# "Root and community inference on the latent growth process of a network"

Authors: Harry Crane, Min Xu

Nilava Metya

nilavam.github.io
nilava.metya@rutgers.edu

Department of Mathematics
Rutgers University

# The real world problems

# The real world problems

- An infection starting from a source, and has spread out.

# The real world problems

- An infection starting from a source, and has spread out.

- An incorrect rumor has been generated by a person and has spread through social networks.

# The real world problems

- An infection starting from a source, and has spread out.

- An incorrect rumor has been generated by a person and has spread through social networks.

We will only observe the structure of spreading after the spreading has been done.

# The real world problems

- An infection starting from a source, and has spread out.

- An incorrect rumor has been generated by a person and has spread through social networks.

We will only observe the structure of spreading after the spreading has been done.
Want to find the source.

# Notation

# Notation

- All graphs are undirected. Standard: $g = (V, E), V = V(g), E = E(g)$.

# Notation

- All graphs are undirected. Standard: $g = (V, E), V = V(g), E = E(g)$.

- Capital letters <-> Random objects
  Lowercase letters <-> Fixed objects

$$\text{APA}(\alpha, \beta)$$

# APA($\alpha, \beta$)

- The affine <u>preferential attachment</u> tree model with parameters $\alpha$, $\beta$ generates an increasing sequence $T_1 \subset T_2 \subset \cdots \subset T_n$ of random trees where $T_i$ is a labelled tree with $i$ nodes and nodes are labelled by their arrival time so that $V(T_i) = [i]$.

# APA($\alpha, \beta$)

- The affine <u>preferential attachment</u> tree model with parameters $\alpha$, $\beta$ generates an increasing sequence $T_1 \subset T_2 \subset \cdots \subset T_n$ of random trees where $T_i$ is a labelled tree with $i$ nodes and nodes are labelled by their arrival time so that $V(T_i) = [i]$.

- The generation looks something like this:

# APA($\alpha, \beta$)

- The affine <u>preferential attachment</u> tree model with parameters $\alpha$, $\beta$ generates an increasing sequence $T_1 \subset T_2 \subset \cdots \subset T_n$ of random trees where $T_i$ is a labelled tree with $i$ nodes and nodes are labelled by their arrival time so that $V(T_i) = [i]$.

- The generation looks something like this:

  > $T_1 = ([1], \{\})$

# APA($\alpha, \beta$)

- The affine <u>preferential attachment</u> tree model with parameters $\alpha, \beta$ generates an increasing sequence $T_1 \subset T_2 \subset \cdots \subset T_n$ of random trees where $T_i$ is a labelled tree with $i$ nodes and nodes are labelled by their arrival time so that $V(T_i) = [i]$.

- The generation looks something like this:

  > $T_1 = ([1], \{\})$

  > Given $T_{t-1}$, add a node labelled $t$ and a random edge $(t, w_t)$ to get $T_t$ where $w_t$ is chosen with probability $\dfrac{\beta \cdot D_{T_{t-1}}(w_t) + \alpha}{2\beta(t-2) + \alpha(t-1)}$.

# Examples for APA($\alpha, \beta$)

# Examples for APA($\alpha, \beta$)

- APA$(1,0)$ gives the probability $\dfrac{1}{t-1}$. So a neighbor is chosen uniformly from $V(T_{t-1})$.

- APA$(0,1)$ gives the probability $\dfrac{D_{T_{t-1}}(w_t)}{2(t-2)}$. So a neighbor is chosen with probability proportional to its degree.

# PAPER$(\alpha, \beta, \theta)$

PAPER = Preferential Attachment Plus Erdös-Rényi

# PAPER$(\alpha, \beta, \theta)$

PAPER = Preferential Attachment Plus Erdös-Rényi

We say that a random graph $G_n$ is distributed accordion to PAPER$(\alpha, \beta, \theta)$ if $G_n = T_n + R_n$ if $T_n \sim \mathrm{APA}(\alpha, \beta)$ and $R_n \sim \mathrm{Erdös-Rényi}(\theta)$.

# PAPER($\alpha, \beta, \theta$)

PAPER = Preferential Attachment Plus Erdös-Rényi

We say that a random graph $G_n$ is distributed accordion to PAPER($\alpha, \beta, \theta$) if $G_n = T_n + R_n$ if $T_n \sim \text{APA}(\alpha, \beta)$ and $R_n \sim \text{Erdös} - \text{Rényi}(\theta)$.

