

Shared Information for a Markov Chain on a Tree

Sagnik Bhattacharya , Prakash Narayan
ECE-ISR@UMD

ISIT 2022

July 1, 2022

#1: Capturing dependence among multiple rvs – *Shared Information*

How to capture dependence among multiple rvs?

Shared information

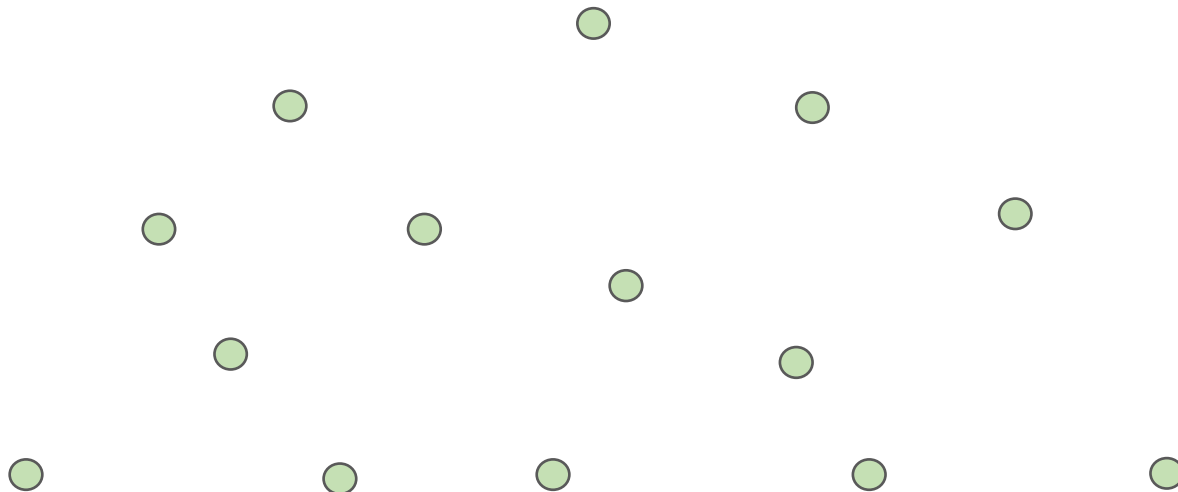
$$\text{SI}(X_{\mathcal{M}}) = \min_{2 \leq k \leq m} \min_{\pi = (\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}})$$

Shared information

$$\begin{aligned} \text{SI}(X_{\mathcal{M}}) &= \min_{2 \leq k \leq m} \min_{\pi = (\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}}) \\ &= \min_{2 \leq k \leq m} \min_{\pi = (\pi_u, u=1, \dots, k)} \frac{1}{k-1} \left[\sum_{u=1}^k H(X_{\pi_u}) - H(X_{\mathcal{M}}) \right] \end{aligned}$$

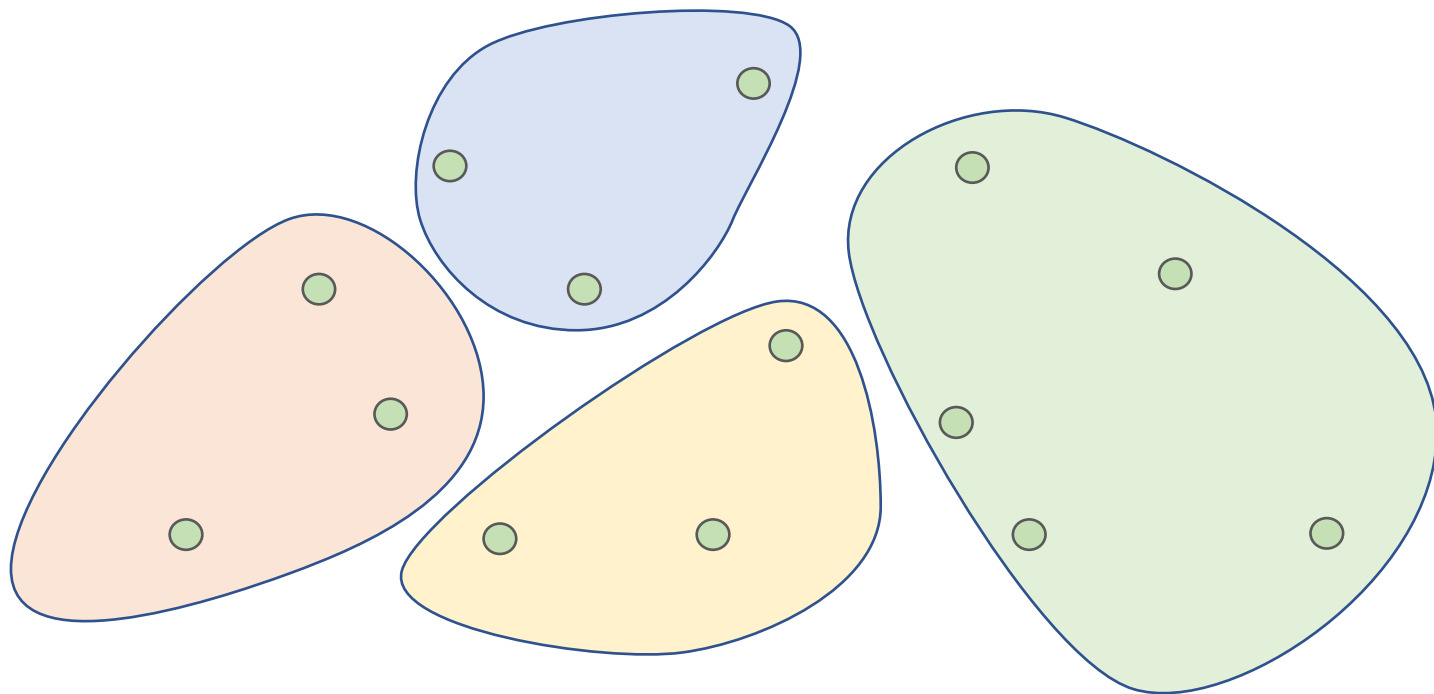
Shared information

$$\begin{aligned} \text{SI}(X_{\mathcal{M}}) &= \min_{2 \leq k \leq m} \min_{\pi=(\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}}) \\ &= \min_{2 \leq k \leq m} \min_{\pi=(\pi_u, u=1, \dots, k)} \frac{1}{k-1} \left[\sum_{u=1}^k H(X_{\pi_u}) - H(X_{\mathcal{M}}) \right] \end{aligned}$$



Shared information

$$\begin{aligned} \text{SI}(X_{\mathcal{M}}) &= \min_{2 \leq k \leq m} \min_{\pi=(\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}}) \\ &= \min_{2 \leq k \leq m} \min_{\pi=(\pi_u, u=1, \dots, k)} \frac{1}{k-1} \left[\sum_{u=1}^k H(X_{\pi_u}) - H(X_{\mathcal{M}}) \right] \end{aligned}$$



Some special cases

- Two rvs

$$SI(X_1, X_2) = \text{mutual information } I(X_1 \wedge X_2)$$

Some special cases

- Two rvs

$$SI(X_1, X_2) = \text{mutual information } I(X_1 \wedge X_2)$$

- Three rvs

- Minimum of

$$I(X_1 \wedge X_2, X_3) \quad I(X_2 \wedge X_1, X_3) \quad I(X_3 \wedge X_1, X_2)$$

$$\frac{1}{2} [H(X_1) + H(X_2) + H(X_3) - H(X_1, X_2, X_3)]$$

Shared information - operational significance

- Secret key capacity of a multiterminal source model
 - [Csiszár-Narayan, 2004]

Shared information - operational significance

- Secret key capacity of a multiterminal source model
 - [Csiszár-Narayan, 2004]
 - [Chan, 2008]
 - [Chan-Zheng, 2016]
 - [Chan, 2011]

Shared information - operational significance

- Secret key capacity of a multiterminal source model
 - [Csiszar-Narayan, 2004]
 - [Chan, 2008]
 - [Chan-Zheng, 2016]
 - [Chan, 2011]
- Hypothesis Testing
 - [Tyagi-Watanabe, 2016]

Shared information - operational significance

- Secret key capacity of a multiterminal source model
 - [Csiszar-Narayan, 2004]
 - [Chan, 2008]
 - [Chan-Zheng, 2016]
 - [Chan, 2011]
- Hypothesis Testing
 - [Tyagi-Watanabe, 2016]
- Extensively studied in [Chan-Bashabsheh-Ebrahimi-Kaced-Liu, 2015]

Shared information - operational significance

- Secret key capacity of a multiterminal source model
 - [Csiszar-Narayan, 2004]
 - [Chan, 2008]
 - [Chan-Zheng, 2016]
 - [Chan, 2011]
- Hypothesis Testing
 - [Tyagi-Watanabe, 2016]
- Extensively studied in [Chan-Bashabsheh-Ebrahimi-Kaced-Liu, 2015]
- Clustering
 - [Chan-Bashabsheh-Zhou-Kaced-Liu, 2016]

Shared information - computation

$$\text{SI}(X_{\mathcal{M}}) = \min_{2 \leq k \leq m} \min_{\pi = (\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}})$$

- When underlying pmf is known, there is an efficient algorithm to compute SI
 - Submodular optimization [\[CBEKL15\]](#)

Shared information - computation

$$\text{SI}(X_{\mathcal{M}}) = \min_{2 \leq k \leq m} \min_{\pi = (\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}})$$

- When underlying pmf is known, there is an efficient algorithm to compute SI
 - Submodular optimization [CBEKL15]
- When pmf is unknown, estimation involves prohibitively massive search space.

Shared information - computation

$$\text{SI}(X_{\mathcal{M}}) = \min_{2 \leq k \leq m} \min_{\pi=(\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}})$$

- When underlying pmf is known, there is an efficient algorithm to compute SI
 - Submodular optimization [CBEKL15]
- When pmf is unknown, estimation involves prohibitively massive search space.

