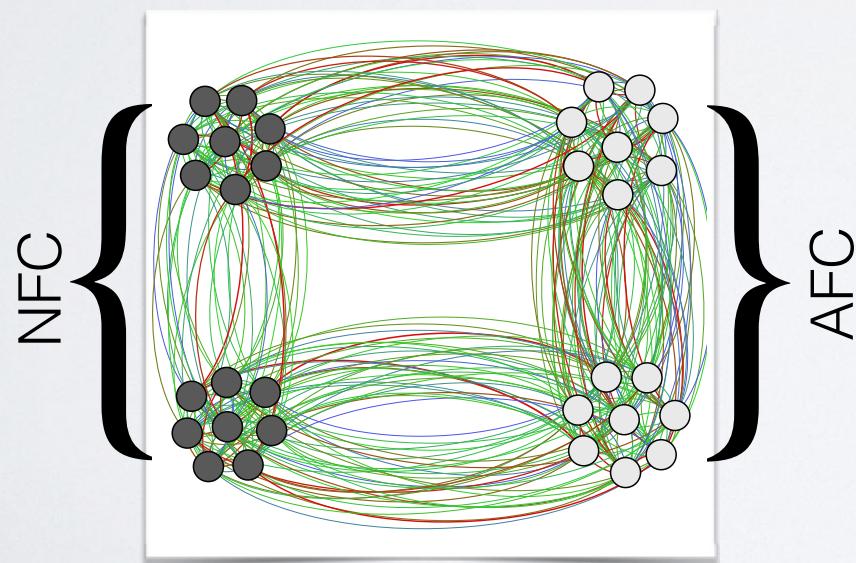


weighted networks



NFL 2009 season

- 32 teams, 2 “divisions”, 4 “subdivisions”
- SBM ($\alpha = 1$) recovers subdivisions perfectly



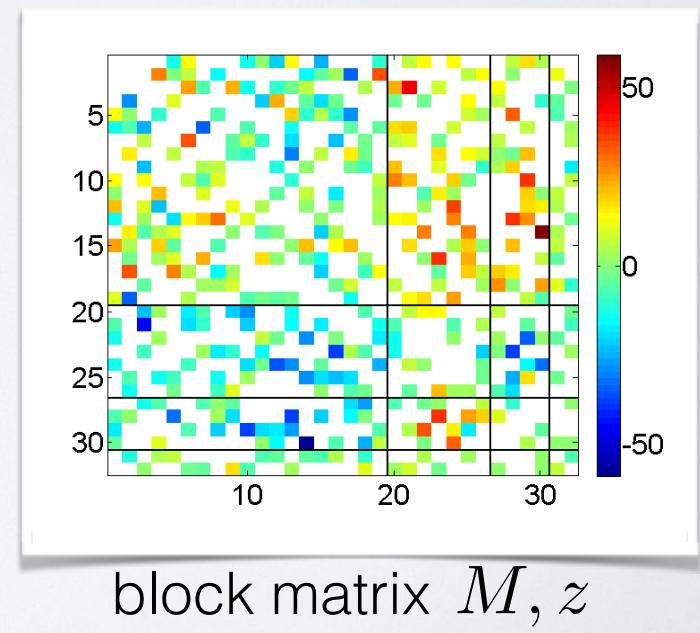
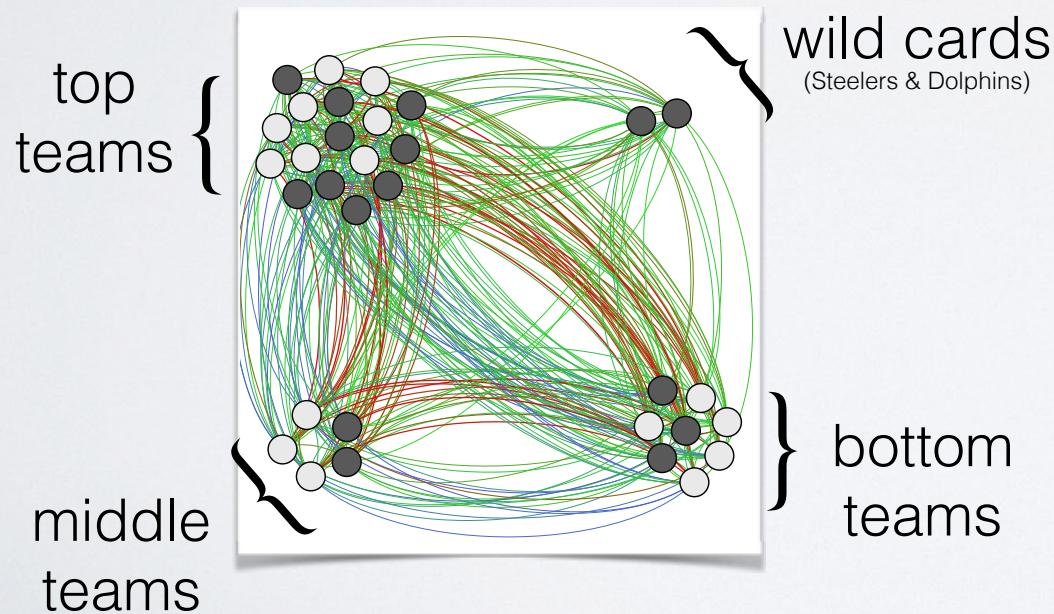
block matrix M, z