Will drop subscript
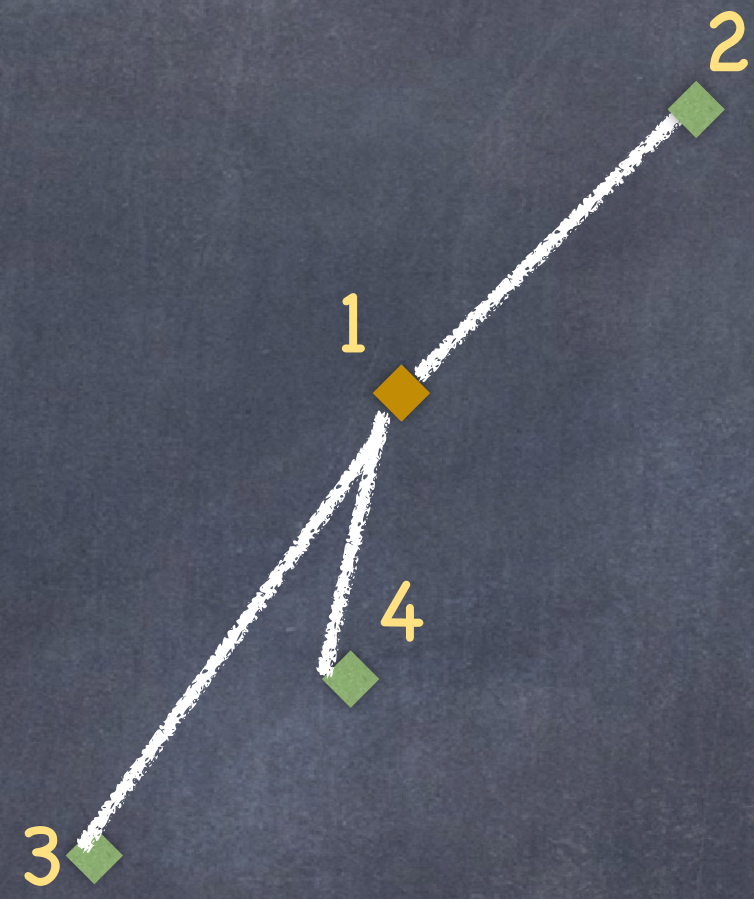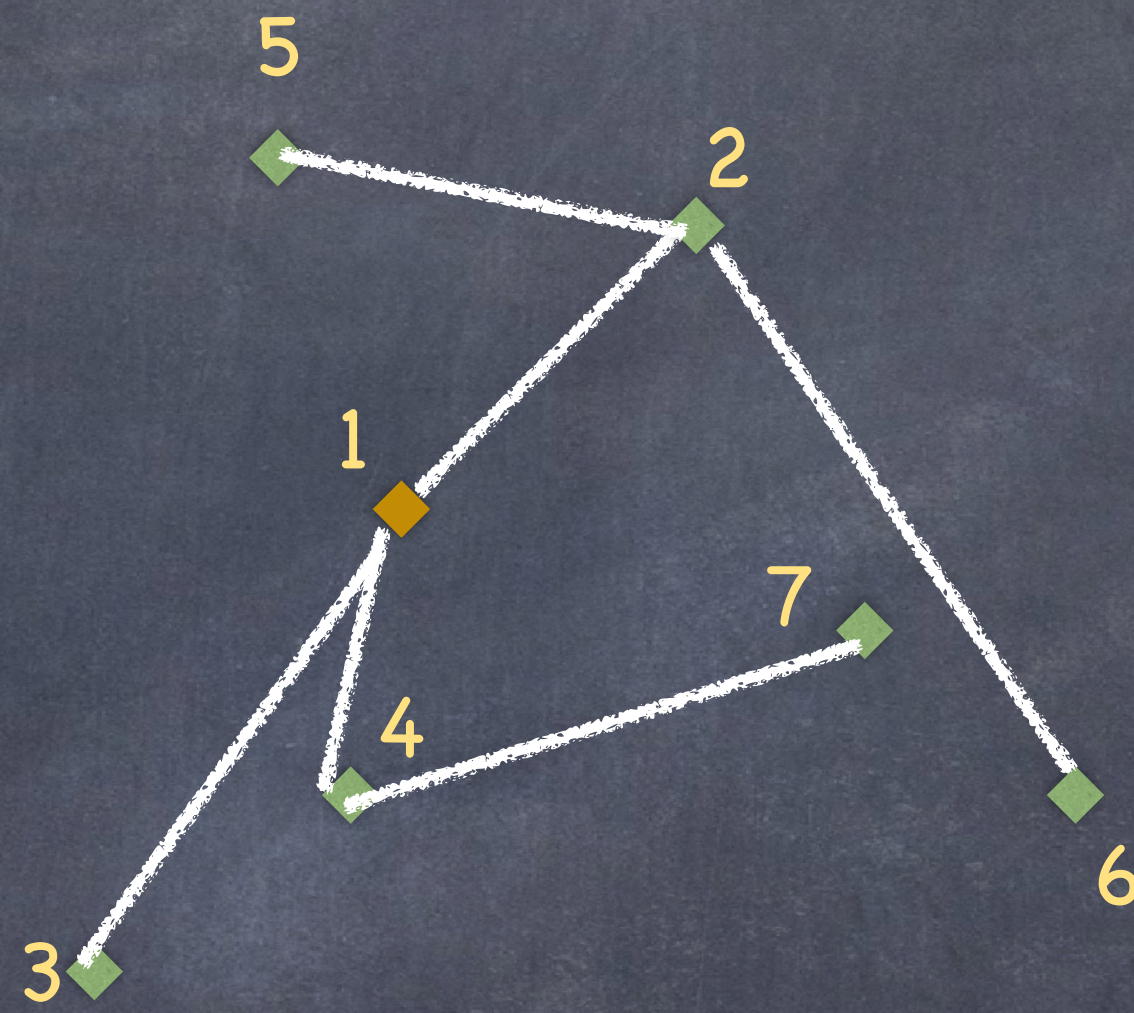