

Khai Thác Dữ Liệu Đồ Thị

PHÁT HIỆN CỘNG ĐỒNG

Giảng viên: Lê Ngọc Thành

Email: lnthanh@fit.hcmus.edu.vn



fit@hcmus

Nội dung

- **Khái niệm**
- Tối ưu cộng đồng
- Các phương pháp phát hiện cộng đồng
 - Phương pháp dựa trên lát cắt tối thiểu
 - Phương pháp dựa trên trung gian
 - Phương pháp dựa trên RandomWalk



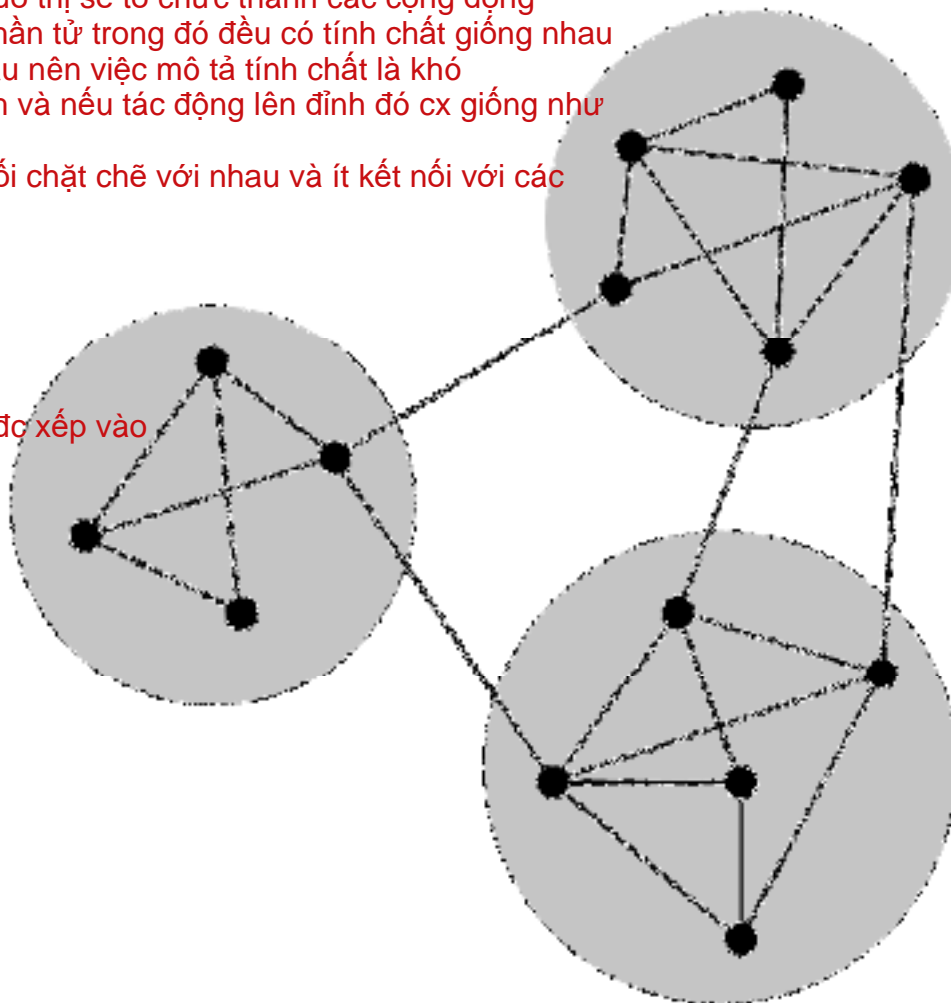
Cộng đồng

- **Cộng đồng** (community) là tập các đỉnh mà mỗi đỉnh có nhiều kết nối bên trong hơn kết nối ra ngoài.

- Các đồ thị trong thực tế ko phải là ngẫu nhiên
- Cta kỳ vọng rằng các đỉnh trong 1 đồ thị sẽ tổ chức thành các cộng đồng
- Phát hiện ra nhóm (CD) vì tất cả phần tử trong đó đều có tính chất giống nhau
- Mỗi nhóm sẽ có tính chất khác nhau nên việc mô tả tính chất là khó
- Gom các ptu trong CD thành 1 đỉnh và nếu tác động lên đỉnh đó cx giống như tác động trên nhóm đó
- Các đỉnh trong CD sẽ có mối kết nối chặt chẽ với nhau và ít kết nối với các cộng đồng khác

- Có 2 loại cộng đồng là:
 - + Cộng đồng giao nhau (*)
 - + Cộng đồng ko giao nhau

=> Bài toán gom nhóm (Clustering) đc xếp vào Unsupervise learning vì ko có nhãn.
Nếu đã có nhãn thì nó thuộc về bài toán phân lớp (Classification)



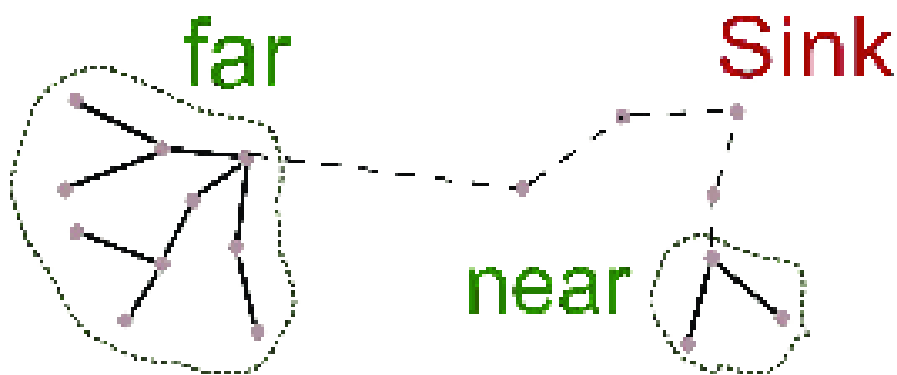
Nội dung

- Khái niệm
- **Tối ưu cộng đồng**
- Các phương pháp phát hiện cộng đồng
 - Phương pháp dựa trên lát cắt tối thiểu
 - Phương pháp dựa trên trung gian
 - Phương pháp dựa trên RandomWalk