0 1 3 6 2 7
 : 13
 : 20
 23 12
 10 22 11 21

THE ON-LINE ENCYCLOPEDIA
 OF INTEGER SEQUENCES[®]

founded in 1964 by N. J. A. Sloane

[Hints](#)

(Greetings from [The On-Line Encyclopedia of Integer Sequences!](#))

A000110	Bell or exponential numbers: number of ways to partition a set of n labeled elements. (Formerly M1484 N0585)	1164
---------	---	------

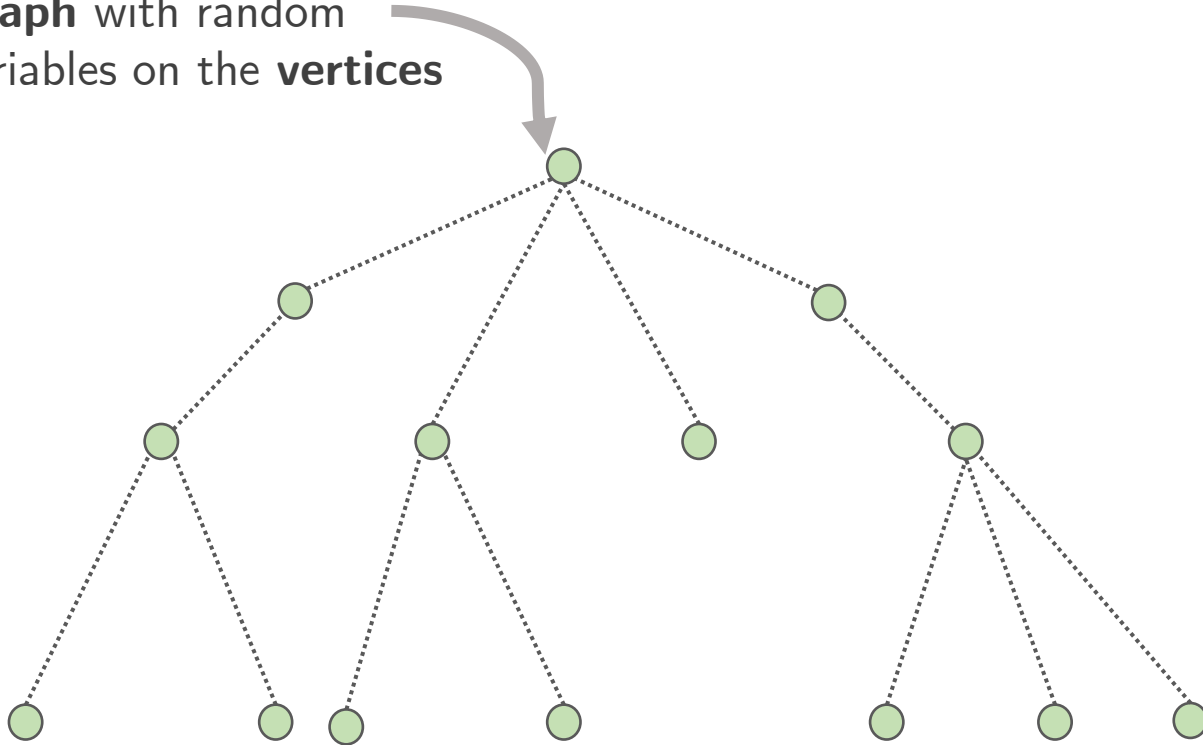
1, 1, 2, 5, 15, 52, 203, 877, 4140, 21147, 115975, 678570, 4213597, 27644437, 190899322,
 1382958545, 10480142147, 82864869804, 682076806159, 5832742205057, 51724158235372, 474869816156751,
 4506715738447323, 44152005855084346, 445958869294805289, 4638590332229999353, 49631246523618756274

Want: simpler forms in special cases

#2: Markov Chain on a Tree (MCT)

Markov chain on a tree (MCT)

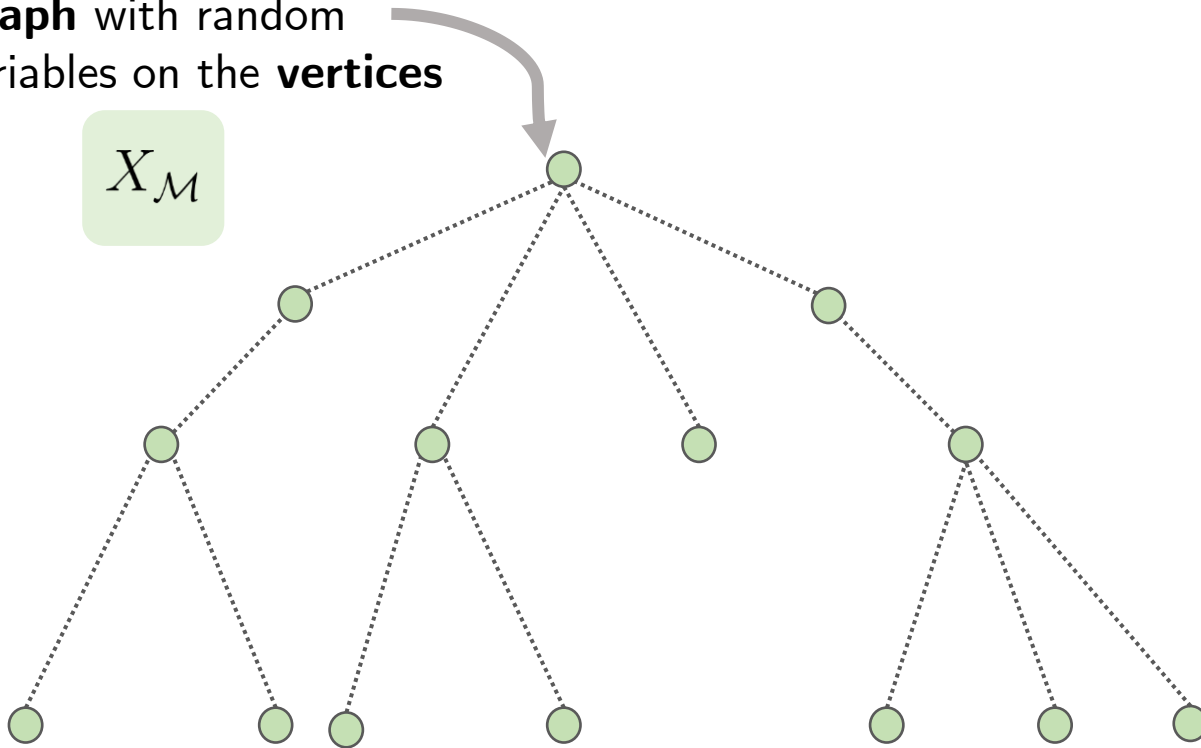
Graph with random variables on the **vertices**



Markov chain on a tree (MCT)

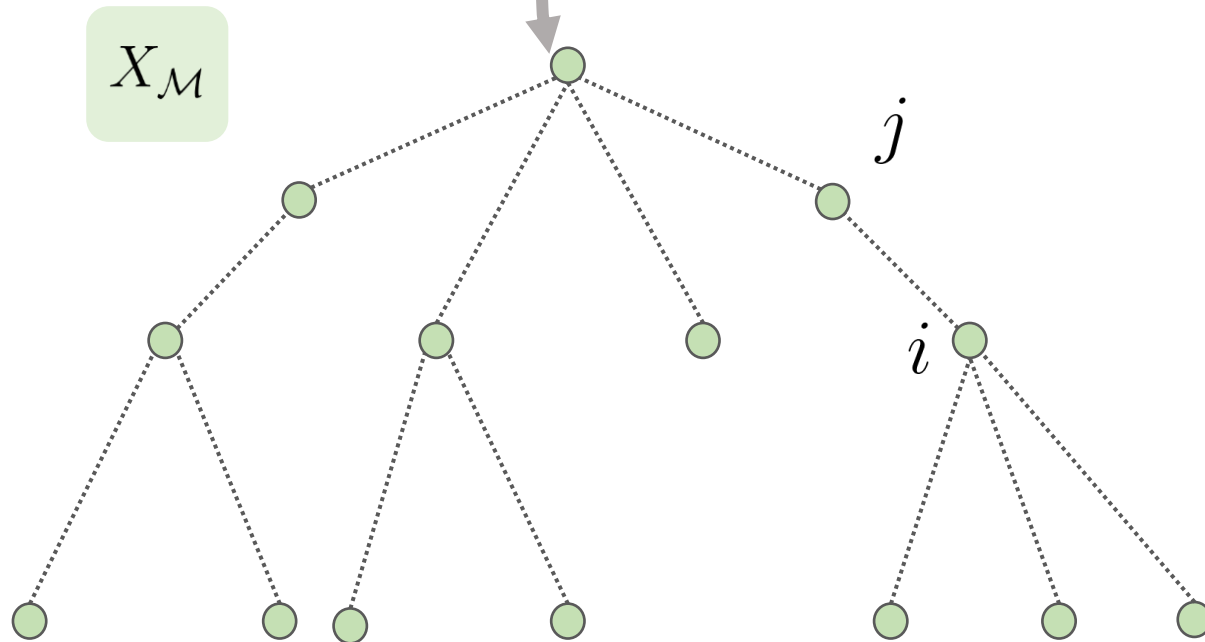
Graph with random variables on the **vertices**

$X_{\mathcal{M}}$



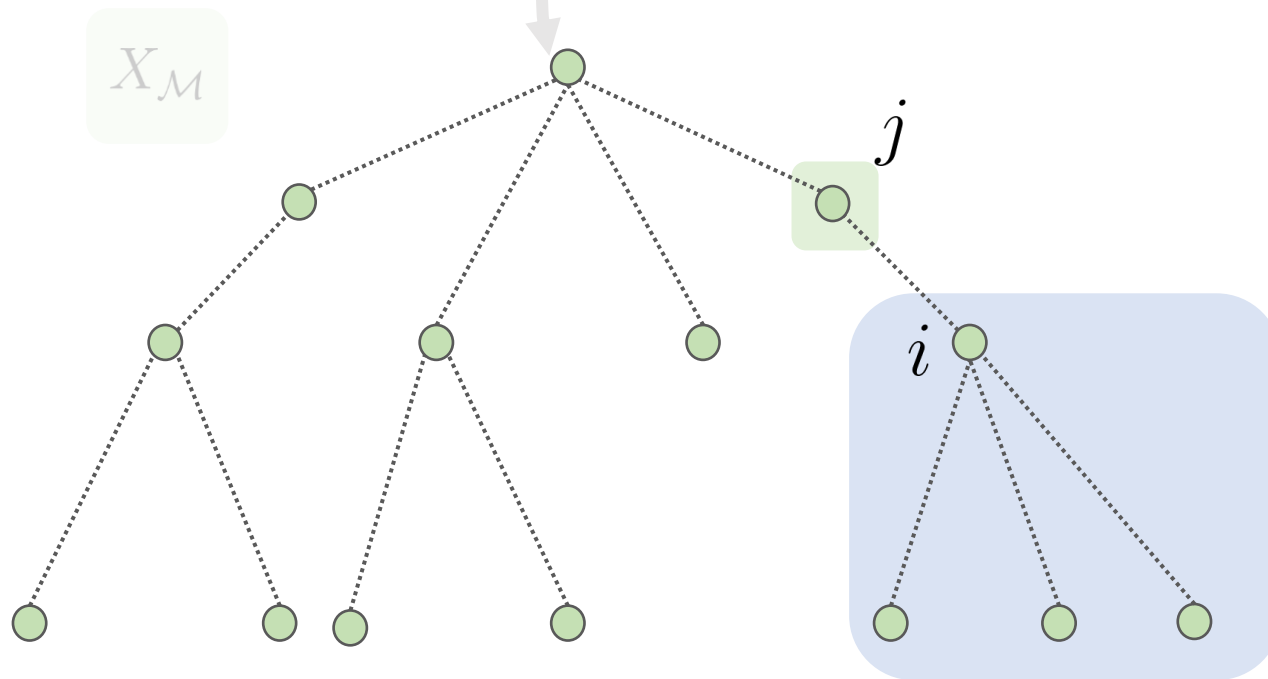
Markov chain on a tree (MCT)

Graph with random variables on the **vertices**



Markov chain on a tree (MCT)