weighted networks



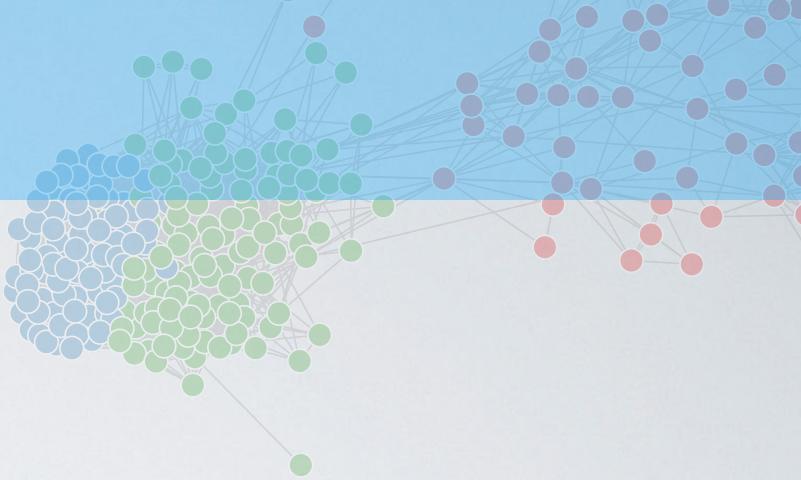
NFL 2009 season

- 32 teams, 2 “divisions”, 4 “subdivisions”
- WSBM ($\alpha=0$) recovers team skill hierarchy

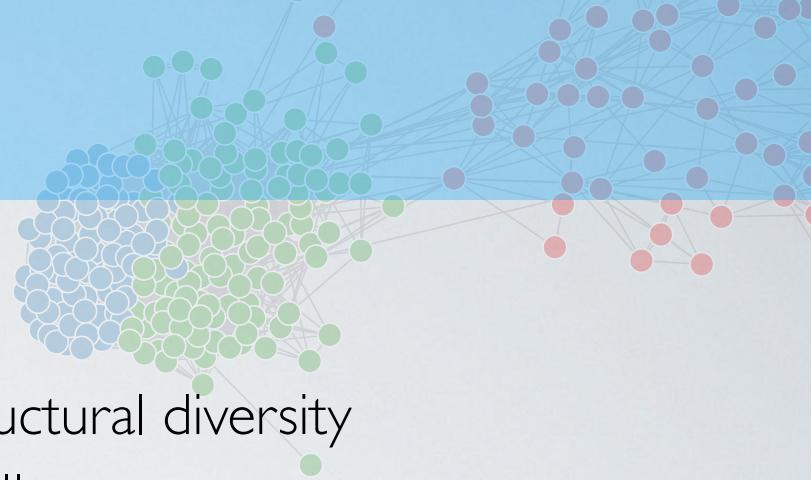


parting thoughts on networks

- networks are cool!



parting thoughts on networks



- **networks are cool!**
but also complicated objects = enormous structural diversity
still figuring out how to describe structure well
- **we have only scratched the surface**
auxiliary data (weights, attributes, time)
multiplex networks
stronger tools for testing hypotheses
applications abound [new ideas often come from these]
- **structure + dynamics = function**
how does structure constrain dynamics, robustness, etc.
to what degree does structure = function?
- **formation mechanisms**
where does structure come from?

thanks



Chris Aicher
(Colorado)



Abigail Z. Jacobs
(Colorado)



Dr. Dan Larremore
(Harvard)



Dr. Leto Peel
(Colorado)



Prof. Cris Moore
(Santa Fe)



Prof. Mark Newman
(Michigan)