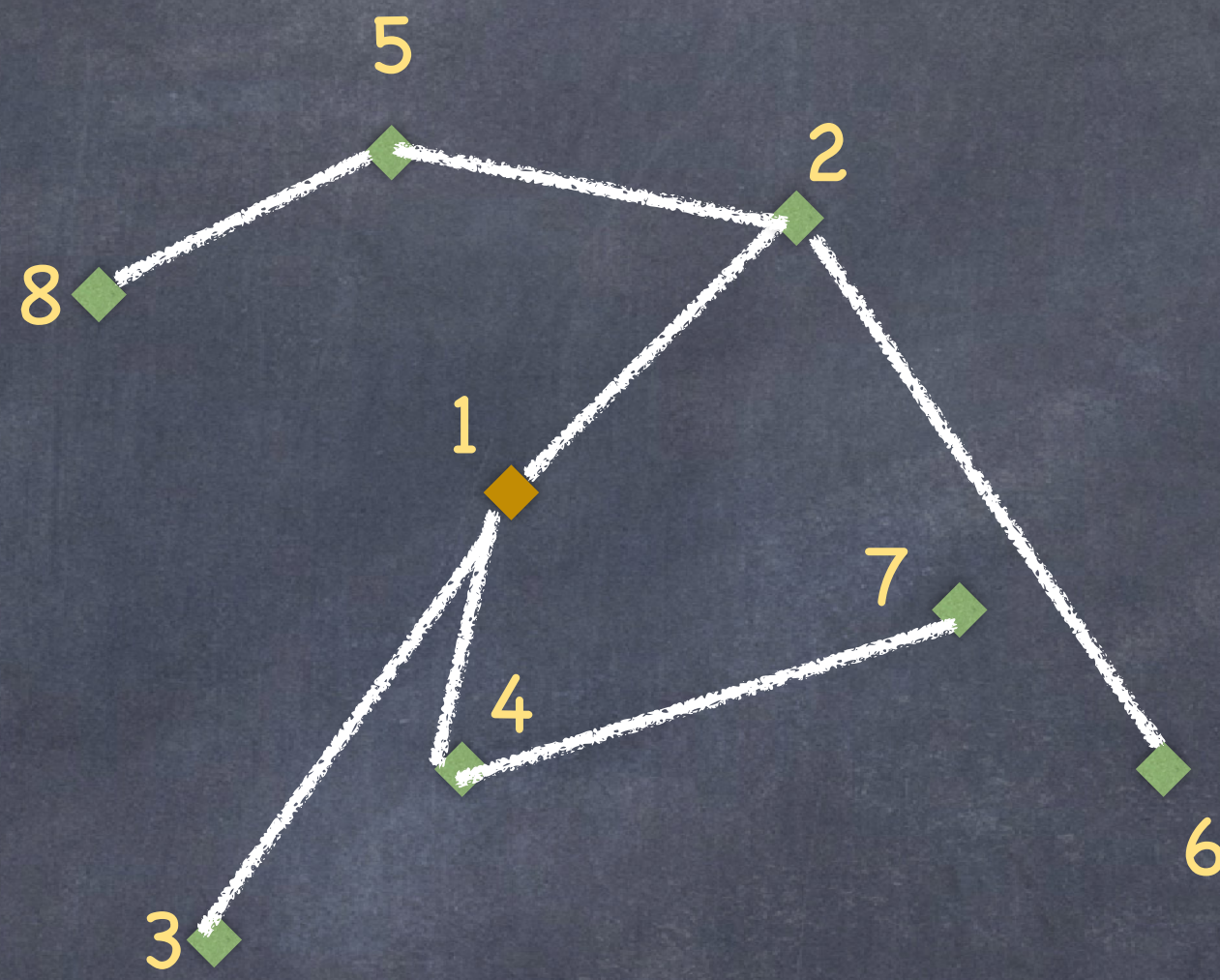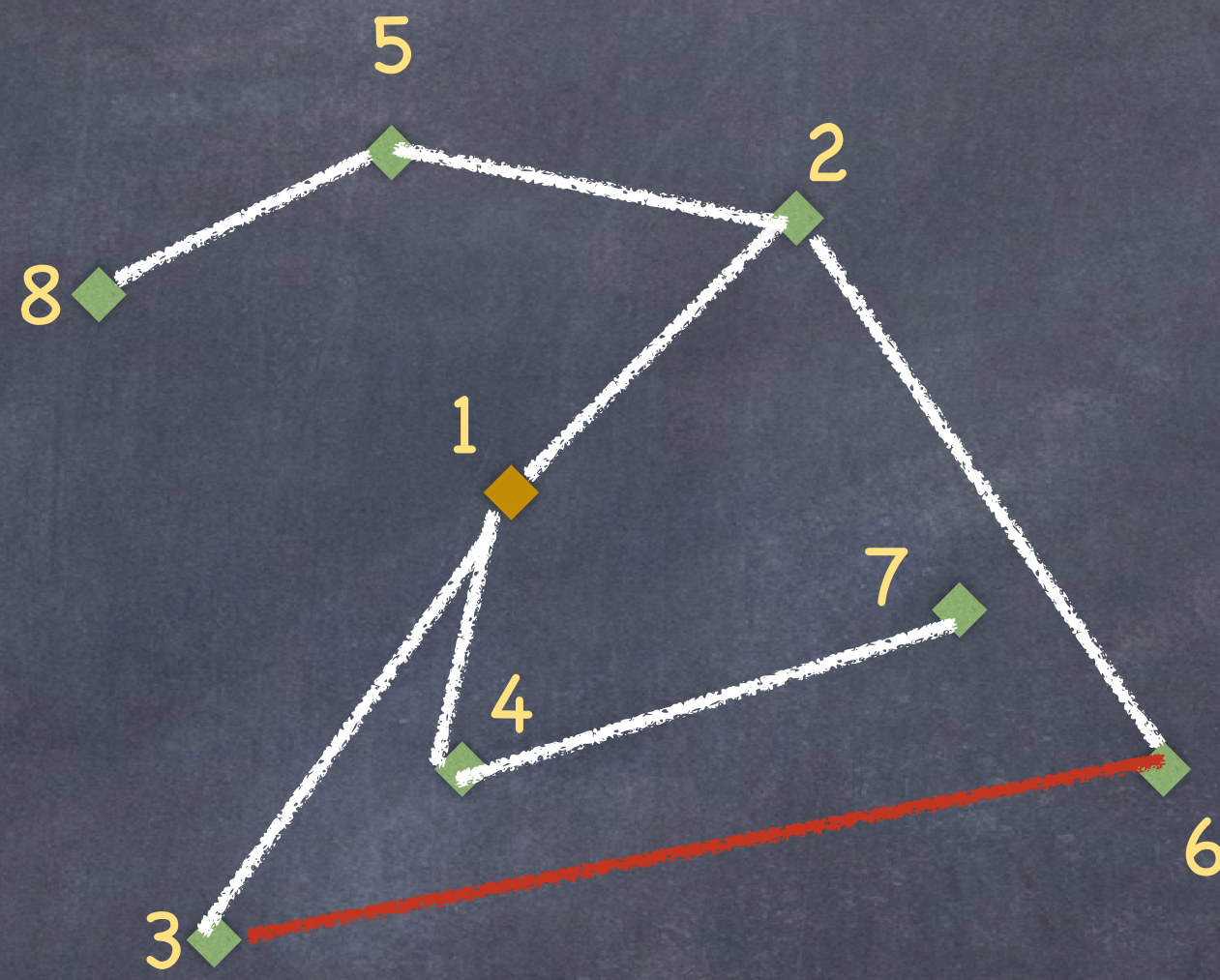
# Actual network