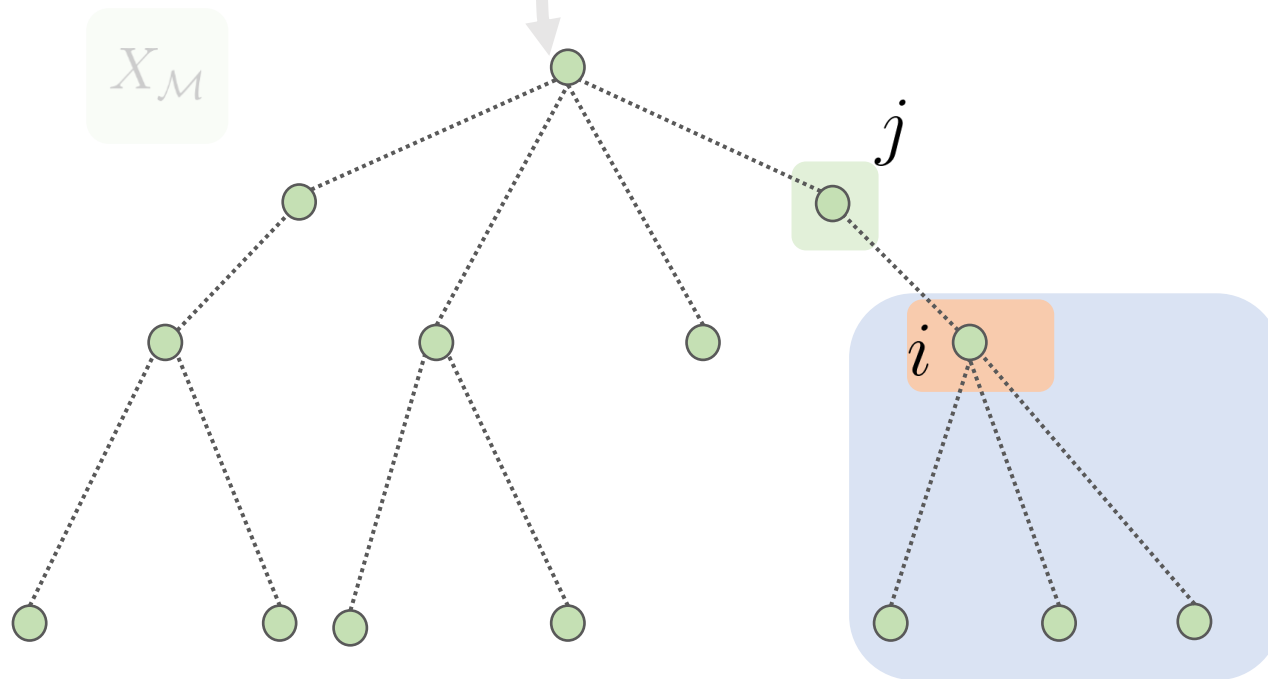
Graph with random variables on the **vertices**



$$P_{X_j} \mid \mathcal{B}(i \leftarrow j)$$

Markov chain on a tree (MCT)

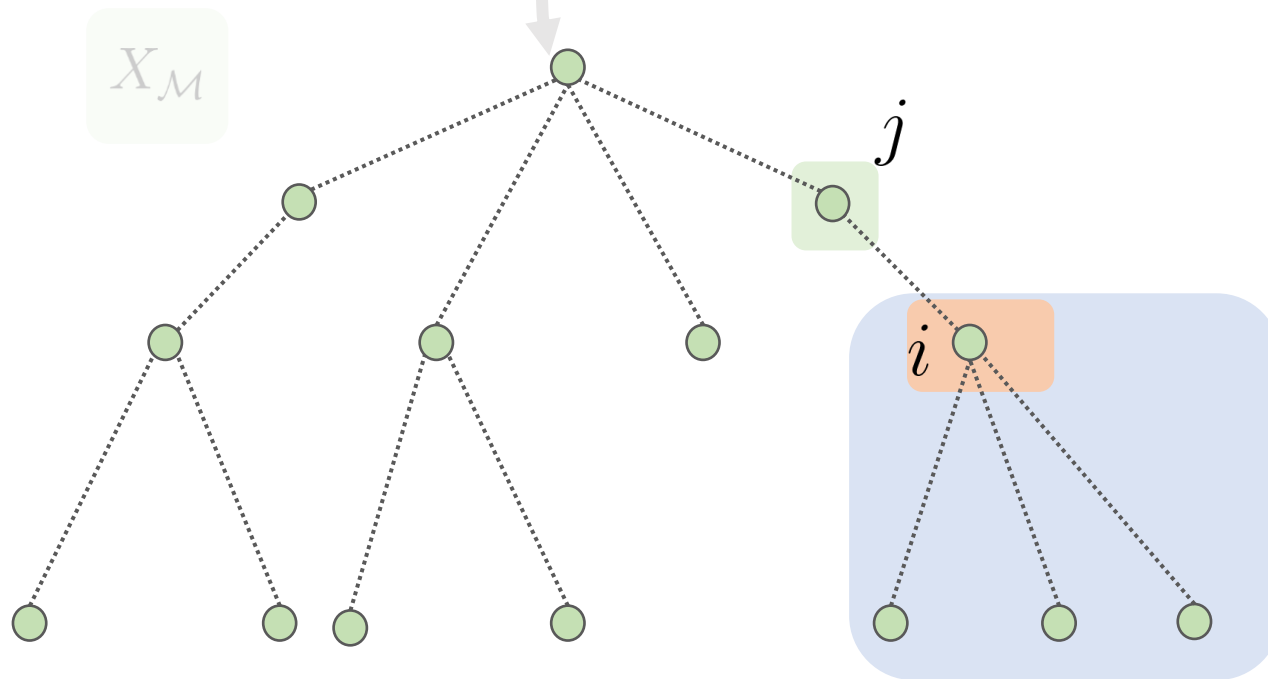
Graph with random variables on the **vertices**



$$P_{X_j} \mid \mathcal{B}(i \leftarrow j)$$

Markov chain on a tree (MCT)

Graph with random variables on the **vertices**



$$P_{X_j \mid \mathcal{B}(i \leftarrow j)} = P_{X_j \mid X_i}$$

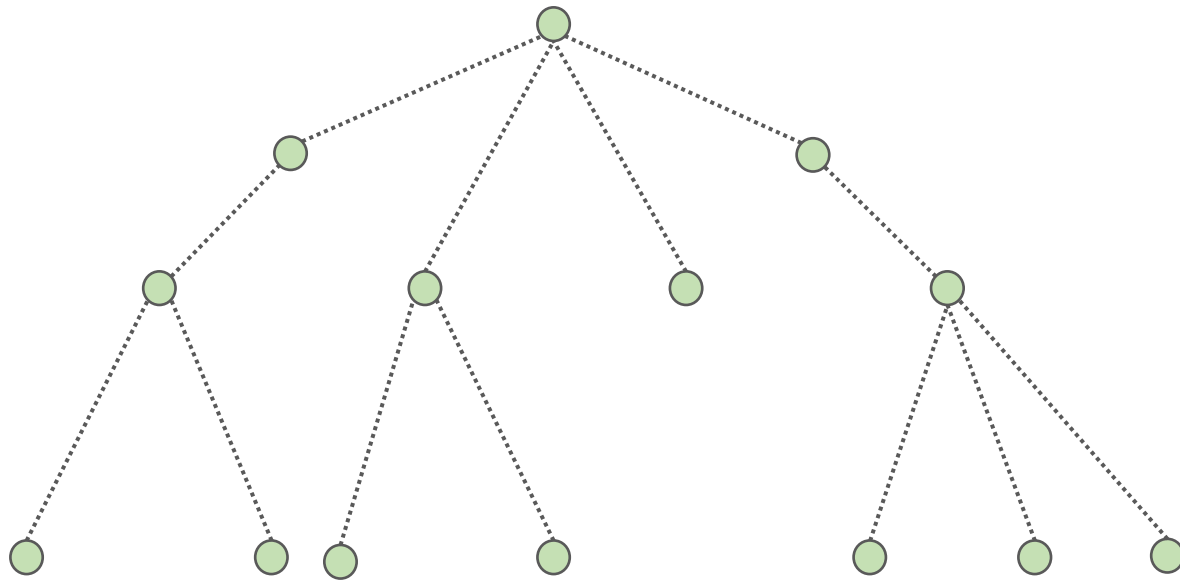
#3: SI in an MCT

SI in an MCT

- 2-partition achieves SI [[Csiszár-Narayan, 2004](#)]

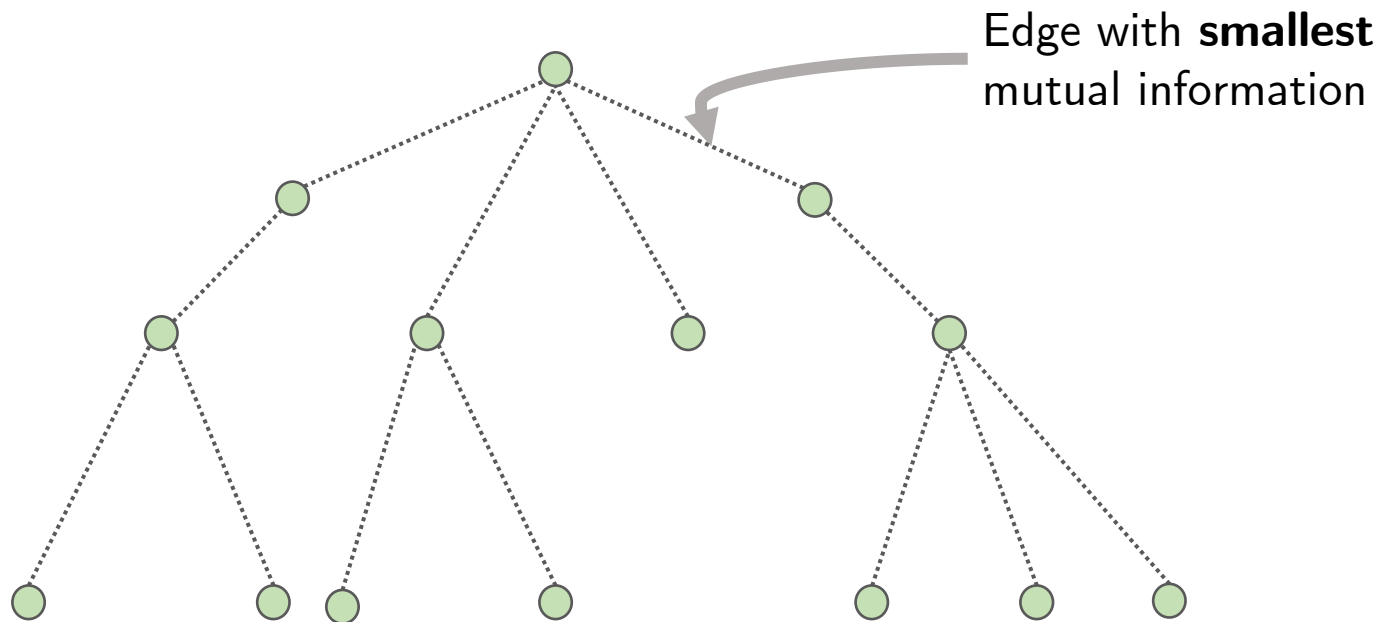
SI in an MCT

- 2-partition achieves SI [[Csiszár-Narayan, 2004](#)]