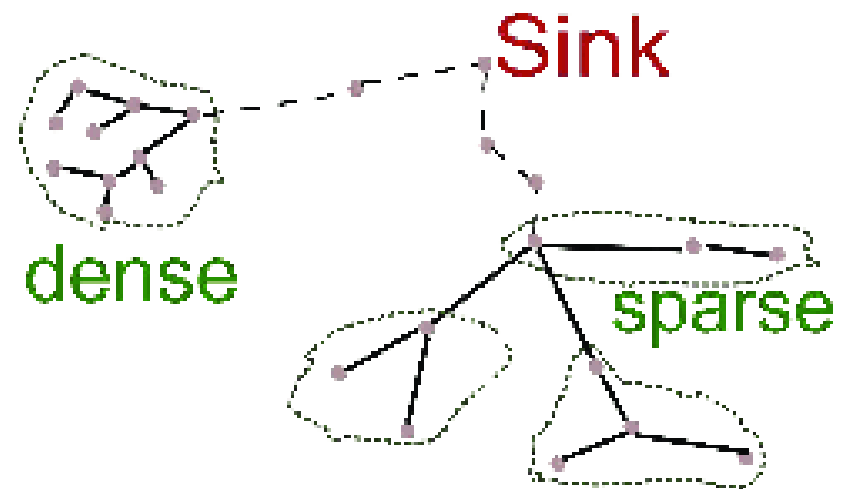


Tối ưu cộng đồng

- Ngoài các độ đo chất lượng và khoảng cách, người ta có thêm 2 cách đo để đánh giá mật độ:
 - **Mật độ trong nhóm** (intracuster density): càng lớn càng tốt
 - **Mật độ ngoài nhóm** (intercluster density): càng nhỏ càng tốt



(a) distance



(b) density

Mật độ trong nhóm => THI

- **Mật độ trong nhóm** là tỉ số giữa số cạnh bên trong nhóm và số cạnh có thể tối đa trong nhóm.

N_c là số đỉnh trong nhóm

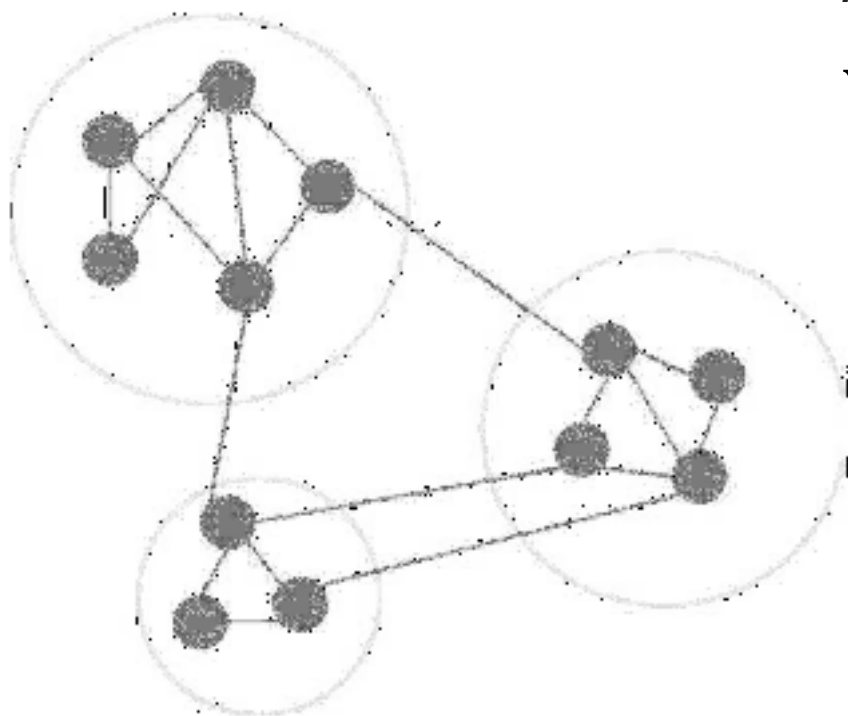
$$\delta(C) = \frac{\text{số cạnh trong nhóm}}{N_c(N_c - 1)/2}$$

Ví dụ:

$$\delta(C_1) = \frac{7}{10} = 0.7$$

$$\delta(C_2) = \frac{4}{6} = 0.75 \quad \Rightarrow 5/6$$

$$\delta(C_3) = \frac{3}{3} = 1.0$$



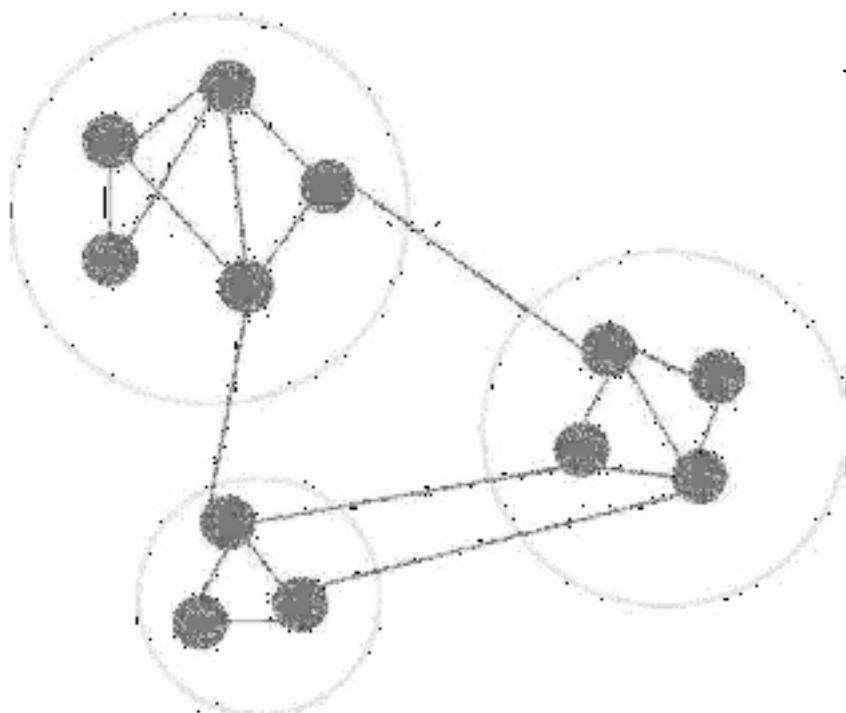
Mật độ ngoài nhóm

- **Mật độ ngoài nhóm** là tỉ số giữa số cạnh ngoài nhóm và số cạnh ngoài nhóm có thể.

$N_c(N - N_c)$: số đỉnh trong nhóm * số đỉnh ngoài nhóm

$$\epsilon(C) = \frac{\text{số cạnh ngoài nhóm}}{N_c(N - N_c)}$$

Ví dụ:



$$\begin{aligned}\epsilon(C_1) &= \frac{2}{35} \\ \epsilon(C_2) &= \frac{3}{32} \\ \epsilon(C_3) &= \frac{3}{27}\end{aligned}$$

