Distinguishing scale-free topology generators using temporal data Naomi Arnold, Raul Mondragón, Richard Clegg



CCS2018 THESSALONIKI GREECE 25 September 2018



QUEELI I VICI Y University of London



Motivation

- The structure of many networks of interest is often dynamic in nature
- Rich temporal data for network topologies is becoming more prevalent
- This allows us to better investigate models for network growth over time

























Action (new node/internal link)



Action (new node/internal link)



Action (new node/internal link)



Seed network

Attachment probabilities

 ν_i

corresponding to $\mathbb{P}($ **choose node** i)



$p_i \propto 1$ Random/neutral model. All nodes equally likely





$p_i \propto 1$ **Random/neutral** model. All nodes equally likely



 $p_i \propto f(k_i)$

Function of node i's degree, e.g. BA model

$p_i \propto 1$ **Random/neutral** model. All nodes equally likely



 $p_i \propto f(k_i)$

Function of node i's degree, e.g. BA model

 $p_i \propto f(\eta_i)$ **Function of some** other intrinsic node property





How do we quantify how good a fit a model is to real data?

How do we quantify how good a fit a model is to real data?

 Traditional approach: generate a network from model of the same size as target network, and compare on different statistics

How do we quantify how good a fit a model is to real data?

- different statistics
- Problem: generating large networks is time consuming

 Traditional approach: generate a network from model of the same size as target network, and compare on

How do we quantify how good a fit a model is to real data?

- different statistics
- Problem: generating large networks is time consuming
- **Problem:** two networks having similar properties in the same way

 Traditional approach: generate a network from model of the same size as target network, and compare on

doesn't necessarily mean their evolution was governed

How do we quantify how good a fit a model is to real data?

- different statistics
- Problem: generating large networks is time consuming
- **Problem:** two networks having similar properties Why? in the same way

 Traditional approach: generate a network from model of the same size as target network, and compare on

doesn't necessarily mean their evolution was governed

Two different models for scale-free networks

[A.-L. Barabasi, R. Albert 1999: Emergence of scaling in random networks]

Two different models for scale-free networks

Barabási Albert model

 $p_i \propto k_i$

Higher degree nodes likelier to attract new links

[A.-L. Barabasi, R. Albert 1999: Emergence of scaling in random networks]

Two different models for scale-free networks

Barabási Albert model

$p_i \propto k_i$

Higher degree nodes likelier to attract new links

Mean field prediction: $P(k) \sim k^{-3}$

[A.-L. Barabasi, R. Albert 1999: Emergence of scaling in random networks]

Two different models for scale-free networks **Rank-preference model** $p_i \propto i^{-\alpha}$

Barabási Albert model

$p_i \propto k_i$

Higher degree nodes likelier to attract new links

Mean field prediction: $P(k) \sim k^{-3}$

[A.-L. Barabasi, R. Albert 1999: Emergence of scaling in random networks]

Older nodes likelier to attract new links

Two different models for scale-free networks **Rank-preference model**

Barabási Albert model

$p_i \propto k_i$

Higher degree nodes likelier to attract new links

Mean field prediction: $P(k) \sim k^{-3}$

[A.-L. Barabasi, R. Albert 1999: Emergence of scaling in random networks]

$p_i \propto i^{-\alpha}$ Older nodes likelier to attract new links

Mean field prediction: $P(k) \sim k^{-(1+1/\alpha)}$



Two different models for scale-free networks **Rank-preference model**

Barabási Albert model

$p_i \propto k_i$

Higher degree nodes likelier to attract new links

Mean field prediction: $P(k) \sim k^{-3}$

[A.-L. Barabasi, R. Albert 1999: Emergence of scaling in random networks]

$p_i \propto i^{-\alpha}$ Older nodes likelier to attract new links

Mean field prediction: $P(k) \sim l^{-(1+1/\alpha)}$



Degree distribution



Network of 100,000 nodes generated from the Barabási-Albert and **Rank Preference** models, each new node having initially 3 neighbours

