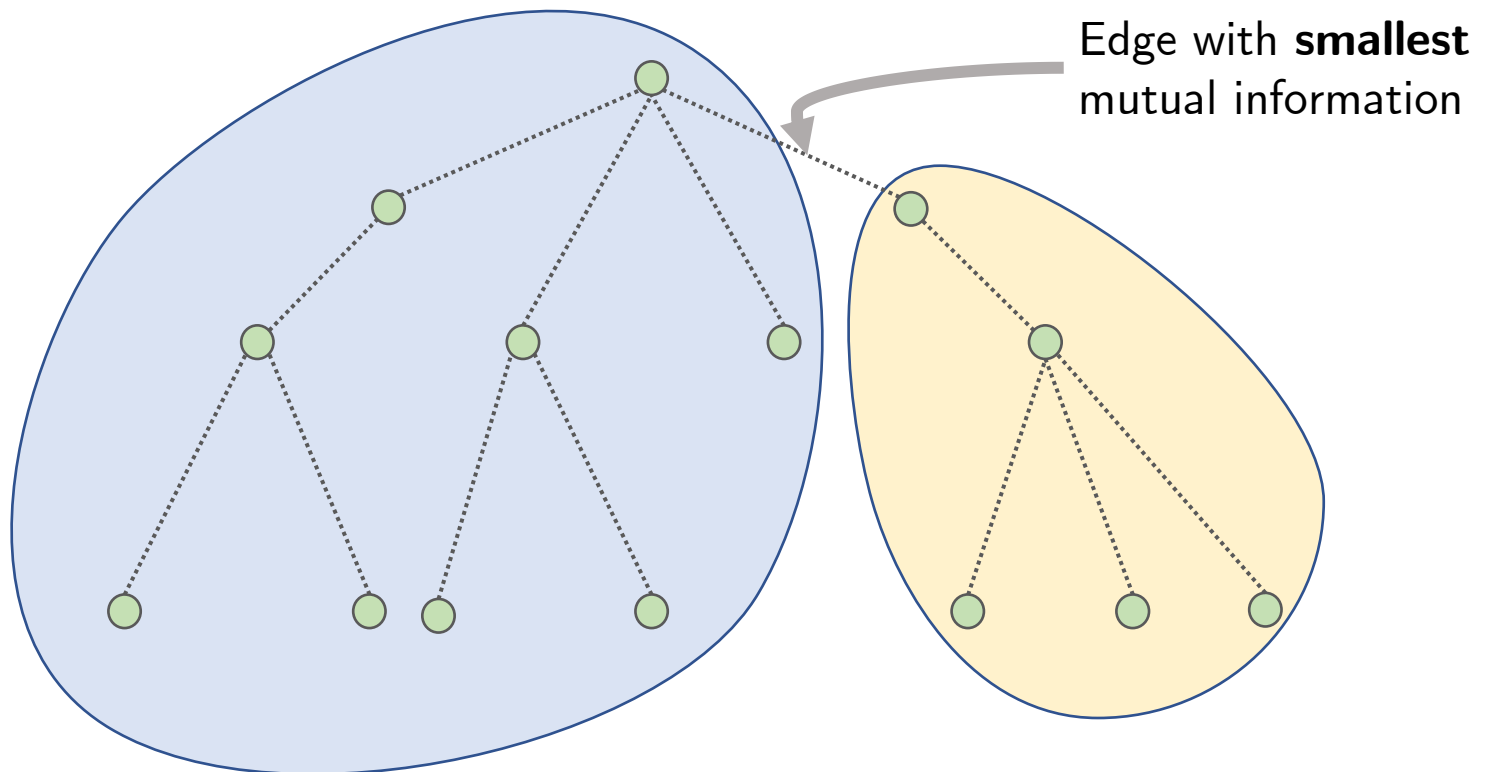
SI in an MCT

- 2-partition achieves SI [[Csiszár-Narayan, 2004](#)]



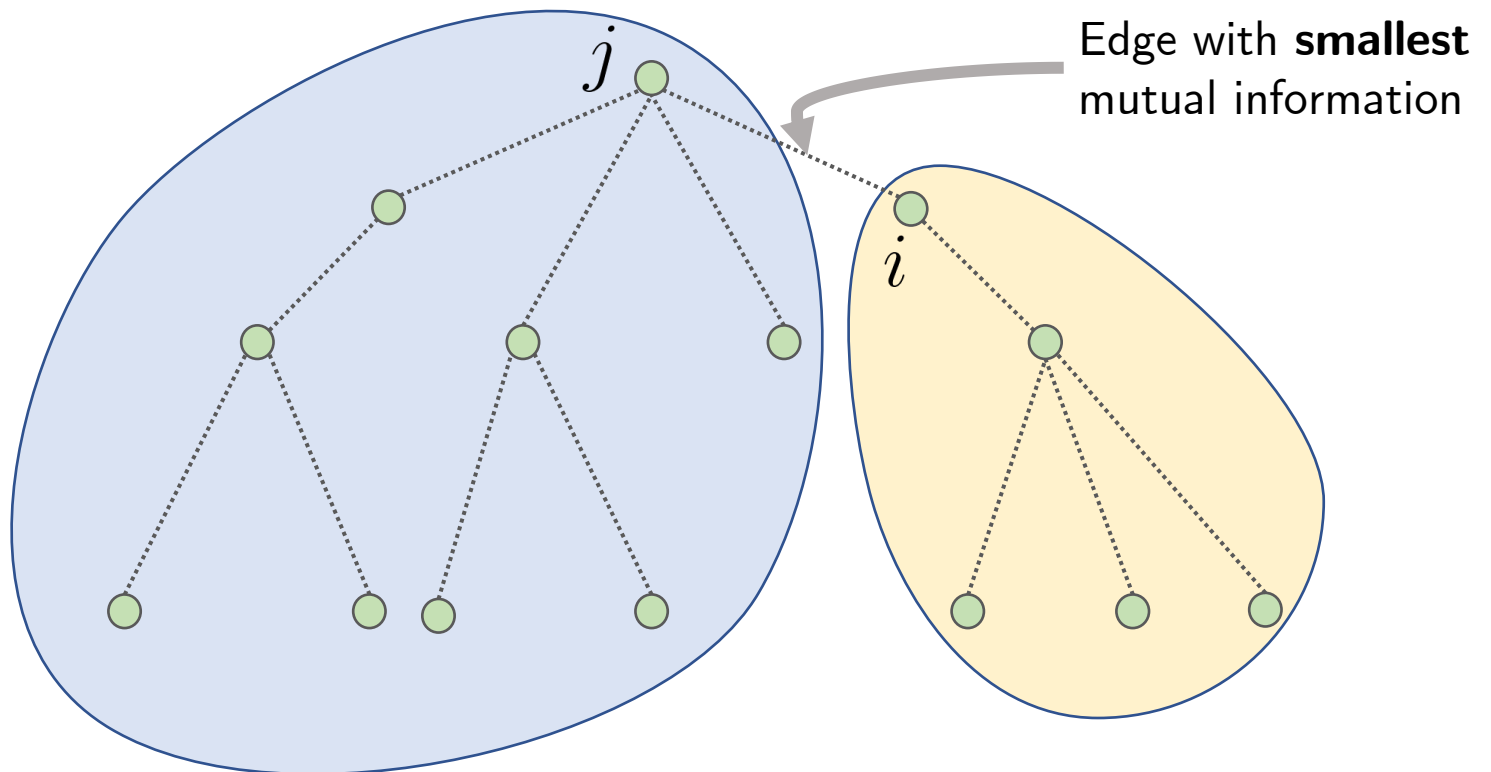
SI in an MCT

- 2-partition achieves SI [[Csiszár-Narayan, 2004](#)]



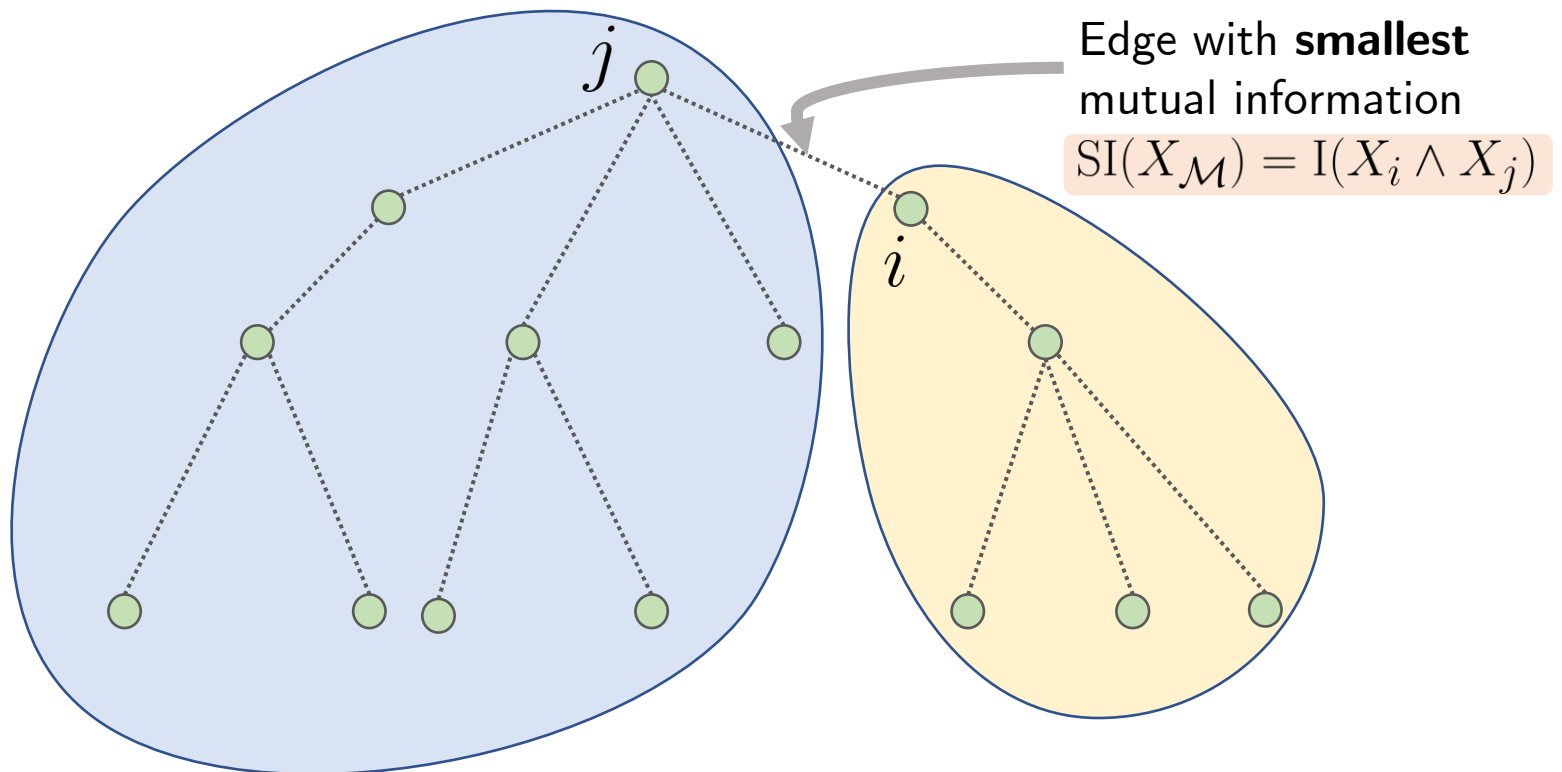
SI in an MCT

- 2-partition achieves SI [[Csiszár-Narayan, 2004](#)]