$2 / (5 * (3 + 4)) = 2 / 35$



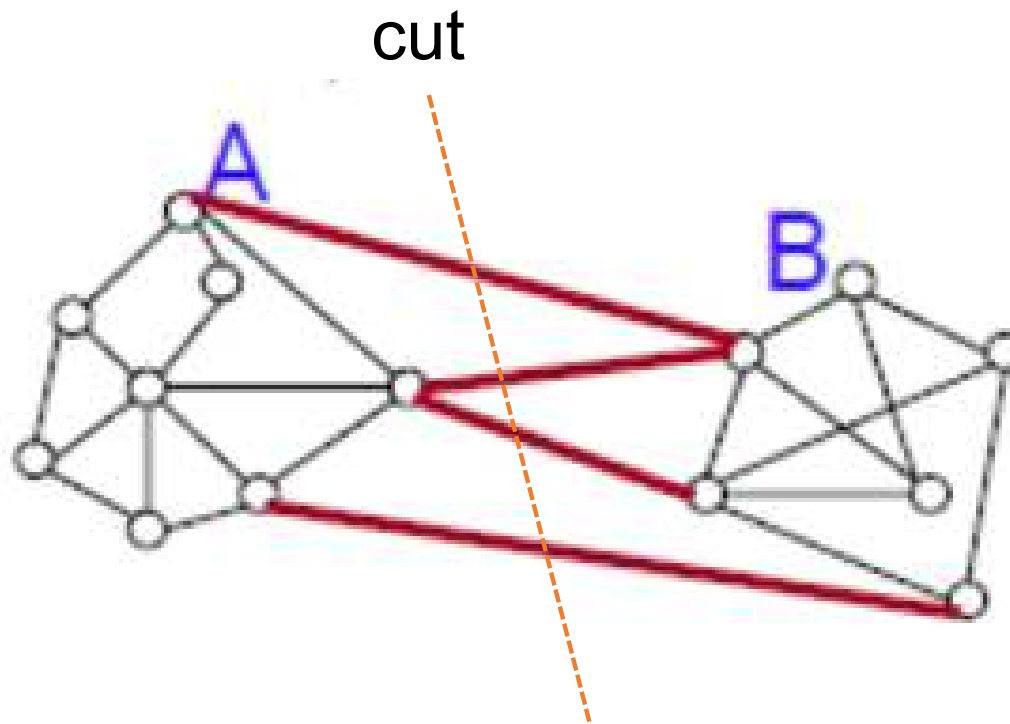
Nội dung

- Khái niệm
- Tối ưu cộng đồng
- **Các phương pháp phát hiện cộng đồng**
 - Phương pháp dựa trên lát cắt tối thiểu
 - Phương pháp dựa trên trung gian
 - Phương pháp dựa trên RandomWalk



Lát cắt tối thiểu

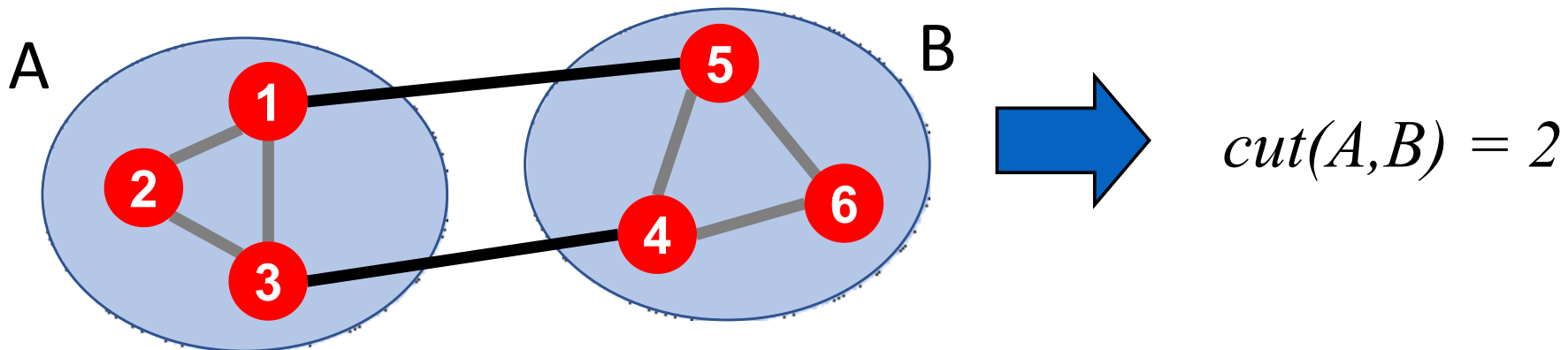
- Mục tiêu của lát cắt tối thiểu là tìm tập cạnh ít nhất mà chặn luồng từ nguồn S đến T.
 - Kích thước cắt là tổng trọng lượng các cạnh đó



Lát cắt tối thiểu

- Lát cắt:

$$\text{cut}(A, B) = \sum_{i \in A, j \in B} w_{ij}$$

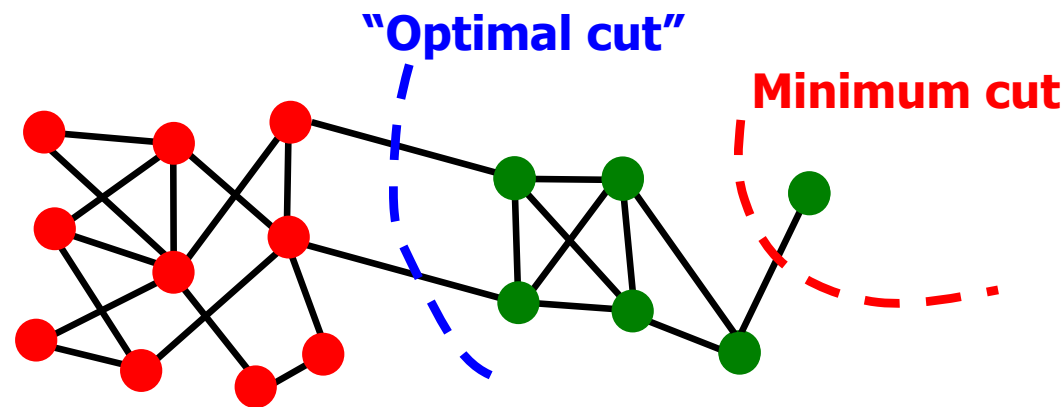


Phát hiện cộng đồng với lát cắt tối thiểu

- Mục tiêu: lát cắt tối thiểu

$$\arg \min_{A,B} \text{cut}(A,B)$$

- Vấn đề:



Vấn đề xảy ra do không xem xét kết nối bên trong

Chuẩn hóa lát cắt

- Lát cắt được chuẩn hóa phụ thuộc vào mật độ trong từng nhóm

$$ncut(A, B) = \frac{cut(A, B)}{vol(A)} + \frac{cut(A, B)}{vol(B)}$$

với $vol(A)$: tổng trọng số các cạnh với ít nhất một đầu ở trong A

$$vol(A) = \sum_{i \in A} k_i$$

Nội dung

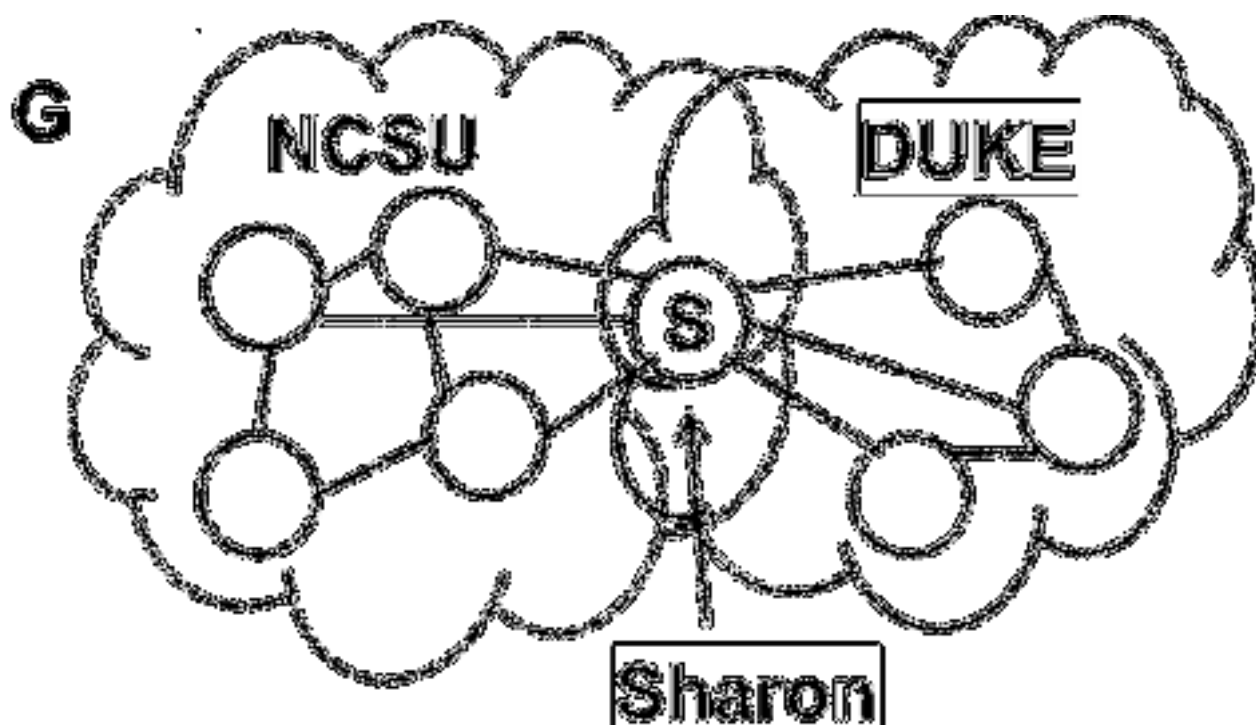
- Khái niệm
- Tối ưu cộng đồng
- **Các phương pháp phát hiện cộng đồng**
 - Phương pháp dựa trên lát cắt tối thiểu
 - **Phương pháp dựa trên trung gian** => Tập trung
 - Phương pháp dựa trên RandomWalk

Phát hiện cộng đồng dựa trên tính trung gian

- Phát hiện cộng đồng dựa trên tính trung gian thực hiện quá trình xác định cộng đồng bằng cách lần lượt bỏ đi đỉnh/cạnh có tính trung gian cao.
- Có hai loại:
 - Dựa trên tính trung gian đỉnh
 - Dựa trên tính trung gian cạnh.

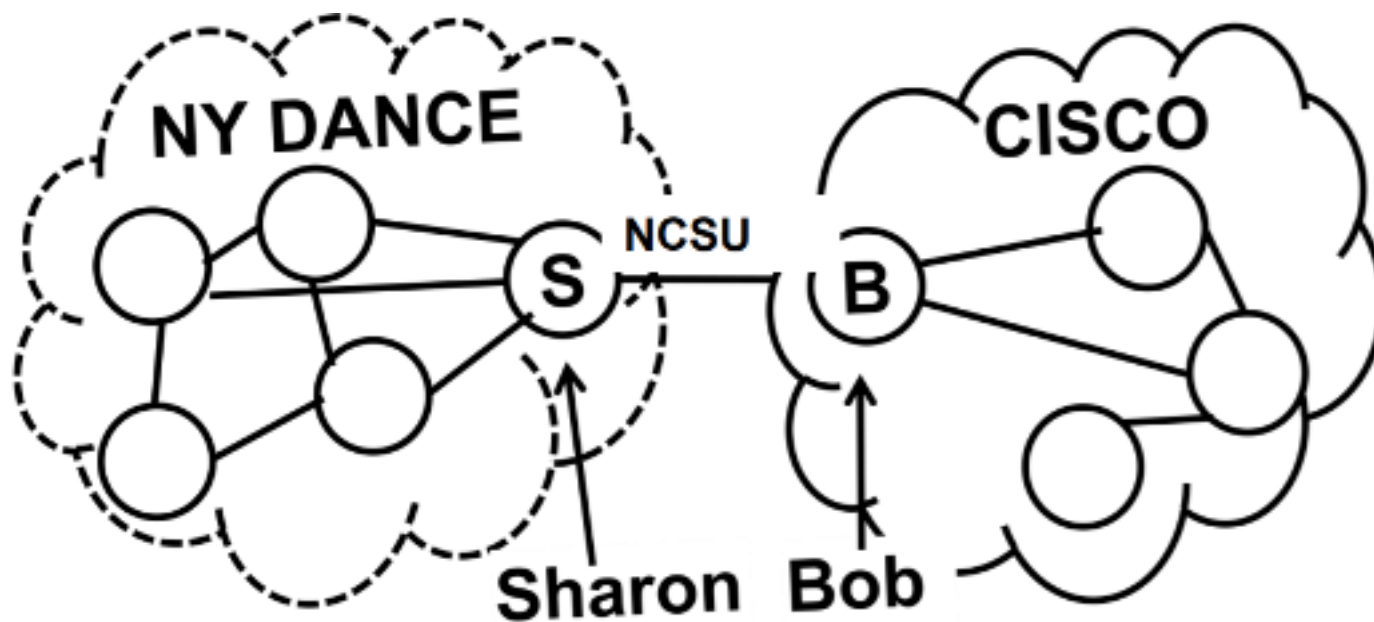
Trung gian đỉnh

- **Trung gian đỉnh** được tính dựa trên số đường đi ngắn nhất trong đồ thị mà phải đi qua đỉnh cho trước.



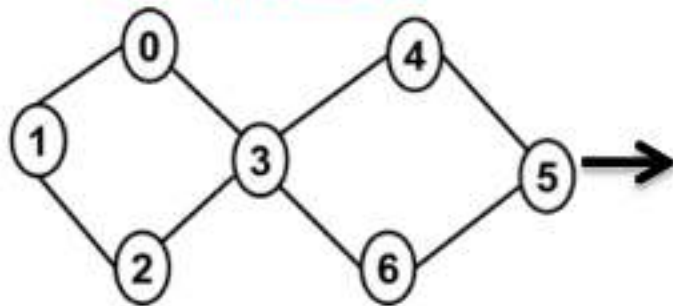
Trung gian cạnh

- **Trung gian cạnh** được đánh giá dựa trên số đường đi ngắn nhất phải đi qua cạnh cho trước.