# Actual network

1
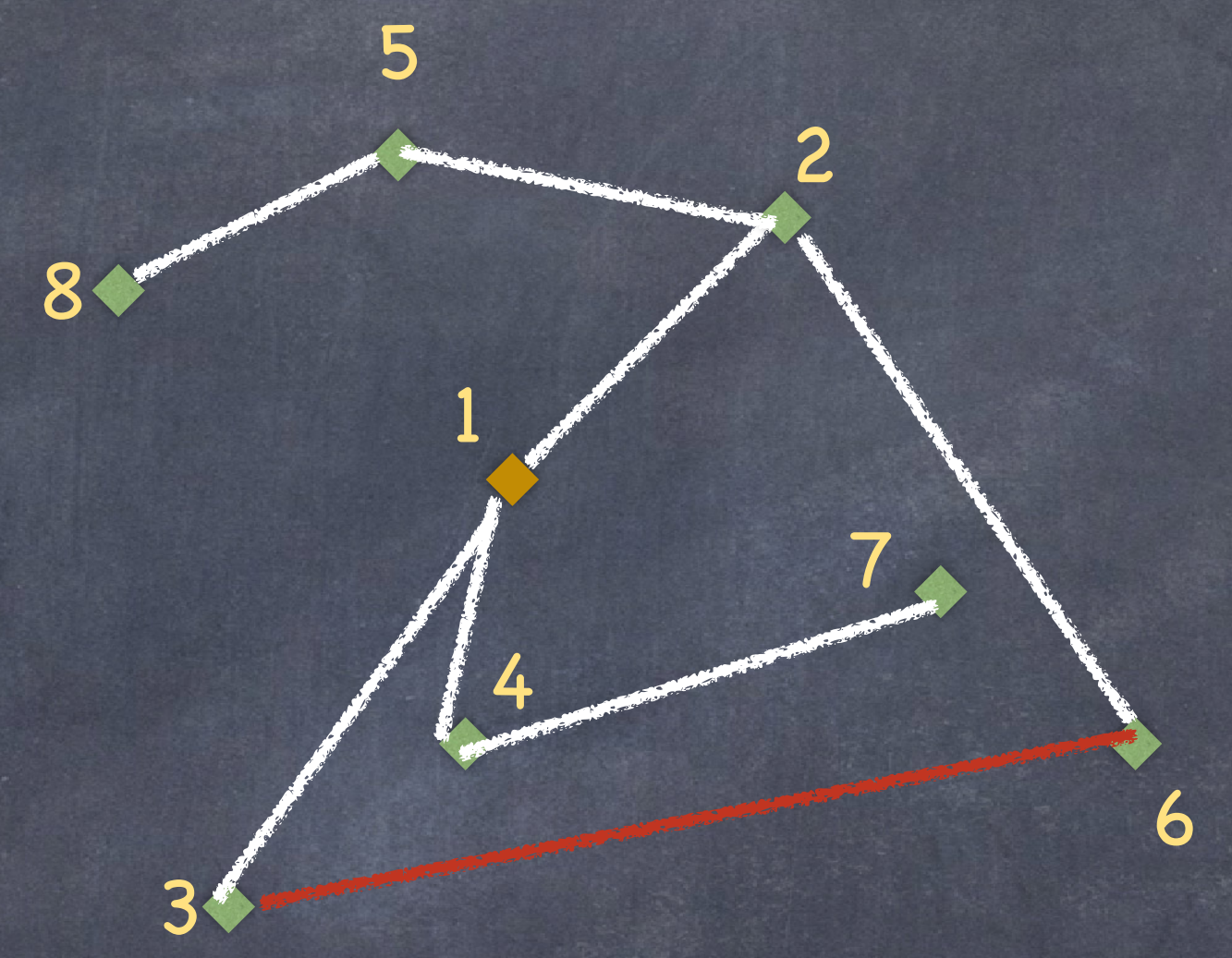
# Actual network

# Actual network

# Actual network

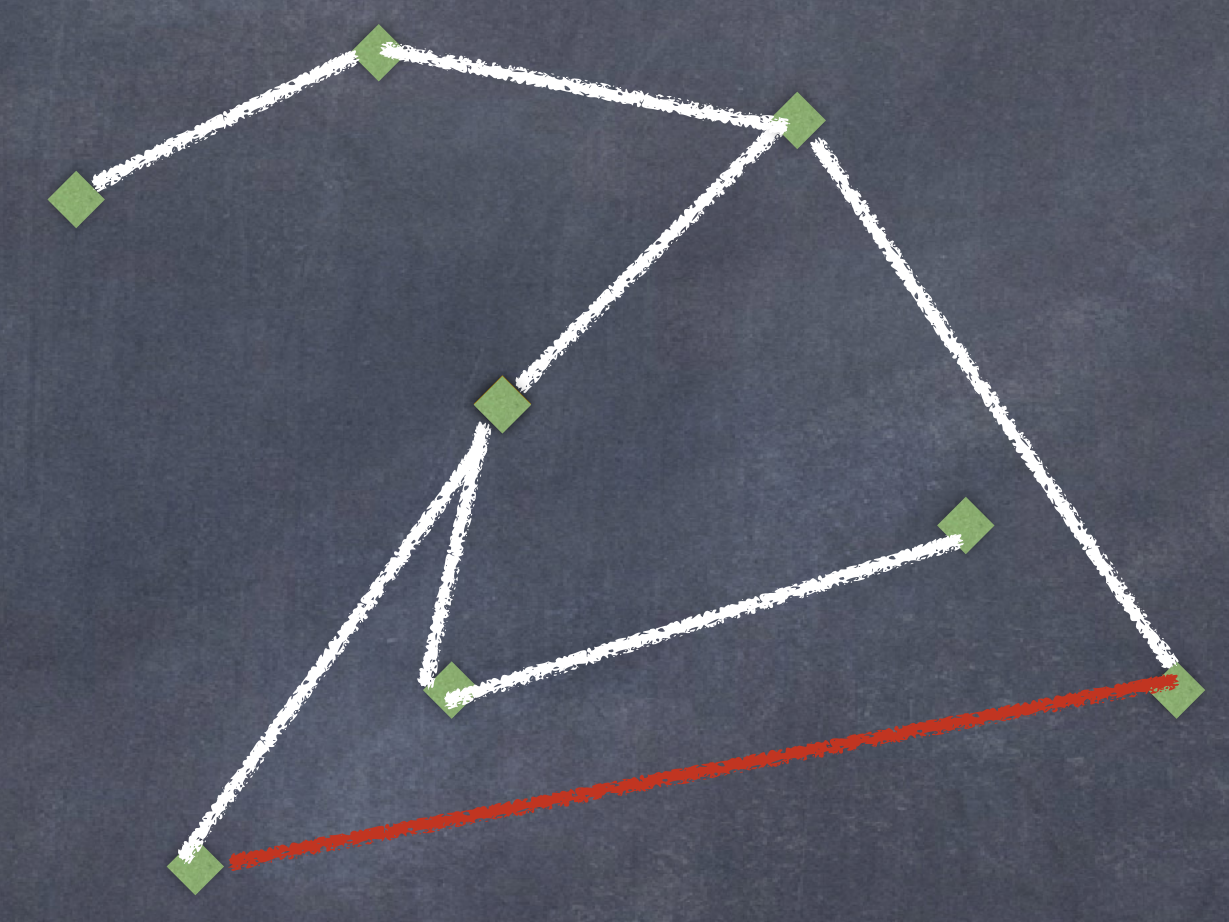# Actual network

To tackle the problem, label ourselves

- **The problem:**
Given such an observed network, tell me about the root.