SI in an MCT

- 2-partition achieves SI [Csiszár-Narayan, 2004]



Shared information – simpler forms in special cases

$$\text{SI}(X_{\mathcal{M}}) = \min_{2 \leq k \leq m} \min_{\pi = (\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}})$$

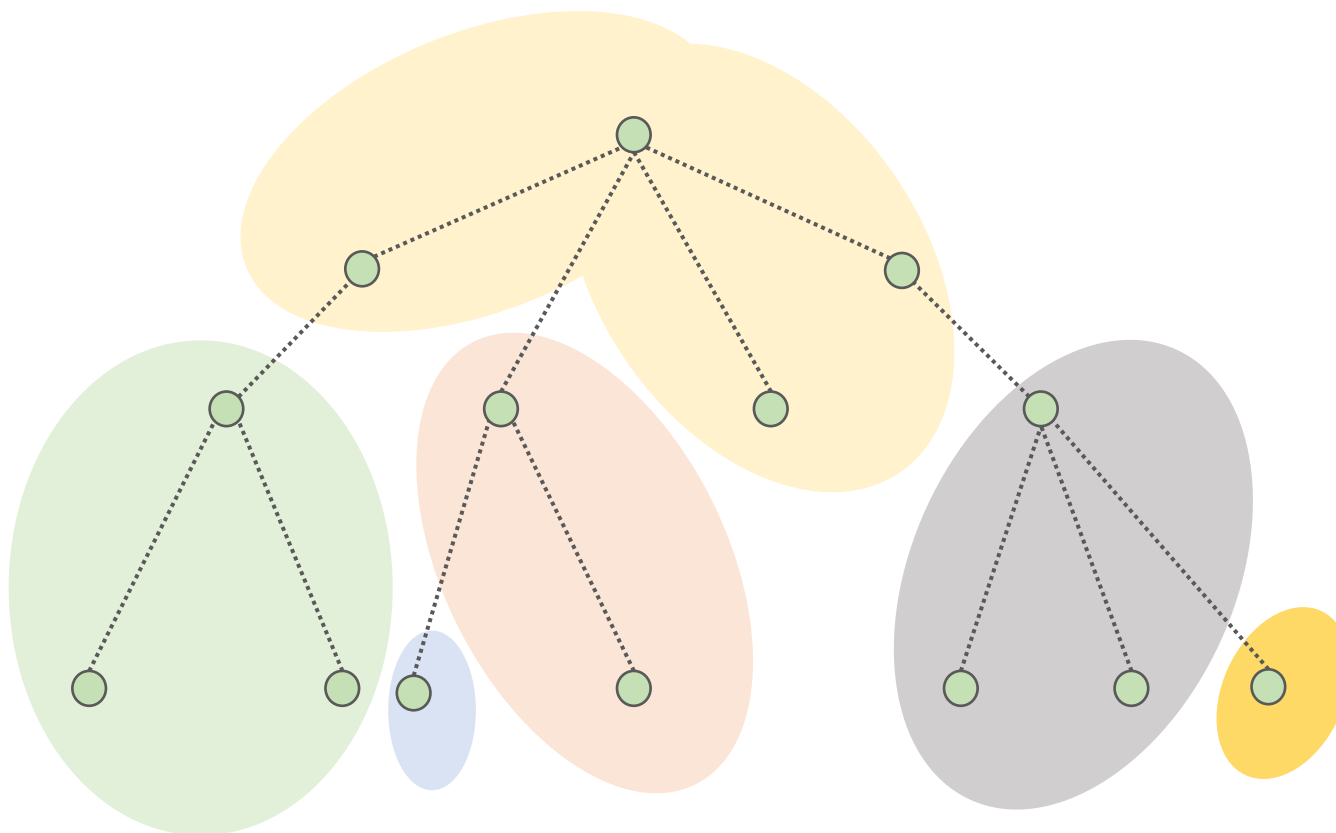
- Upper bound is clear*
- Choose that special partition

Shared information – simpler forms in special cases

$$\text{SI}(X_{\mathcal{M}}) = \min_{2 \leq k \leq m} \min_{\pi = (\pi_u, u=1, \dots, k)} \frac{1}{k-1} D(P_{X_{\mathcal{M}}} \parallel \prod_{u=1}^k P_{X_{\pi_u}})$$

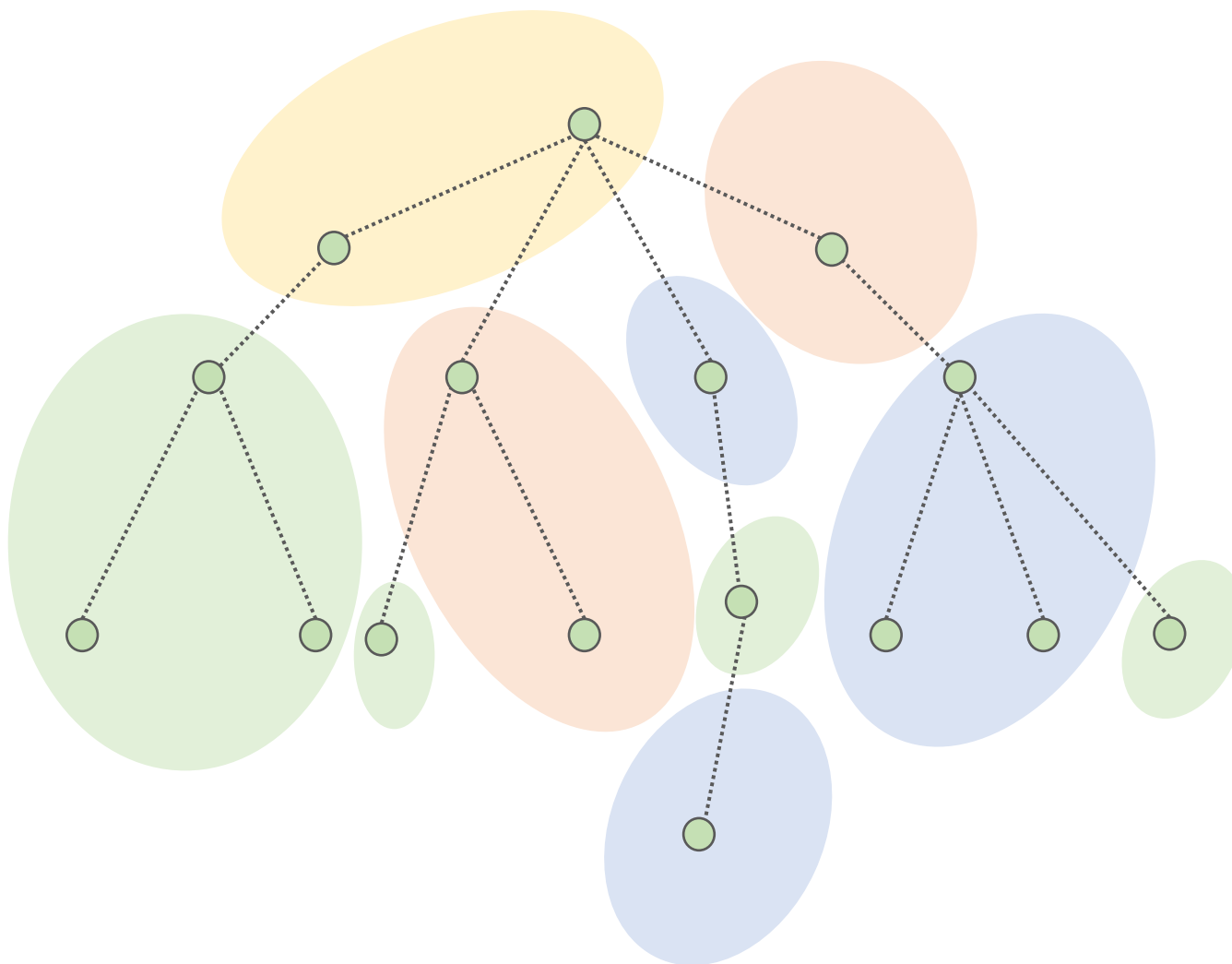
- Upper bound is clear*
- Choose that special partition
- Need: lower bound
 - Original proof from secret-key capacity [Csiszár-Narayan, 2004]
 - Different in spirit from [Chan-Bashabsheh-Zhou-Kaced-Liu, 2016]
 - Easy* when atoms of the partition are connected