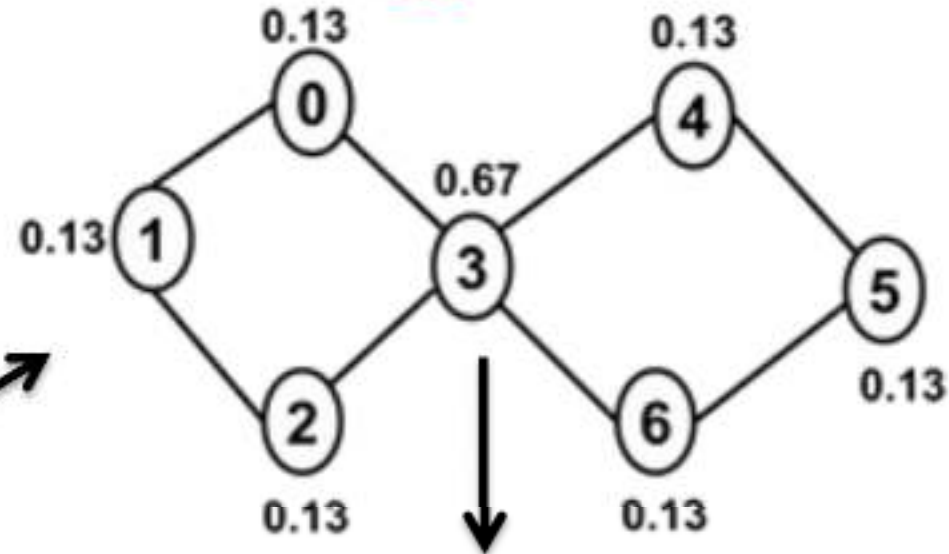


Thuật toán trung gian đỉnh

Đồ thị ban đầu

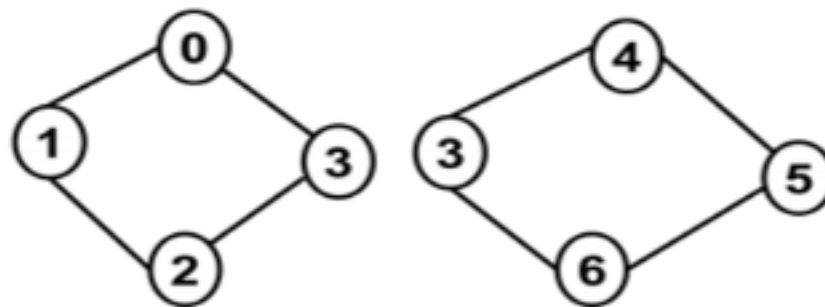


Trung gian đỉnh



Lặp lại cho đến khi trung gian đỉnh nhỏ hơn ngưỡng cho trước

Ngắt đồ thị tại đỉnh này



Chọn đỉnh có giá trị trung gian đỉnh lớn nhất

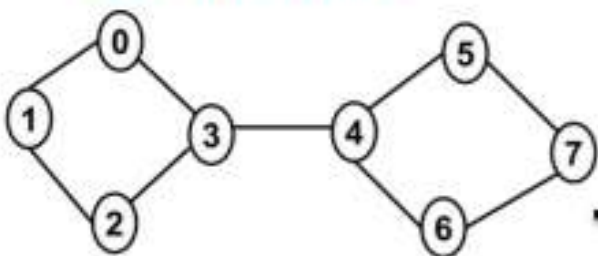
Thuật toán trung gian cạnh

=> Thuật toán Girvan Newman

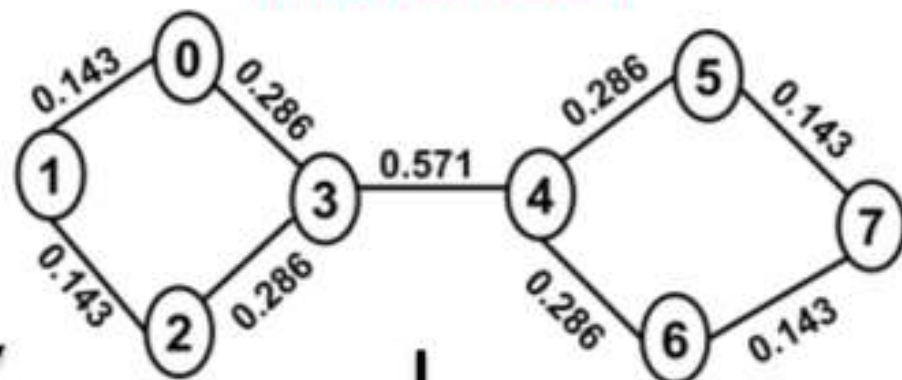
=> Tập trung vào trung gian cạnh

=> Tập trung vào thuật toán Girvan Newman

Đồ thị ban đầu

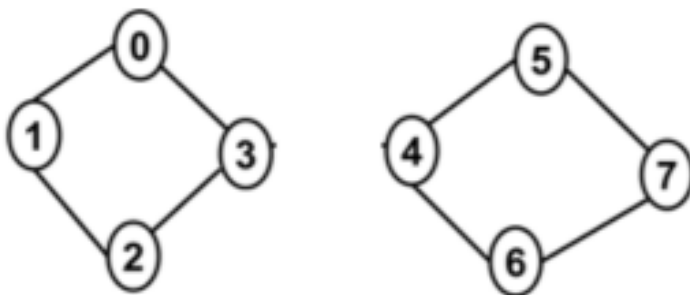


Trung gian cạnh



Lặp lại cho đến khi độ
cạnh trung gian nhỏ
hơn ngưỡng cho trước

Ngắt đồ thị tại
cạnh này



Chọn cạnh có
giá trị trung gian
đỉnh lớn nhất

Cho đồ thị cho trước và ngưỡng modularity
và áp dụng thuật toán

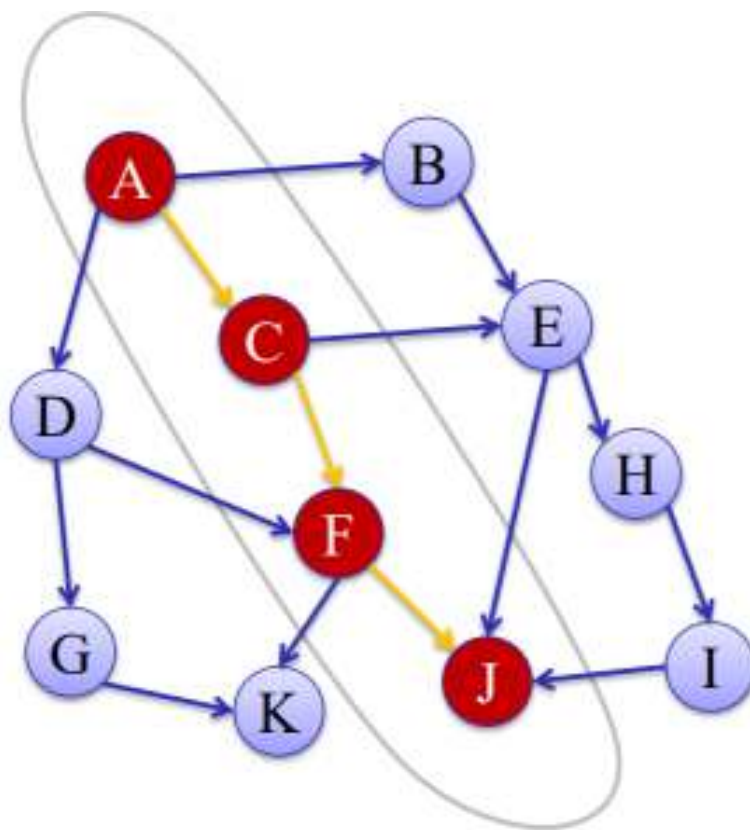
=> Hội tụ khi mà các nhóm chỉ có 1 đỉnh => Dg cắt đến từ tính Modularity, đến 1 ngưỡng Modularity thì cắt