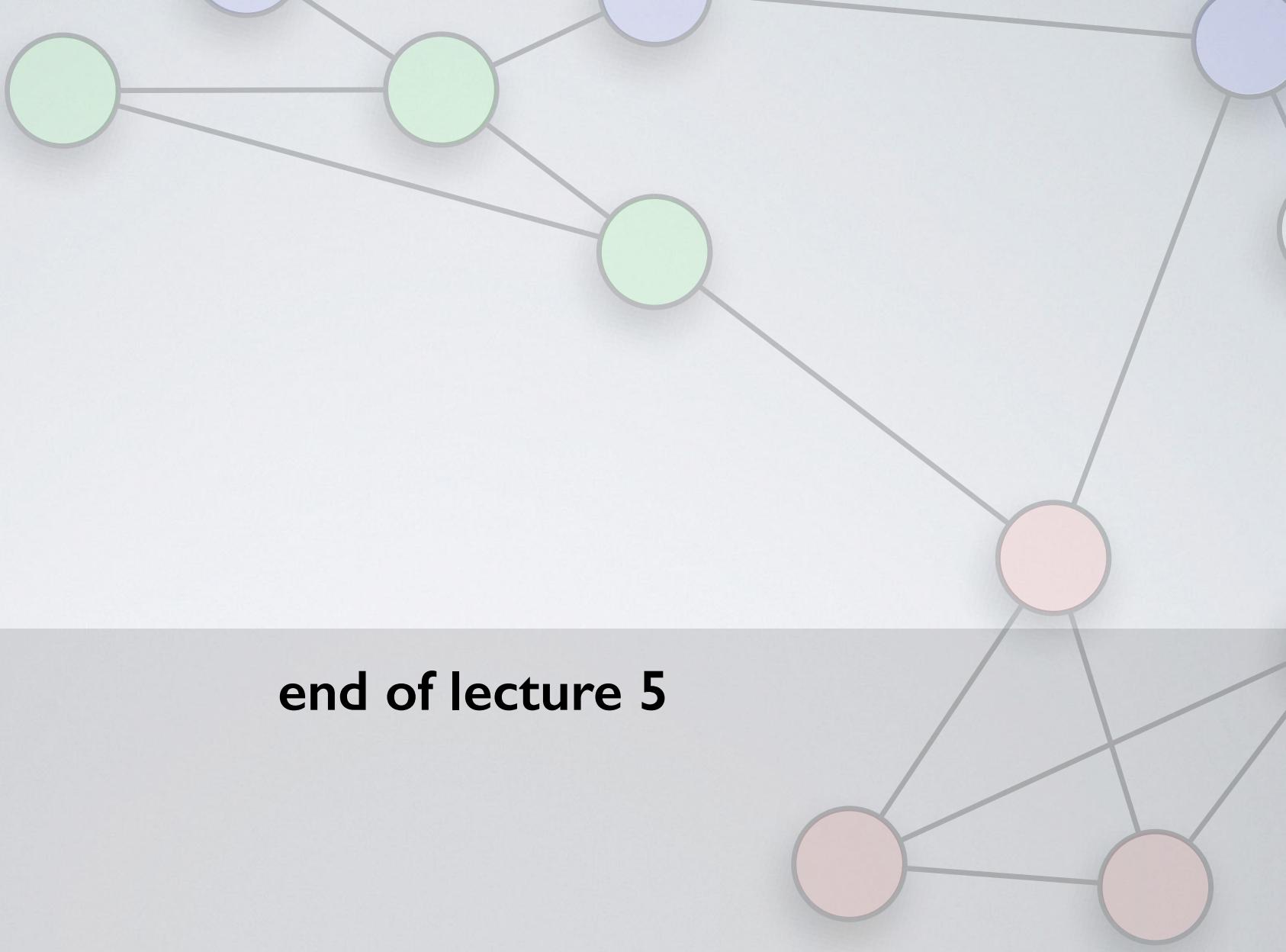
Prof. Caroline Buckee
(Harvard)



James S. McDonnell Foundation



CENTER for
COMMUNICABLE
DISEASE DYNAMICS



A network graph is displayed against a light gray background. The graph consists of several circular nodes connected by thin gray lines. There are three distinct clusters of nodes: a top-left cluster of three green nodes, a bottom-right cluster of three pink nodes, and a large, sparse cluster of numerous small, semi-transparent gray nodes. The text 'end of lecture 5' is centered over the green cluster.