Connected atoms

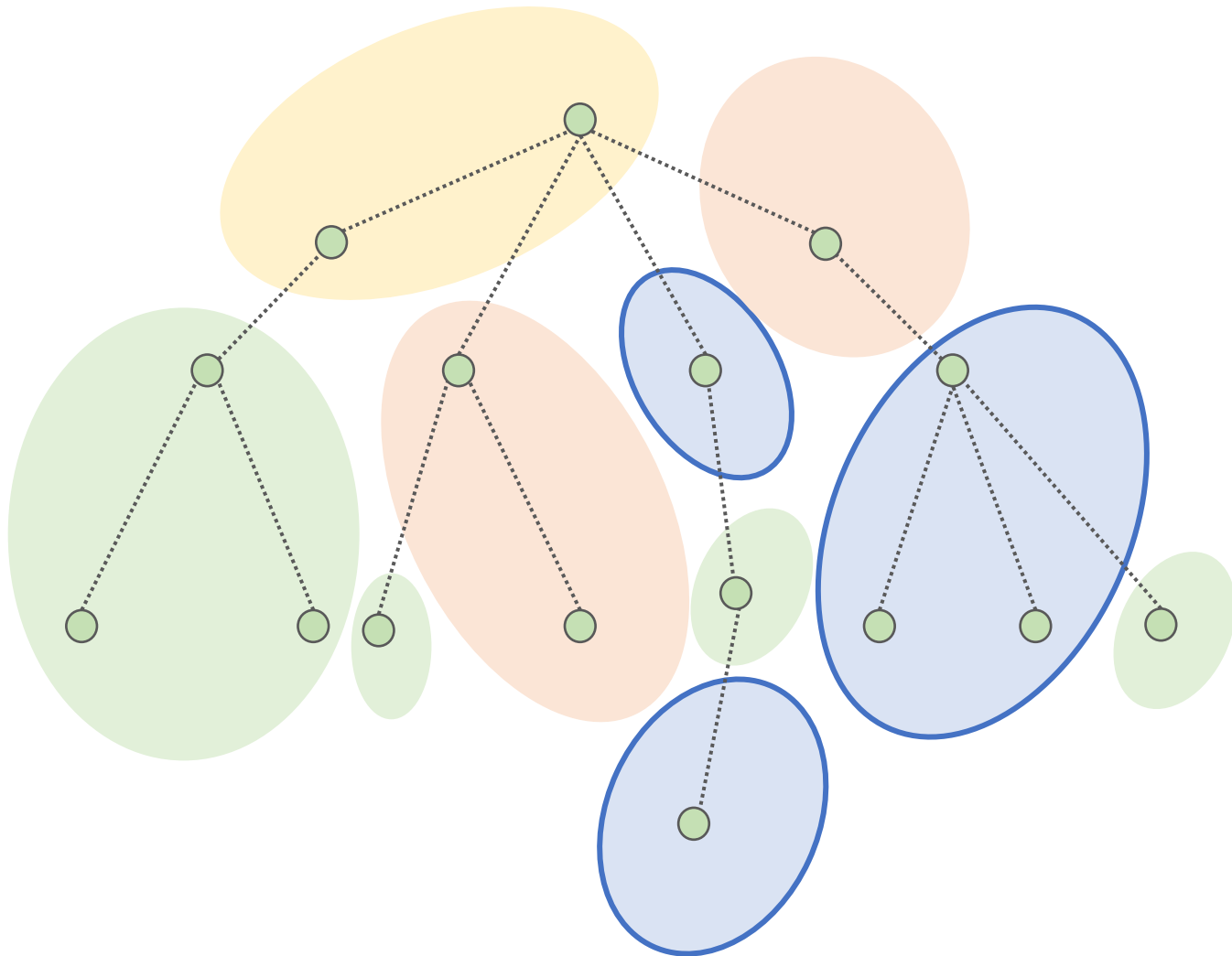


Idea: reduce nonconnected atom case to this case

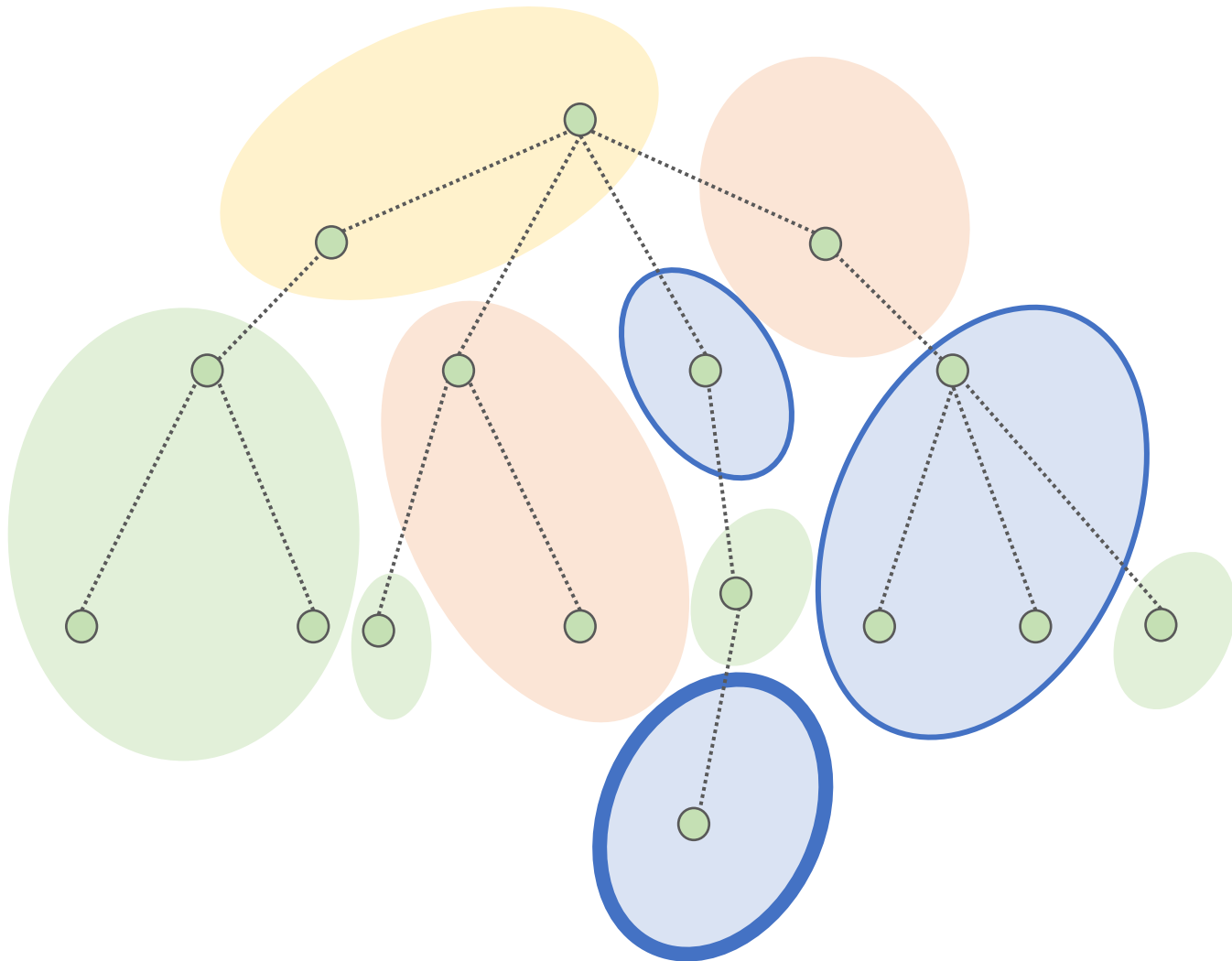
Idea: reduce nonconnected atom case to this case



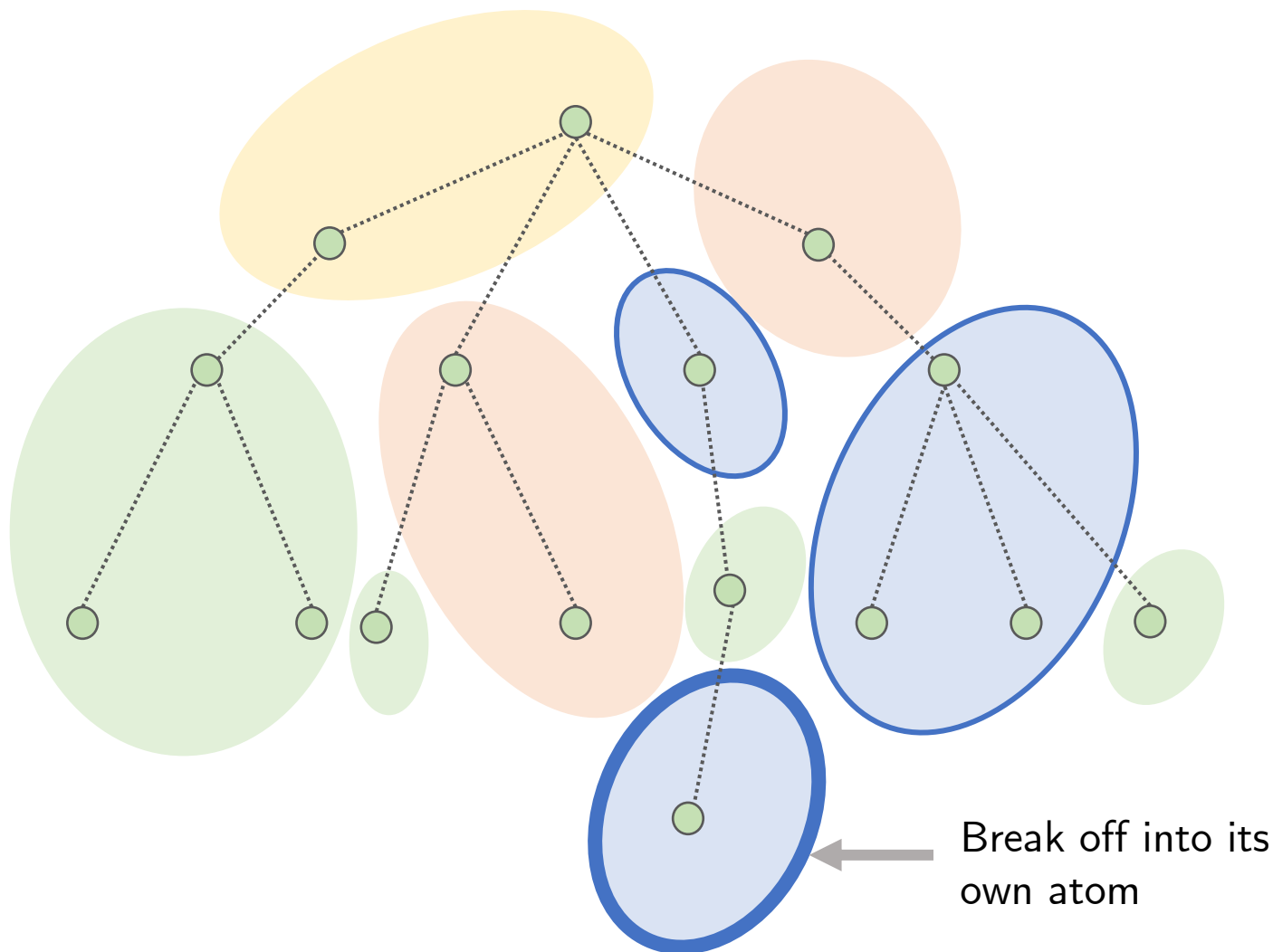
Idea: reduce nonconnected atom case to this case