- **The problem:**
  Given such an observed network, tell me about the root.


- **Somewhat flexible:**
  Give me a set of vertices that contains the root with some confidence.

- **The problem:**
  Given such an observed network, tell me about the root.

- **Somewhat flexible:**
  Give me a set of vertices that contains the root with some confidence.

- **More concretely:**
  Give me a set of vertices $C(G^*) \subseteq V(G^*)$ such that $\mathbb{P}(\blacklozenge \in C(G^*)) \geq 95\,\%$.

- So our goal now:
Given $\epsilon \in (0,1)$, find $C_\epsilon \subseteq V = \{A, B, \cdots\}$ such that $\mathbb{P}(\blacklozenge \in C_\epsilon(\boldsymbol{G}^*)) \geq 1 - \epsilon$.

- So our goal now:

Given $\epsilon \in (0,1)$, find $C_\epsilon \subseteq V = \{A, B, \cdots\}$ such that $\mathbb{P}(\blacklozenge \in C_\epsilon(G^*)) \geq 1 - \epsilon$.

- Trivial:

Take $C_\epsilon(G^*) = V(G^*)$. Works for all $\epsilon$.

- So our goal now:

  Given $\epsilon \in (0,1)$, find $C_\epsilon \subseteq V = \{A, B, \cdots\}$ such that $\mathbb{P}(\blacklozenge \in C_\epsilon(G^*)) \geq 1 - \epsilon$.

- Trivial:

  Take $C_\epsilon(G^*) = V(G^*)$. Works for all $\epsilon$.

- Really, the problem asks for:

  Smallest possible $C_\epsilon$.

One issue with $C_\epsilon(G^*)$:

One issue with $C_\epsilon(G^*)$:

One issue with $C_\epsilon(G^*)$:



$C_\epsilon$ either contains all ◆ or contains none of them.

$C_\epsilon(G^*)$ should be labelling equivariant

$C_\epsilon(G^*)$ should be labelling equivariant

which means that $\tau C_\epsilon(G^*) = C_\epsilon(\tau G^*)$ for all relabelings $\tau$ of $G$.

$C_\epsilon(G^*)$ should be labelling equivariant

which means that $\tau C_\epsilon(G^*) = C_\epsilon(\tau G^*)$ for all relabelings $\tau$ of $G$.

Use randomization to break ties.
In that case, $\tau C_\epsilon(\boldsymbol{G}^*) \underset{d}{=} C_\epsilon(\tau \boldsymbol{G}^*)$ for every relabelling $\tau$ of $\boldsymbol{G}^*$.

# The construction for $C_\epsilon(\cdot)$

# The construction for $C_\epsilon(\cdot)$

- Say we have a labelled observed graph $\tilde{G} = g$ with randomized labels from $G*$.

# The construction for $C_\epsilon(\cdot)$

- Say we have a labelled observed graph $\tilde{G} = g$ with randomized labels from $G^*$.

- Say $u_i$ is the node which is $i^{th}$ most likely to be the root. That is,
$\mathbb{P}(\blacklozenge = u_1 \mid \tilde{G} = g) \geq \mathbb{P}(\blacklozenge = u_2 \mid \tilde{G} = g) \geq \cdots$

# The construction for $C_\epsilon(\cdot)$

- Say we have a labelled observed graph $\tilde{G} = g$ with randomized labels from $G*$.

- Say $u_i$ is the node which is $i^{\text{th}}$ most likely to be the root. That is,
$\mathbb{P}(\blacklozenge = u_1 \mid \tilde{G} = g) \geq \mathbb{P}(\blacklozenge = u_2 \mid \tilde{G} = g) \geq \cdots$

- Take smallest $k$ such that $\displaystyle\sum_{i=1}^{k} \mathbb{P}(\blacklozenge = u_i \mid \tilde{G} = g) \geq 1 - \epsilon$. This is our Bayesian

coverage set: $B_\epsilon(g) = \{u_1, \cdots, u_k\}$.

# The construction for $C_\epsilon(\cdot)$

- Say we have a labelled observed graph $\tilde{G} = g$ with randomized labels from $G*$.

- Say $u_i$ is the node which is $i^{th}$ most likely to be the root. That is,
  $$\mathbb{P}(\blacklozenge = u_1 \mid \tilde{G} = g) \geq \mathbb{P}(\blacklozenge = u_2 \mid \tilde{G} = g) \geq \cdots$$

- Take smallest $k$ such that $\displaystyle\sum_{i=1}^{k} \mathbb{P}(\blacklozenge = u_i \mid \tilde{G} = g) \geq 1 - \epsilon$. This is our Bayesian coverage set: $B_\epsilon(g) = \{u_1, \cdots, u_k\}$.

- This is, in fact, an honest coverage set: if $G*$ is an alphabetically labelled observation of $G$ (whose root is $\blacklozenge$), then $\mathbb{P}(\blacklozenge \in B_\epsilon(G*)) \geq 1 - \epsilon$.

# The construction for $C_\epsilon(\cdot)$

- Say we have a labelled observed graph $\tilde{G} = g$ with randomized labels from $G*$.

- Say $u_i$ is the node which is $i^{\text{th}}$ most likely to be the root. That is,
$$\mathbb{P}(\blacklozenge = u_1 \mid \tilde{G} = g) \geq \mathbb{P}(\blacklozenge = u_2 \mid \tilde{G} = g) \geq \cdots$$

- Take smallest $k$ such that $\displaystyle\sum_{i=1}^{k} \mathbb{P}(\blacklozenge = u_i \mid \tilde{G} = g) \geq 1 - \epsilon$. This is our Bayesian coverage set: $B_\epsilon(g) = \{u_1, \cdots, u_k\}$.

- This is, in fact, an honest coverage set: if $G*$ is an alphabetically labelled observation of $G$ (whose root is $\blacklozenge$), then $\mathbb{P}(\blacklozenge \in B_\epsilon(G*)) \geq 1 - \epsilon$.