Nội dung

- Khái niệm
- Tối ưu cộng đồng
- **Các phương pháp phát hiện cộng đồng**
 - Phương pháp dựa trên lát cắt tối thiểu
 - Phương pháp dựa trên trung gian
 - **Phương pháp dựa trên RandomWalk**

Random Walk

- Cho một đồ thị, từ một đỉnh bắt đầu, ta chọn ngẫu nhiên một đỉnh để đi tiếp. Sau t lần lặp lại như vậy, ta sẽ có một **đường đi ngẫu nhiên** (random walk) kích thước t .

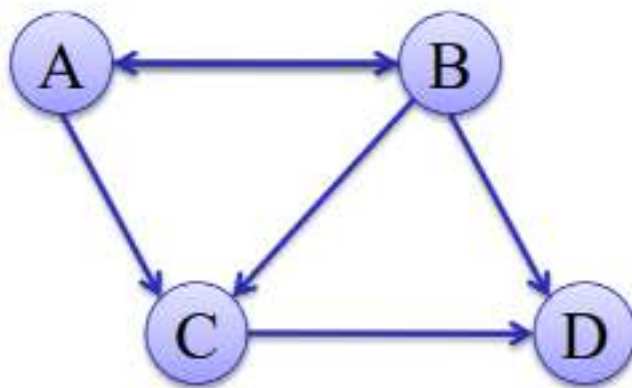


Random Walk

- Xác suất chọn cạnh có thể đánh giá dựa trên độ tương quan giữa hai văn bản và được chuẩn hóa.

Transition Matrix

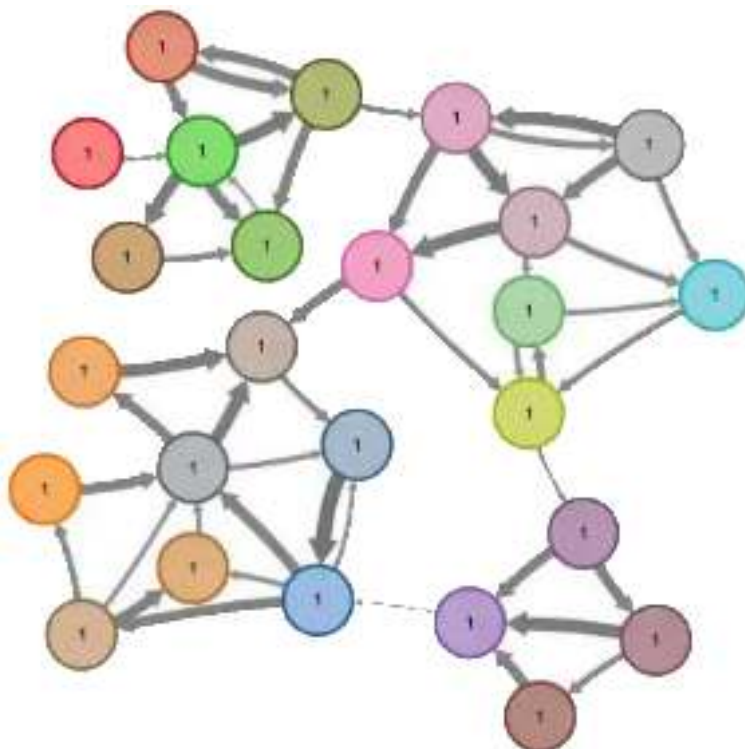
$$\begin{bmatrix} 0 & 1/2 & 1/2 & 0 \\ 1/3 & 0 & 1/3 & 1/3 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$



Phát hiện cộng đồng dựa trên random walk

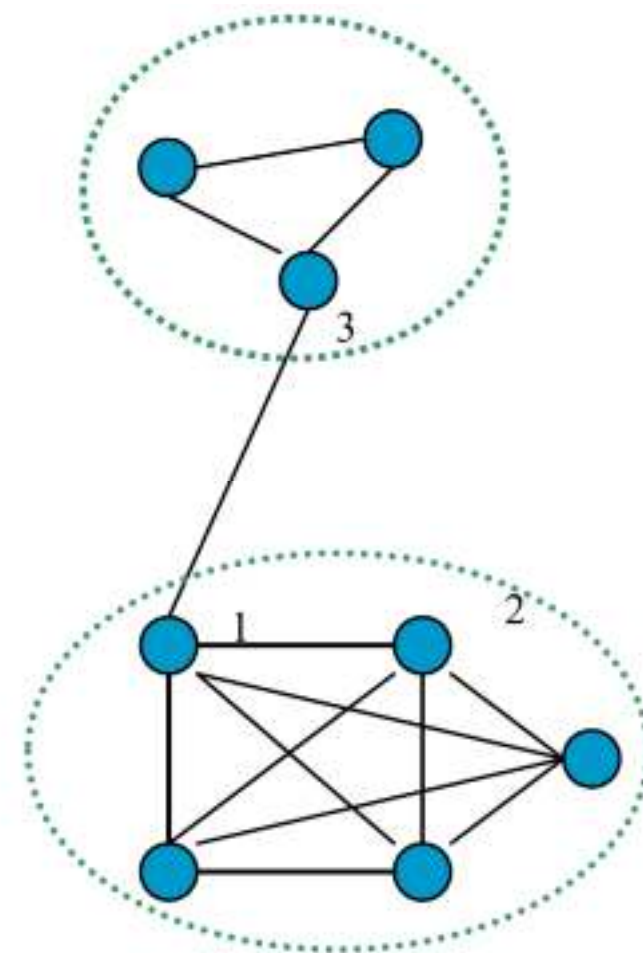
- Việc phát hiện cộng đồng dựa trên random walk được thực hiện dựa trên nhận định:

Một đường đi ngẫu nhiên bắt đầu ở một đỉnh nhiều khả năng sẽ di chuyển trong cộng đồng hơn là di chuyển trong cộng đồng khác.