[R. Clegg, B. Parker, M. Rio 2016: Likelihood-based assessment of dynamic network models]

Likelihood of model given observation = probability of seeing observation given model













[R. Clegg, B. Parker, M. Rio 2016: Likelihood-based assessment of dynamic network models]

Likelihood of model given observation = probability of seeing observation given model













[R. Clegg, B. Parker, M. Rio 2016: Likelihood-based assessment of dynamic network models]

Likelihood of model given observation = probability of seeing observation given model

> Likelihood of random/uniform model given by:

 $\mathbb{P}_{rand}(\textbf{Choose node } 2) = \frac{1}{3}$



















[R. Clegg, B. Parker, M. Rio 2016: Likelihood-based assessment of dynamic network models]

Likelihood of model given observation = probability of seeing observation given model

> Likelihood of random/uniform model given by:

 \mathbb{P}_{rand} (**Choose node** 2) = $\frac{1}{3}$

Likelihood of BA preferential attachment model given by:

Quickly calculated, compared to generating networks

- Quickly calculated, compared to generating networks
- Given a number of models, can define the 'best' as that which has the highest likelihood

- Quickly calculated, compared to generating networks
- Given a number of models, can define the 'best' as that which has the highest likelihood
- For models with parameters, can find maximum likelihood estimators for params

Barabási-Albert

Model Mixtures

- Hypothesis: network growth likely to be governed by a mixture of mechanisms, not just one
 - Example: Mixture of BA and rank preference model

 $p_i = (1 - \beta)p_i^{BA} + \beta p_i^{RP}$

Rank-Preference

Barabási-Albert

Model Mixtures

- Hypothesis: network growth likely to be governed by a mixture of mechanisms, not just one
 - Example: Mixture of BA and rank preference model

 $p_i = (1 - \beta)p_i^{BA} + \beta p_i^{RP}$

Rank-Preference

Barabási-Albert

Model Mixtures

- Hypothesis: network growth likely to be governed by a mixture of mechanisms, not just one
 - Example: Mixture of BA and rank preference model

 $p_i = (1 - \beta)p_i^{BA} + \beta p_i^{RP}$

Rank-Preference

Distinguishing using maximum likelihood estimation

Having generated artificial networks using:

$$p_i = \beta p_i^{RP} + (1 - \beta) p_i^{BA}$$

We can accurately recover the proportion β as an MLE!

Distinguishing using maximum likelihood estimation

Having generated artificial networks using:

$$p_i = \beta p_i^{RP} + (1 - \beta) p_i^{BA}$$

We can accurately recover the proportion β as an MLE!

- Q&A site for mathematical problems
- Nodes are users
- An undirected edge between node A and B if A answers a question by B, A comments on B's answer or question
- Multiple edges collapsed
- Models components tested: BA, static rank preference

Real data example: Routeviews AS topology dataset

- links represent peering relationship
- Model components tested: random/uniform model, positive-feedback preference model

Note - positive-feedback preference model gives:

PFP model: [S. Zhou, R. Mondragón 2004: Accurately modelling the Internet Topology]

Timestamped dataset: nodes autonomous systems and

$$p_i \propto k_i^{1+\delta \log_{10} k_i}$$

Real data example: Routeviews AS topology dataset

Real data example: Routeviews AS topology dataset

Remarks

- Often (**not always**) the highest likelihood mixture of model components may generate networks with better matching statistics than any single component alone
- But this highest likelihood mix is not guaranteed to be a good model still need to look at network statistics
- Finding maximum likelihood mix of more than two model components may become expensive candidate for ML techniques

Conclusions and future directions

- Temporal data allows calculation of likelihoods of dynamic models given observed network evolution
- similar in structure
- network growth
- change over time

• With this measure we can distinguish models which generate networks that are

Fitting mixed models helps uncover the roles of different processes governing

Work in progress: using likelihood measure to analyse how such processes may

Thank you for listening! What are your questions?

github.com/narnolddd

n.a.arnold@qmul.ac.uk