Computing $\mathbb{P}(\blacklozenge = u \ \tilde{G} = g)$

# Computing $\mathbb{P}(\blacklozenge = u \mid \tilde{G} = g)$

- Let $\Pi$ be a randomized labelling such that $\Pi G^* = \tilde{G}$.

# Computing $\mathbb{P}(\blacklozenge = u \mid \tilde{G} = g)$

- Let $\Pi$ be a randomized labelling such that $\Pi G^* = \tilde{G}$.

- Key observation 1:
$$\mathbb{P}(\blacklozenge = u \mid \tilde{G} = g) = \mathbb{P}\left(\Pi(1) = u \mid \tilde{G} = g\right) = \sum_{\pi} \mathbf{1}_{\{\pi(1)=u\}} \mathbb{P}\left(\Pi = \pi \mid \tilde{G} = g\right).$$

# Computing $\mathbb{P}(\diamond = u \mid \tilde{G} = g)$

- Let $\Pi$ be a randomized labelling such that $\Pi G^* = \tilde{G}$.

- Key observation 1:
$$\mathbb{P}(\diamond = u \mid \tilde{G} = g) = \mathbb{P}\left(\Pi(1) = u \mid \tilde{G} = g\right) = \sum_{\pi} \mathbf{1}_{\{\pi(1)=u\}} \mathbb{P}\left(\Pi = \pi \mid \tilde{G} = g\right).$$

- Key observation 2: $\mathbb{P}\left(\Pi = \pi \mid \tilde{G} = g\right) = \dfrac{\mathbb{P}\left(\tilde{G} = g \mid \Pi = \pi\right)}{\sum_{\pi'} \mathbb{P}\left(\tilde{G} = g \mid \Pi = \pi'\right)}.$

# Computing $\mathbb{P}(\blacklozenge = u \mid \tilde{G} = g)$

- Let $\Pi$ be a randomized labelling such that $\Pi G^* = \tilde{G}$.

- Key observation 1:
$$\mathbb{P}(\blacklozenge = u \mid \tilde{G} = g) = \mathbb{P}\left(\Pi(1) = u \mid \tilde{G} = g\right) = \sum_{\pi} \mathbf{1}_{\{\pi(1) = u\}} \mathbb{P}\left(\Pi = \pi \mid \tilde{G} = g\right).$$

- Key observation 2: $\mathbb{P}\left(\Pi = \pi \mid \tilde{G} = g\right) = \dfrac{\mathbb{P}\left(\tilde{G} = g \mid \Pi = \pi\right)}{\sum_{\pi'} \mathbb{P}\left(\tilde{G} = g \mid \Pi = \pi'\right)}.$

- Useful (because of the noise): $\mathbb{P}\left(\Pi = \pi, \tilde{T} = t \mid \tilde{G} = g\right).$

# Computing $\mathbb{P}(\blacklozenge = u \mid \tilde{G} = g)$

- Let $\Pi$ be a randomized labelling such that $\Pi G^* = \tilde{G}$.

- Key observation 1:
$$\mathbb{P}(\blacklozenge = u \mid \tilde{G} = g) = \mathbb{P}\left(\Pi(1) = u \mid \tilde{G} = g\right) = \sum_{\pi} \mathbf{1}_{\{\pi(1)=u\}} \mathbb{P}\left(\Pi = \pi \mid \tilde{G} = g\right).$$

- Key observation 2: $\mathbb{P}\left(\Pi = \pi \mid \tilde{G} = g\right) = \dfrac{\mathbb{P}\left(\tilde{G} = g \mid \Pi = \pi\right)}{\sum_{\pi'} \mathbb{P}\left(\tilde{G} = g \mid \Pi = \pi'\right)}.$

- Useful (because of the noise): $\mathbb{P}\left(\Pi = \pi, \tilde{T} = t \mid \tilde{G} = g\right).$

Will compute this

# A primer on Gibbs sampling

Interested in: $h(i,j) = \mathbf{1}_{(t,\pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

Interested in: $h(i, j) = \mathbf{1}_{(t, \pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

$$X : \tilde{T}$$
$$Y : \Pi$$

Interested in: $h(i, j) = \mathbf{1}_{(t, \pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

$$X : \tilde{T}$$
$$Y : \Pi$$

The math?

Interested in: $h(i, j) = \mathbf{1}_{(t, \pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

The math?

$$X: \tilde{T}$$
$$Y: \Pi$$

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y=j) \to (X=k) \to (Y=l)\right] = p(X=k \mid Y=j) \cdot p(Y=l \mid X=k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.$$

Interested in: $h(i, j) = \mathbf{1}_{(t, \pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

The math?

$$
\boxed{\begin{array}{l} X: \ \tilde{T} \\ Y: \ \Pi \end{array}}
$$

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$
\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y = j) \to (X = k) \to (Y = l)\right] = p(X = k \mid Y = j) \cdot p(Y = l \mid X = k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.
$$

Interested in: $h(i, j) = \mathbf{1}_{(t, \pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

The math?

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y = j) \to (X = k) \to (Y = l)\right] = p(X = k \mid Y = j) \cdot p(Y = l \mid X = k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.$$

$\{Z_n\}$ is an irreducible aperiodic Markov chain with stationary distribution $q$.

Interested in: $h(i, j) = \mathbf{1}_{(t,\pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

The math?

$$\begin{array}{l} X : \tilde{T} \\ Y : \Pi \end{array}$$

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y = j) \to (X = k) \to (Y = l)\right] = p(X = k \mid Y = j) \cdot p(Y = l \mid X = k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.$$

$\{Z_n\}$ is an irreducible aperiodic Markov chain with stationary distribution $q$.

$\xrightarrow{\text{LLN}}$ if $h$ is a bounded function then $\sum_{i,j} h(i, j) p(i, j)$ can be approximated by $\frac{1}{n} \sum_{i=1}^{n} h(X_i, Y_i)$.

Interested in: $h(i, j) = \mathbf{1}_{(t,\pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

$$\begin{aligned} X &: \ \tilde{T} \\ Y &: \ \Pi \end{aligned}$$

The math?

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y = j) \rightarrow (X = k) \rightarrow (Y = l)\right] = p(X = k \mid Y = j) \cdot p(Y = l \mid X = k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.$$

$\{Z_n\}$ is an irreducible aperiodic Markov chain with stationary distribution $q$.