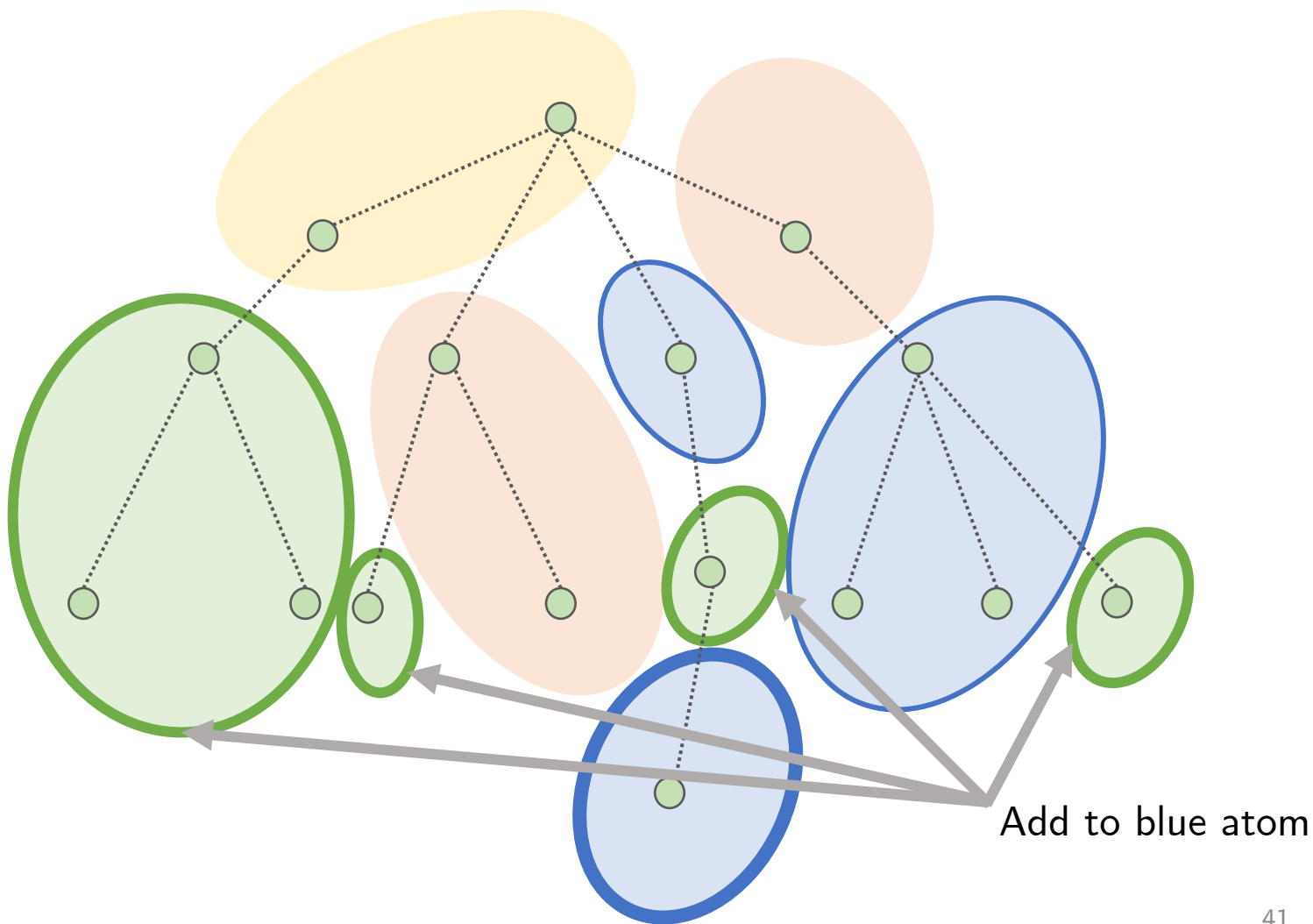
Idea: reduce nonconnected atom case to this case

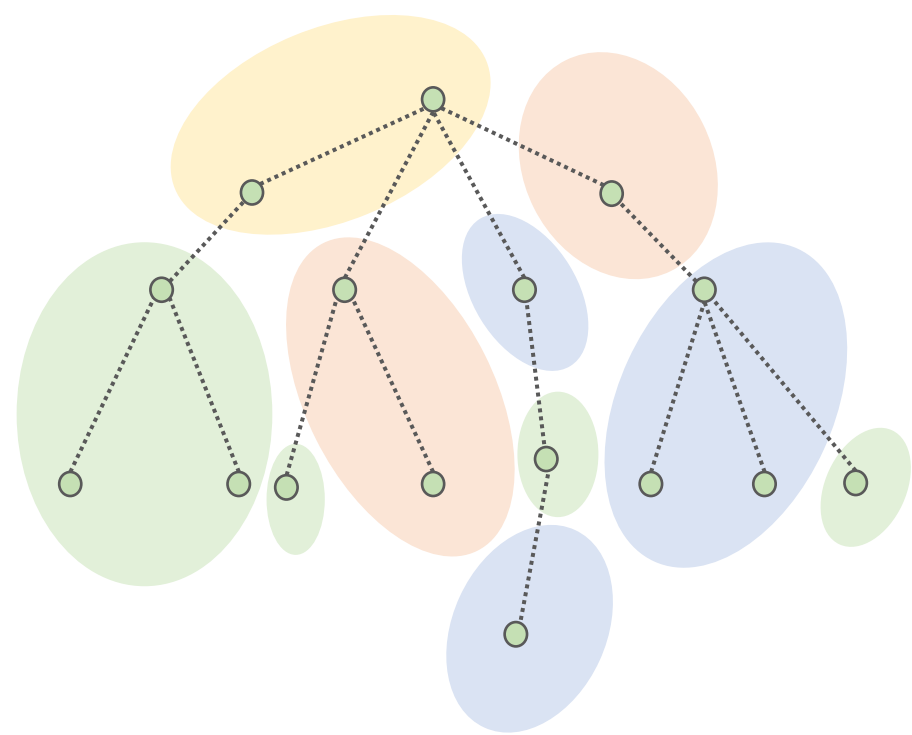


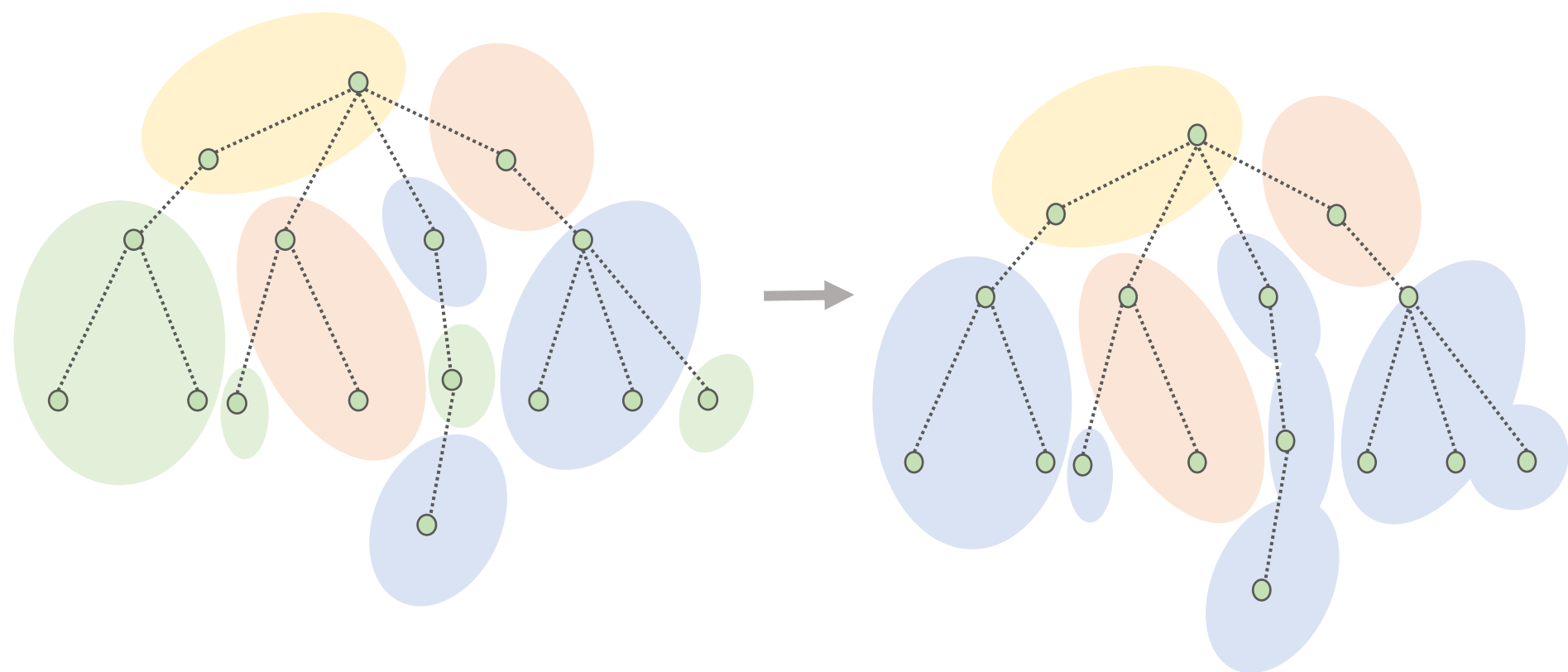
Idea: reduce nonconnected atom case to this case