end of lecture 5

more on inference

code + data available at

hierarchical SBM santafe.edu/~aaronc/hierarchy/

weighted SBM santafe.edu/~aaronc/wsrbm/

bipartite SBM danlarremore.com/bipartiteSBM/

change-point detection SBM gdrive.es/letopeel/code.html

further reading

- Larremore, Clauset and Jacobs, "Efficiently inferring community structure in bipartite networks." Preprint (2014) [arxiv:1403.2933]
- Peel and Clauset, "Detecting change points in the large-scale structure of evolving networks." Preprint (2014) [arxiv:1403.0989]
- Aicher, Jacobs and Clauset, "Learning latent block structure in weighted networks." To appear, *Journal of Complex Networks* (2014) [arxiv:1404.0431]
- Larremore, Clauset and Buckee, "A network approach to analyzing highly recombinant malaria parasite genes." *PLOS Computational Biology* 9, e1003268 (2013) [arxiv:1308.5254]
- Aicher, Jacobs and Clauset, "Adapting the stochastic block model to edge-weighted networks." *ICML Ws* (2013) [arxiv:1305.5782]
- Clauset, Moore, and Newman, "Hierarchical structure and the prediction of missing links in networks" *Nature* 453, 98-101 (2008) [arxiv:0811.0484]

selected references

- The structure and function of complex networks. M. E. J. Newman, *SIAM Review* **45**, 167–256 (2003).
- *The Structure and Dynamics of Networks*. M. E. J. Newman, A.-L. Barabási, and D. J. Watts, Princeton University Press (2006).
- Hierarchical structure and the prediction of missing links in networks. A. Clauset, C. Moore, and M. E. J. Newman, *Nature* **453**, 98–101 (2008).
- Modularity and community structure in networks. M. E. J. Newman, *Proc. Natl. Acad. Sci. USA* **103**, 8577–8582 (2006).
- Why social networks are different from other types of networks. M. E. J. Newman and J. Park, *Phys. Rev. E* **68**, 036122 (2003)
- Random graphs with arbitrary degree distributions and their applications. M. E. J. Newman, S. H. Strogatz, and D. J. Watts, *Phys. Rev. E* **64**, 026118 (2001).
- Comparing community structure identification. L. Danon, A. Diaz-Guilera, J. Duch and A. Arenas. *J. Stat. Mech.* P09008 (2005).
- Characterization of Complex Networks: A Survey of measurements. L. daF. Costa, F. A. Rodrigues, G. Travieso and P. R. VillasBoas. arxiv:cond-mat/050585 (2005).
- Evolution in Networks. S.N. Dorogovtsev and J. F. F. Mendes. *Adv. Phys.* **51**, 1079 (2002).
- Revisiting “scale-free” networks. E. F. Keller. *BioEssays* **27**, 1060-1068 (2005).
- Currency metabolites and network representations of metabolism. P. Holme and M. Huss. arxiv:0806.2763 (2008).
- Functional cartography of complex metabolic networks. R. Guimera and L. A. N. Amaral. *Nature* **433**, 895 (2005).
- Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations. J. Leskovec, J. Kleinberg and C. Faloutsos. *Proc. 11th ACM SIGKDD Intl. Conf. on Knowledge Discovery and Data Mining* 2005.
- The Structure of the Web. J. Kleinberg and S. Lawrence. *Science* **294**, 1849 (2001).
- Navigation in a Small World. J. Kleinberg. *Nature* **406** (2000), 845.
- Towards a Theory of Scale-Free Graphs: Definitions, Properties and Implications. L. Li, D. Alderson, J. Doyle, and W. Willinger. *Internet Mathematics* **2**(4), 2006.
- A First-Principles Approach to Understanding the Internet’s Router-Level Topology. L. Li, D. Alderson, W. Willinger, and J. Doyle. *ACM SIGCOMM* 2004.
- Inferring network mechanisms: The *Drosophila melanogaster* protein interaction network. M. Middendorf, E. Ziv and C. H. Wiggins. *Proc. Natl. Acad. Sci. USA* **102**, 3192 (2005).
- Robustness Can Evolve Gradually in Complex Regulatory Gene Networks with Varying Topology. S. Ciliberti, O. C. Martin and A. Wagner. *PLoS Comp. Bio.* **3**, e15 (2007).
- Simple rules yield complex food webs. R. J. Williams and N. D. Martinez. *Nature* **404**, 180 (2000).
- A network analysis of committees in the U.S. House of Representatives. M. A. Porter, P. J. Mucha, M. E. J. Newman and C. M. Warmbrand. *Proc. Natl. Acad. Sci. USA* **102**, 7057 (2005).
- On the Robustness of Centrality Measures under Conditions of Imperfect Data. S. P. Borgatti, K. M. Carley and D. Krackhardt. *Social Networks* **28**, 124 (2006).