$\overset{\text{LLN}}{\longrightarrow}$ if $h$ is a bounded function then $\sum_{i,j} h(i,j)p(i,j)$ can be approximated by $\frac{1}{n}\sum_{i=1}^{n} h(X_i, Y_i)$.

Interested in: $h(i,j) = \mathbf{1}_{(t,\pi)}$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

The math?

$$\boxed{\begin{array}{l} X : \tilde{T} \\ Y : \Pi \end{array}}$$

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y = j) \to (X = k) \to (Y = l)\right] = p(X = k \mid Y = j) \cdot p(Y = l \mid X = k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.$$

$\{Z_n\}$ is an irreducible aperiodic Markov chain with stationary distribution $q$.

$\xrightarrow{\text{LLN}}$ if $h$ is a bounded function then $\sum_{i,j} h(i,j)p(i,j)$ can be approximated by $\dfrac{1}{n}\sum_{i=1}^{n} h(X_i, Y_i)$.

Interested in: $h(i,j) = \mathbf{1}_{(t,\pi)}$

Start with $X = x_0$

Sample $Y = y_0$ from $p(y \mid x_0)$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.

The math?

$X : \tilde{T}$
$Y : \Pi$

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y = j) \to (X = k) \to (Y = l)\right] = p(X = k \mid Y = j) \cdot p(Y = l \mid X = k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.$$

$\{Z_n\}$ is an irreducible aperiodic Markov chain with stationary distribution $q$.

$\xrightarrow{\text{LLN}}$ if $h$ is a bounded function then $\sum_{i,j} h(i,j) p(i,j)$ can be approximated by $\frac{1}{n} \sum_{i=1}^{n} h(X_i, Y_i)$.

Interested in: $h(i,j) = \mathbf{1}_{(t,\pi)}$

Start with $X = x_0$

Sample $Y = y_0$ from $p(y \mid x_0)$

Sample $X = x_1$ from $p(x \mid y_0)$

# A primer on Gibbs sampling

When? Compute $p(X, Y)$ when $p(X \mid Y)$ and $p(Y \mid X)$ are known/easy to know.
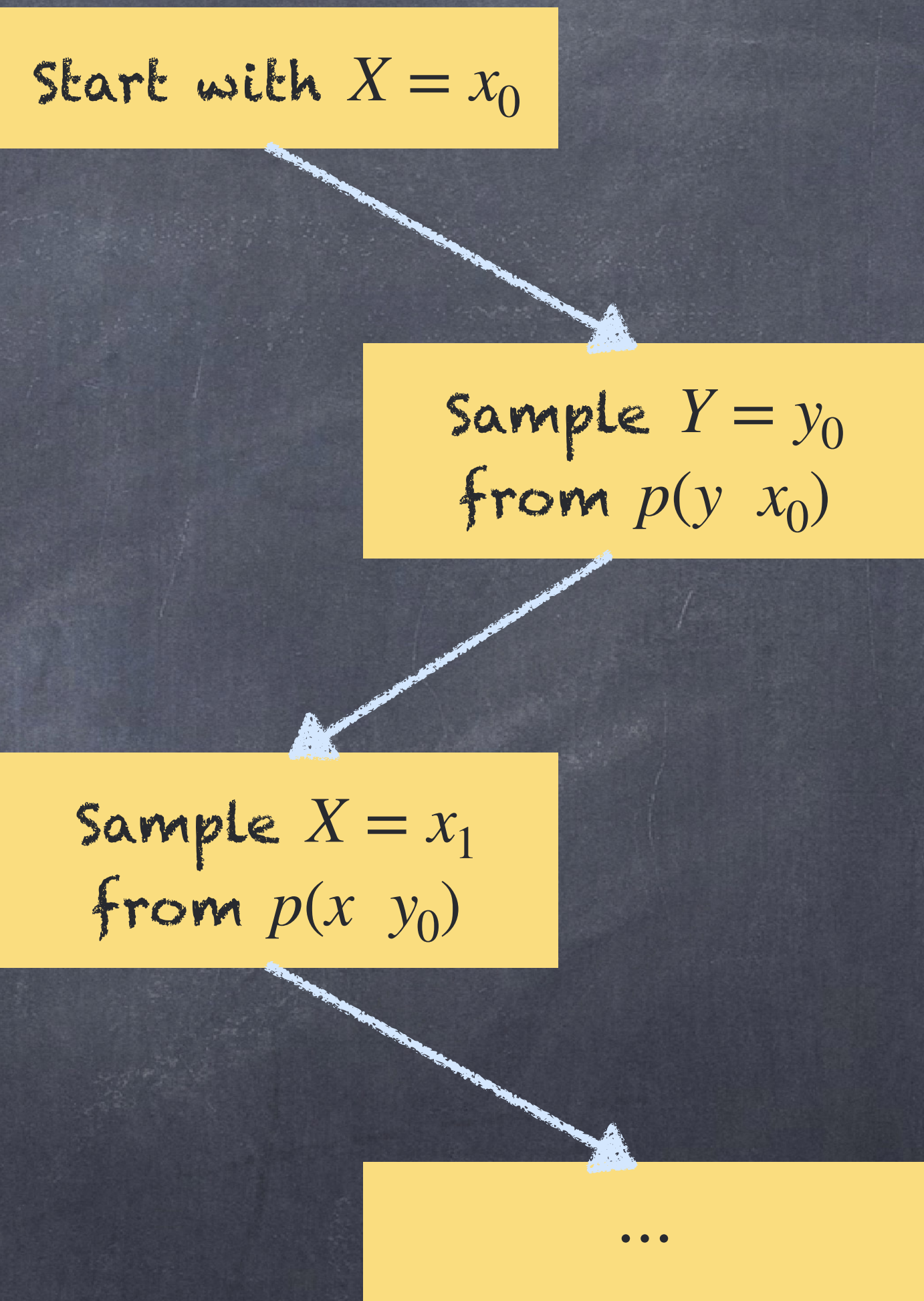
The math?

$$X : \tilde{T}$$
$$Y : \Pi$$

Let $q$ be the distribution of $Z = (X, Y)$. Generate a Markov chain $Z_n = (X_n, Y_n)$ with transition

$$\tilde{q}_{(ij),(kl)} = \mathbb{P}\left[(Y=j) \to (X=k) \to (Y=l)\right] = p(X=k \mid Y=j) \cdot p(Y=l \mid X=k) = \frac{q_{kj}}{\sum_t q_{tj}} \cdot \frac{q_{kl}}{\sum_t q_{kt}}.$$

$\{Z_n\}$ is an irreducible aperiodic Markov chain with stationary distribution $q$.