Ví dụ

Node	Prob. Next Step within cluster	Prob. Next Step between clusters
1	80%	20%
2	100%	0%
3	67%	33%

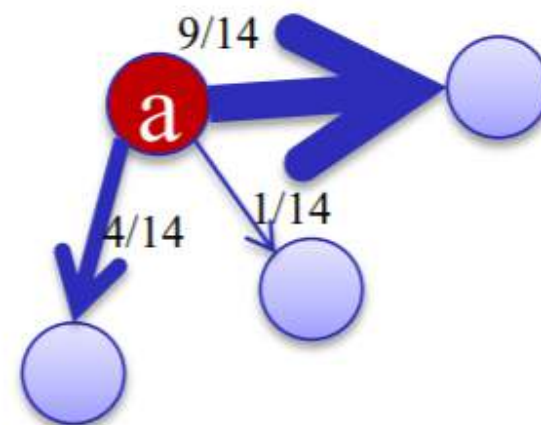
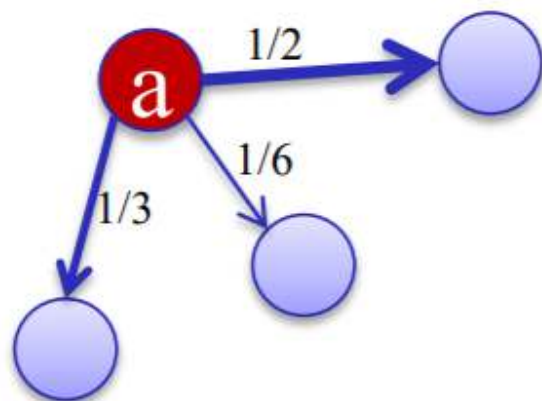


Ý tưởng thuật toán

- Điều chỉnh trọng số để sau bước đi ngẫu nhiên kích thước cho thước, khả năng di chuyển trong nhóm sẽ cao.
- Trọng số được điều chỉnh để mà:
 - Các láng giềng mạnh sẽ càng mạnh lên
 - Các láng giềng yếu sẽ càng yếu đi
 - Tiến trình này được gọi là sự làm phồng (inflation)



Lạm phát hó □



$$\begin{bmatrix} 0 \\ 1/2 \\ 0 \\ 1/6 \\ 1/3 \end{bmatrix}$$

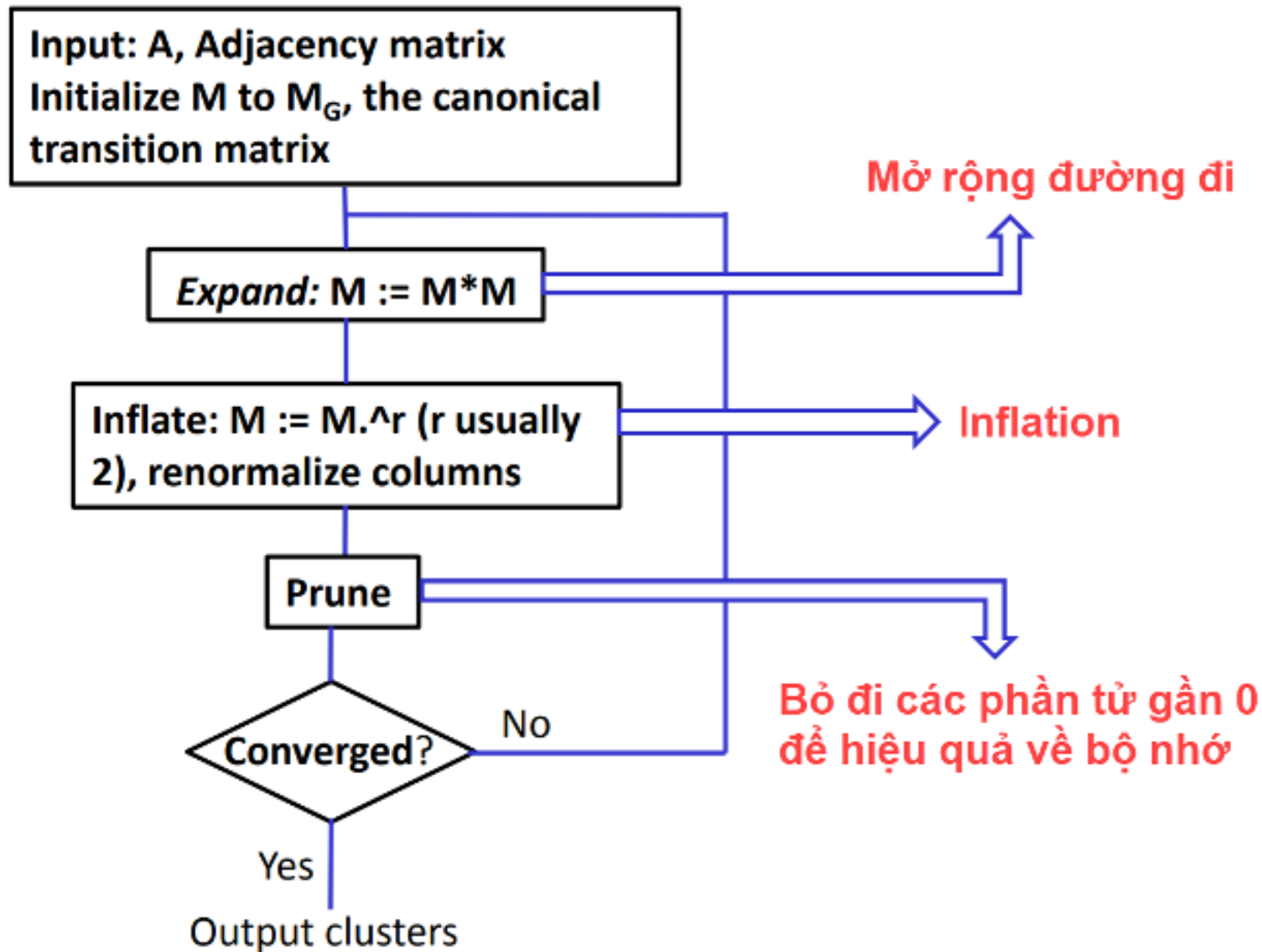
Squaring

$$\begin{bmatrix} 0 \\ 1/4 \\ 0 \\ 1/36 \\ 1/9 \end{bmatrix}$$

Normalization

$$\begin{bmatrix} 0 \\ 9/14 \\ 0 \\ 1/14 \\ 4/14 \end{bmatrix}$$

Dự đoán trên random walk



Dự đoán trên random walk

Input: A, Adjacency matrix
Initialize M to M_G , the canonical transition matrix

Expand: $M := M * M$

Inflate: $M := M.^r$ (r usually 2), renormalize columns

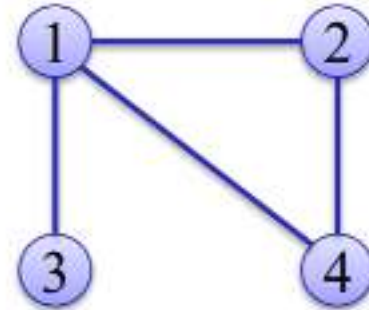
Prune

Converged?