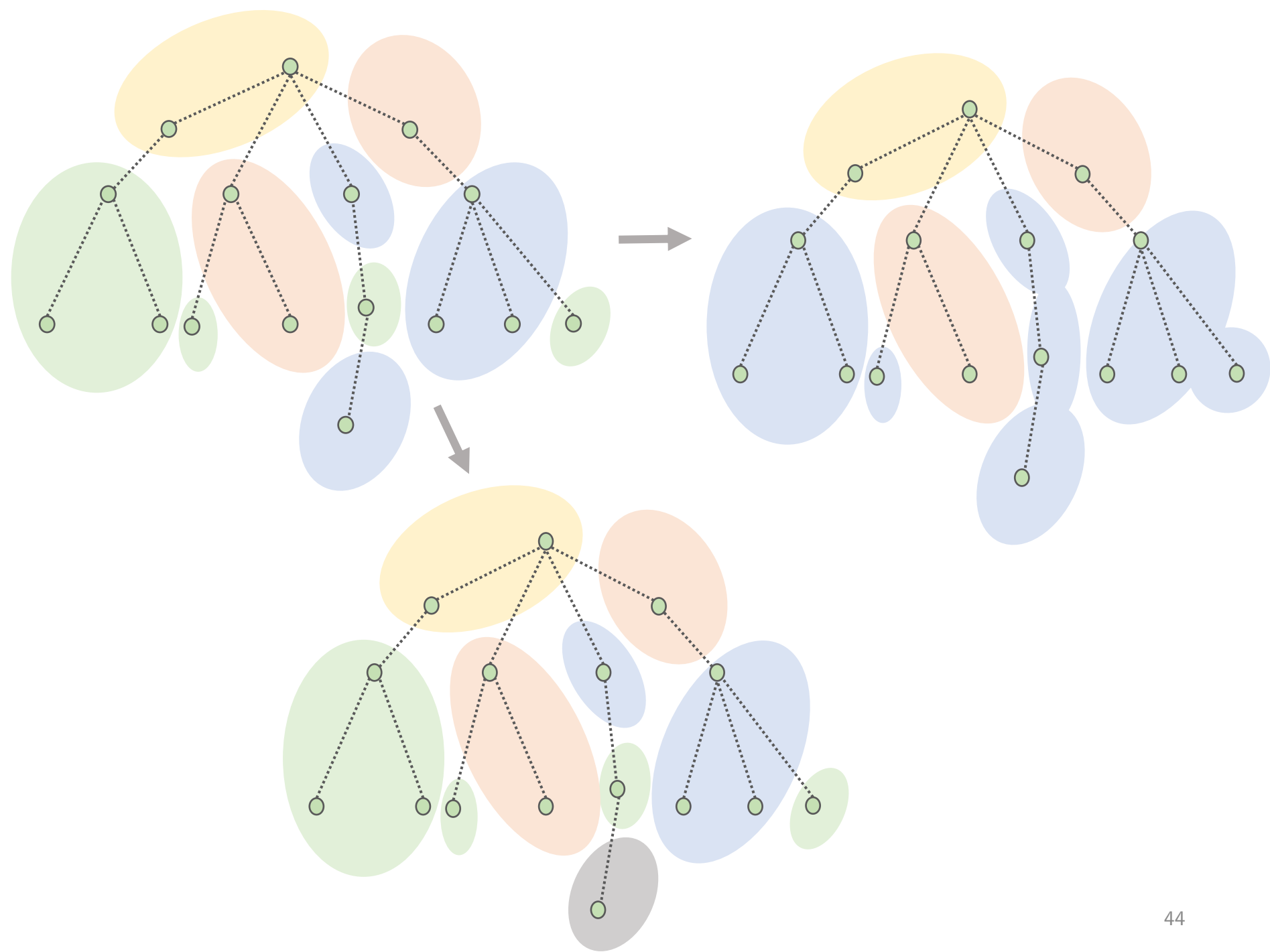


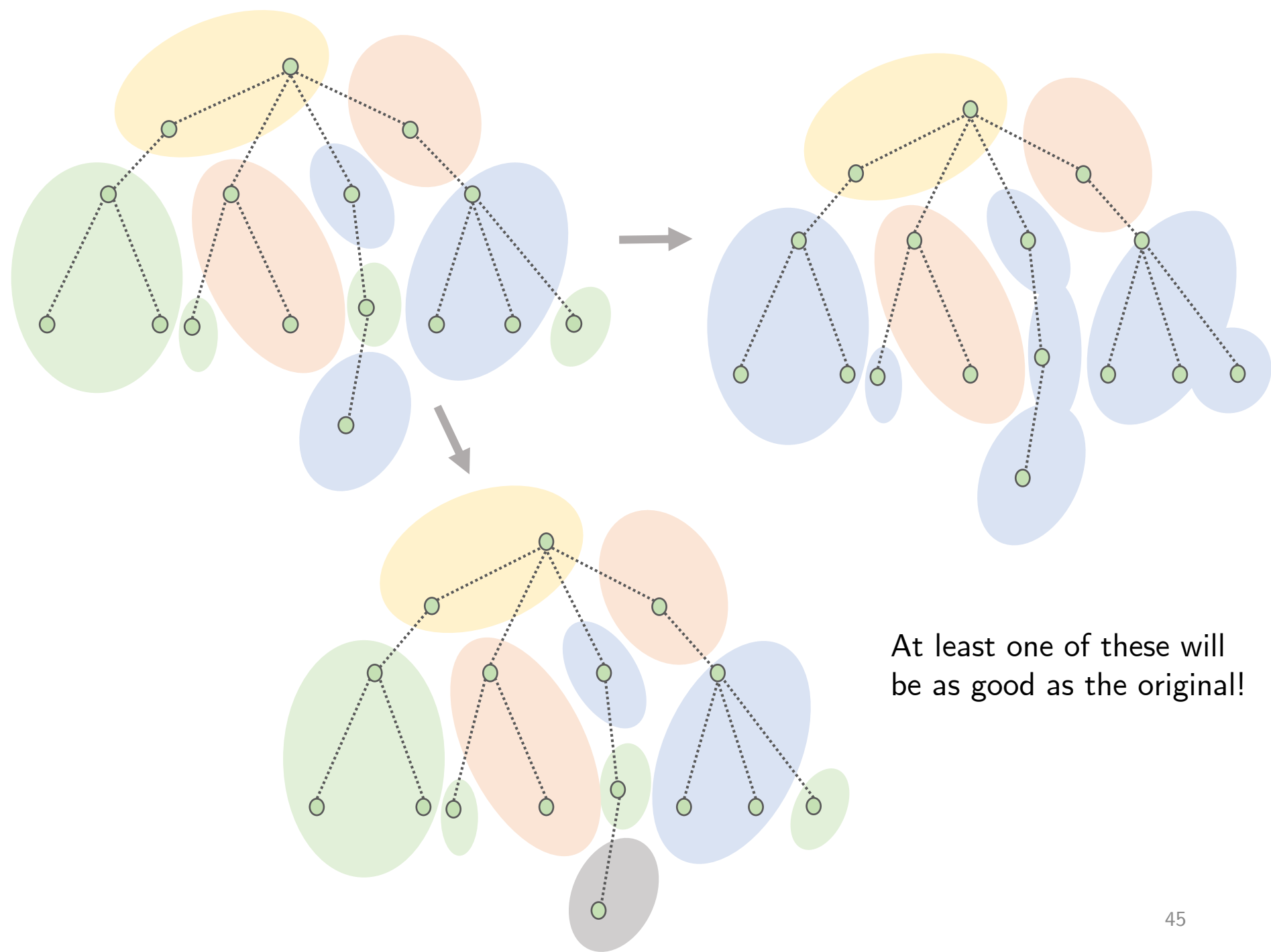
Idea: reduce nonconnected atom case to this case











#4: Estimating SI in an MCT

Estimating SI in an MCT

- Must find the edge with the lowest mutual information
 - Reduces to best arm identification [\[Audibert-Bubeck-Munos, 2010\]](#)

Estimating SI in an MCT

- Must find the edge with the lowest mutual information
 - Reduces to best arm identification [\[Audibert-Bubeck-Munos, 2010\]](#)
- Proposal: a **correlated** bandits algorithm

Estimating SI in an MCT

- Must find the edge with the lowest mutual information
 - Reduces to best arm identification [\[Audibert-Bubeck-Munos, 2010\]](#)
- Proposal: a **correlated** bandits algorithm
 - For every edge, sample both random variables at the same time
 - **Uniform sampling:** each pair is sampled equally often

Estimating SI in an MCT

- Must find the edge with the lowest mutual information
 - Reduces to best arm identification [\[Audibert-Bubeck-Munos, 2010\]](#)
- Proposal: a **correlated** bandits algorithm
 - For every edge, sample both random variables at the same time
 - **Uniform sampling:** each pair is sampled equally often
- **Key challenge:** estimators for mutual information are always **biased**.

MI estimation

- The **empirical mutual information** [\[Goppa, 1975\]](#)

$$I_{\text{EMI}}^{(n)}(\mathbf{x} \wedge \mathbf{y}) = H(P_{\mathbf{x}}^{(n)}) + H(P_{\mathbf{y}}^{(n)}) - H(P_{\mathbf{xy}}^{(n)})$$

MI estimation

- The empirical mutual information [Goppa, 1975]

$$I_{\text{EMI}}^{(n)}(\mathbf{x} \wedge \mathbf{y}) = H(P_{\mathbf{x}}^{(n)}) + H(P_{\mathbf{y}}^{(n)}) - H(P_{\mathbf{xy}}^{(n)})$$

- [Paninski, 2003]

$$-\log \left(1 + \frac{|\mathcal{X}| - 1}{n} \right) \left(1 + \frac{|\mathcal{Y}| - 1}{n} \right) \leq \text{Bias}(I_{\text{EMI}}^{(n)}(\mathbf{X} \wedge \mathbf{Y})) \leq \log \left(1 + \frac{|\mathcal{X}| |\mathcal{Y}| - 1}{n} \right)$$

MI estimation

- The empirical mutual information [Goppa, 1975]

$$I_{\text{EMI}}^{(n)}(\mathbf{x} \wedge \mathbf{y}) = H(P_{\mathbf{x}}^{(n)}) + H(P_{\mathbf{y}}^{(n)}) - H(P_{\mathbf{xy}}^{(n)})$$

- [Paninski, 2003]

$$-\log \left(1 + \frac{|\mathcal{X}| - 1}{n} \right) \left(1 + \frac{|\mathcal{Y}| - 1}{n} \right) \leq \text{Bias}(I_{\text{EMI}}^{(n)}(\mathbf{X} \wedge \mathbf{Y})) \leq \log \left(1 + \frac{|\mathcal{X}| |\mathcal{Y}| - 1}{n} \right)$$

- Using McDiarmid's inequality, similar to [Antos-Kontoyiannis, 2001]

$$P_{XY} \left(I_{\text{EMI}}^{(n)}(\mathbf{X} \wedge \mathbf{Y}) - \mathbb{E}_{P_{XY}} \left[I_{\text{EMI}}^{(n)}(\mathbf{X} \wedge \mathbf{Y}) \right] \geq \epsilon \right) \leq \exp \left(-\frac{2n\epsilon^2}{36 \log^2 n} \right)$$

Misidentification probability

$$P_{X_{\mathcal{M}}} \left(\hat{e}_N(X_{\mathcal{M}}^N) \neq (\bar{i}, \bar{j}) \right) \leq 2 |\mathcal{E}| \exp \left(\frac{-(N/|\mathcal{E}|)\Delta_1^2}{648 \log^2(N/|\mathcal{E}|)} \right)$$

if

$$N > |\mathcal{E}| \max \left\{ \frac{|\mathcal{X}|^2 - 1}{2^{\Delta_1/3} - 1}, \frac{|\mathcal{X}| - 1}{2^{\Delta_1/6} - 1} \right\},$$

Misidentification probability

Number of edges

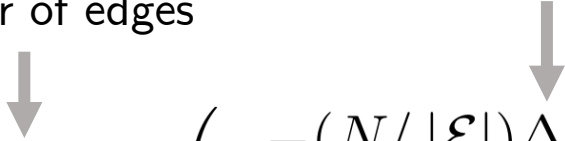


$$P_{X_{\mathcal{M}}} \left(\hat{e}_N(X_{\mathcal{M}}^N) \neq (\bar{i}, \bar{j}) \right) \leq 2 |\mathcal{E}| \exp \left(\frac{-(N/|\mathcal{E}|)\Delta_1^2}{648 \log^2(N/|\mathcal{E}|)} \right)$$

if

$$N > |\mathcal{E}| \max \left\{ \frac{|\mathcal{X}|^2 - 1}{2^{\Delta_1/3} - 1}, \frac{|\mathcal{X}| - 1}{2^{\Delta_1/6} - 1} \right\},$$

Misidentification probability

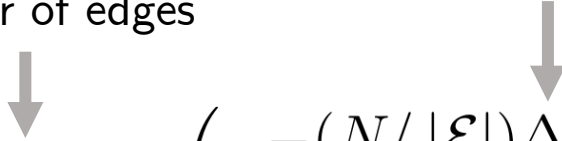
Number of edges Difference between the smallest and second-smallest MIs


$$P_{X_{\mathcal{M}}} \left(\hat{e}_N(X_{\mathcal{M}}^N) \neq (\bar{i}, \bar{j}) \right) \leq 2 |\mathcal{E}| \exp \left(\frac{-(N/|\mathcal{E}|)\Delta_1^2}{648 \log^2(N/|\mathcal{E}|)} \right)$$

if

$$N > |\mathcal{E}| \max \left\{ \frac{|\mathcal{X}|^2 - 1}{2^{\Delta_1/3} - 1}, \frac{|\mathcal{X}| - 1}{2^{\Delta_1/6} - 1} \right\},$$


Misidentification probability

Number of edges Difference between the smallest and second-smallest MIs


$$P_{X_{\mathcal{M}}} \left(\hat{e}_N(X_{\mathcal{M}}^N) \neq (\bar{i}, \bar{j}) \right) \leq 2 |\mathcal{E}| \exp \left(\frac{-(N/|\mathcal{E}|)\Delta_1^2}{648 \log^2(N/|\mathcal{E}|)} \right)$$

if

$$N > |\mathcal{E}| \max \left\{ \frac{|\mathcal{X}|^2 - 1}{2^{\Delta_1/3} - 1}, \frac{|\mathcal{X}| - 1}{2^{\Delta_1/6} - 1} \right\},$$


 Reduces bias below Δ_1

Key takeaways

- New proof for SI in an MCT
- Best-arm identification using biased estimators

Directions for future work

- **Better estimators** lead to better bounds
- **Lower bounds** for best-arm identification using biased estimators