$\xrightarrow{\text{LLN}}$ if $h$ is a bounded function then $\sum_{i,j} h(i,j)p(i,j)$ can be approximated by $\frac{1}{n}\sum_{i=1}^{n} h(X_i, Y_i)$.

Interested in: $h(i,j) = \mathbf{1}_{(t,\pi)}$

Start with $X = x_0$

Sample $Y = y_0$ from $p(y \mid x_0)$

Sample $X = x_1$ from $p(x \mid y_0)$

...

# The Gibbs sampler algorithm for our case

# The Gibbs sampler algorithm for our case

- Remember that $X = \tilde{T}, Y = \Pi$.

# The Gibbs sampler algorithm for our case

- Remember that $X = \tilde{T}, Y = \Pi$.

- So we alternate between two stages:

# The Gibbs sampler algorithm for our case

- Remember that $X = \tilde{T}, Y = \Pi$.

- So we alternate between two stages:

  - Fix $t$ and generate $\pi$ from distribution $\mathbb{P}\left(\Pi = \pi \mid \tilde{T} = t, \tilde{G} = g\right)$.

# The Gibbs sampler algorithm for our case

- Remember that $X = \tilde{T}, Y = \Pi$.

- So we alternate between two stages:

    - Fix $t$ and generate $\pi$ from distribution $\mathbb{P}\left(\Pi = \pi \mid \tilde{T} = t, \tilde{G} = g\right)$.

    - Fix $\pi$ and generate $t$ from the distribution $\mathbb{P}\left(\tilde{T} = t \mid \Pi = \pi, \tilde{G} = g\right)$.

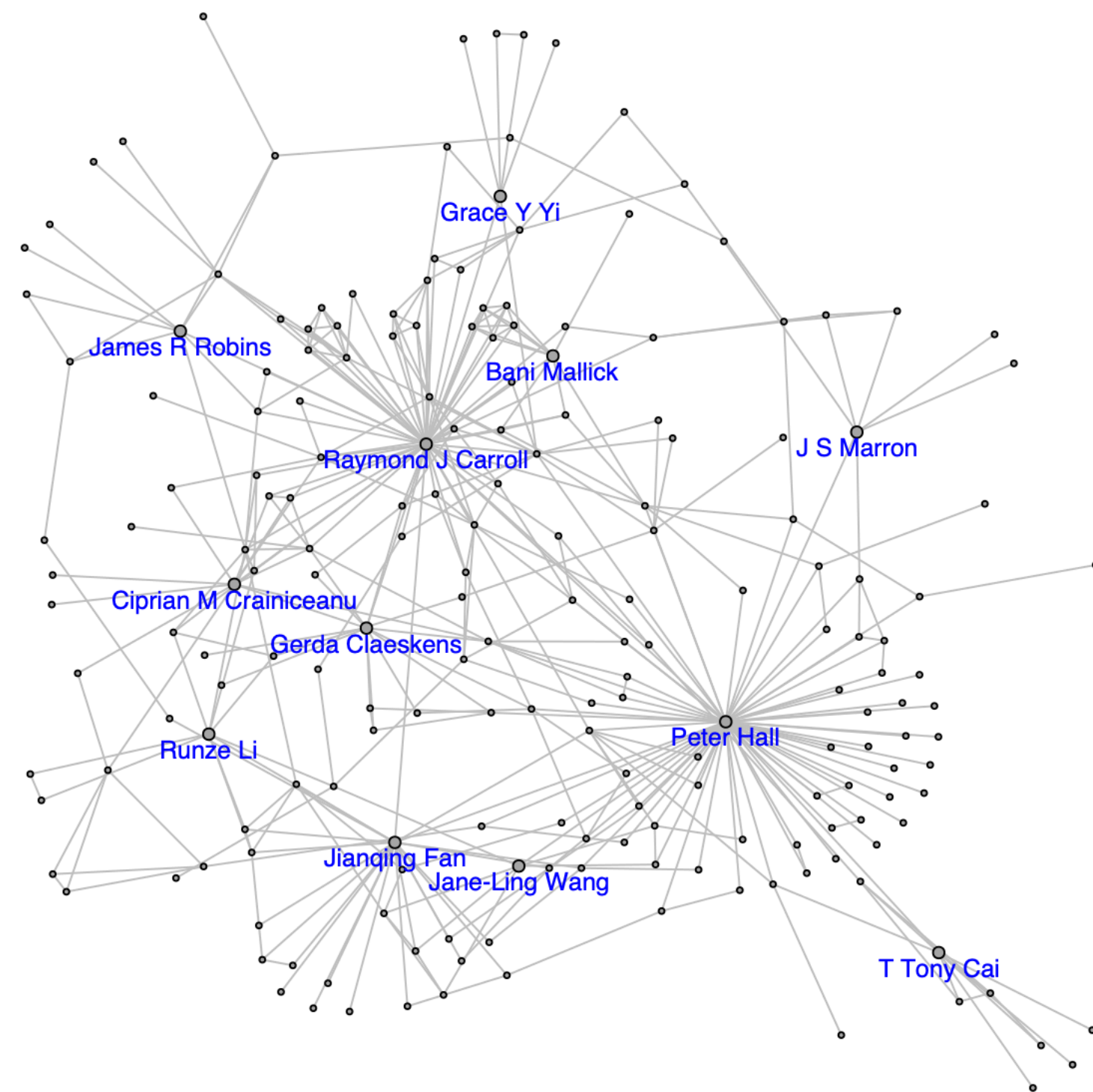# An example (taken from the manuscript)



Figure 20: Subgraph of the co-authorship graph comprising the 200 nodes with the highest posterior root probabilities. We label the 12 nodes with the highest root probabilities.

# Bibliography

[CX23]   Harry Crane and Min Xu. Root and community inference on the latent growth process of a network. 2023. arXiv: 2107.00153 [stat.ME].

[CX21]   Harry Crane and Min Xu. Inference on the history of a randomly growing tree. Journal of the Royal Society of Statistics Series B, Volume 83, Issue 4, September 2021, Pages 639-668, https://doi.org/10.1111/rssb.12428.