No

Yes

Output clusters



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

$$\begin{bmatrix} 1/4 & 1/3 & 1/2 & 1/3 \\ 1/4 & 1/3 & 0 & 1/3 \\ 1/4 & 0 & 1/2 & 0 \\ 1/4 & 1/3 & 0 & 1/3 \end{bmatrix}$$

Dự đoán trên random walk

Input: A, Adjacency matrix
Initialize M to MG, the canonical transition matrix

Expand: $M := M * M$

Inflate: $M := M.^r$ (r usually 2), renormalize columns

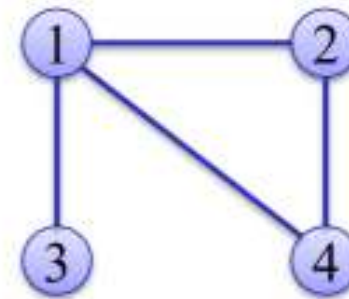
Prune

Converged?

No

Yes

Output clusters

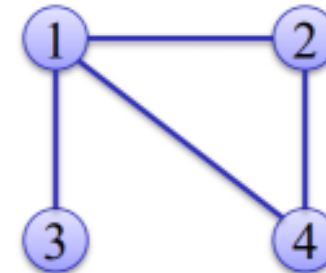
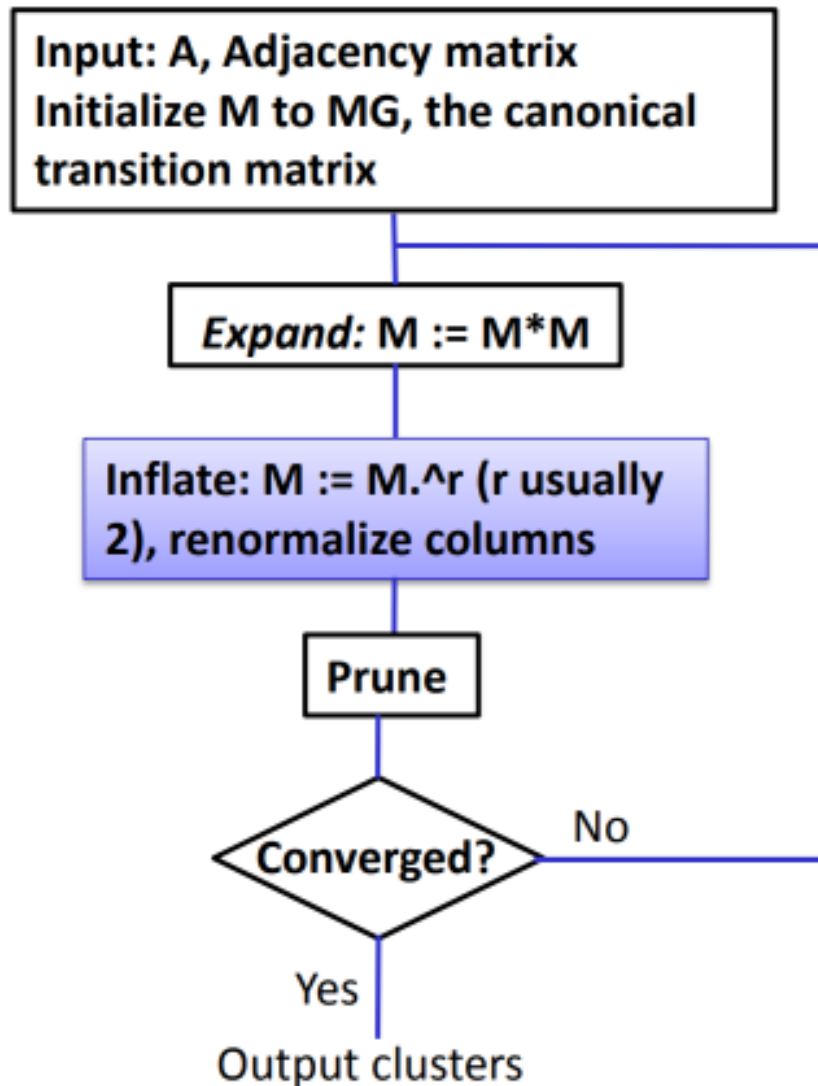


$$\begin{bmatrix} 1/4 & 1/3 & 1/2 & 1/3 \\ 1/4 & 1/3 & 0 & 1/3 \\ 1/4 & 0 & 1/2 & 0 \\ 1/4 & 1/3 & 0 & 1/3 \end{bmatrix} * \begin{bmatrix} 1/4 & 1/3 & 1/2 & 1/3 \\ 1/4 & 1/3 & 0 & 1/3 \\ 1/4 & 0 & 1/2 & 0 \\ 1/4 & 1/3 & 0 & 1/3 \end{bmatrix}$$

=

$$\begin{bmatrix} 0.35 & 0.31 & 0.38 & 0.31 \\ 0.23 & 0.31 & 0.13 & 0.31 \\ 0.19 & 0.08 & 0.38 & 0.08 \\ 0.23 & 0.31 & 0.13 & 0.31 \end{bmatrix}$$

Dự đoán trên random walk



$$\begin{bmatrix} 0.35 & 0.31 & 0.38 & 0.31 \\ 0.23 & 0.31 & 0.13 & 0.31 \\ 0.19 & 0.08 & 0.38 & 0.08 \\ 0.23 & 0.31 & 0.13 & 0.31 \end{bmatrix}$$

inflation

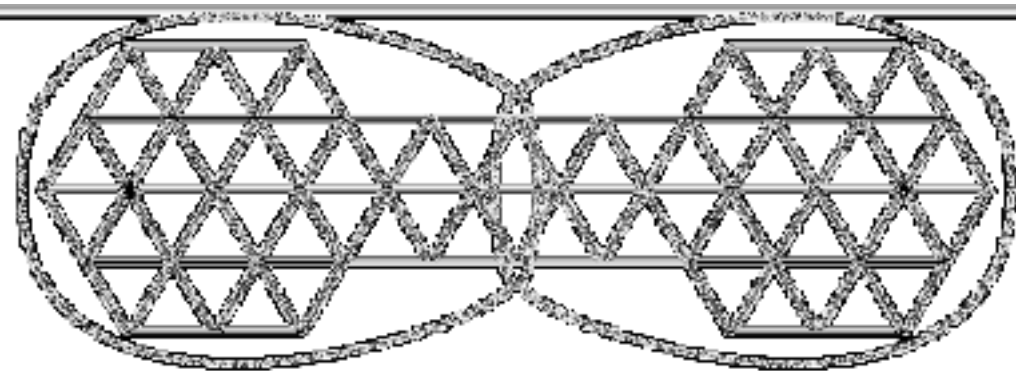
(Bình phương)

$$\begin{bmatrix} 0.13 & 0.09 & 0.14 & 0.09 \\ 0.05 & 0.09 & 0.02 & 0.09 \\ 0.04 & 0.01 & 0.14 & 0.01 \\ 0.05 & 0.09 & 0.02 & 0.09 \end{bmatrix}$$

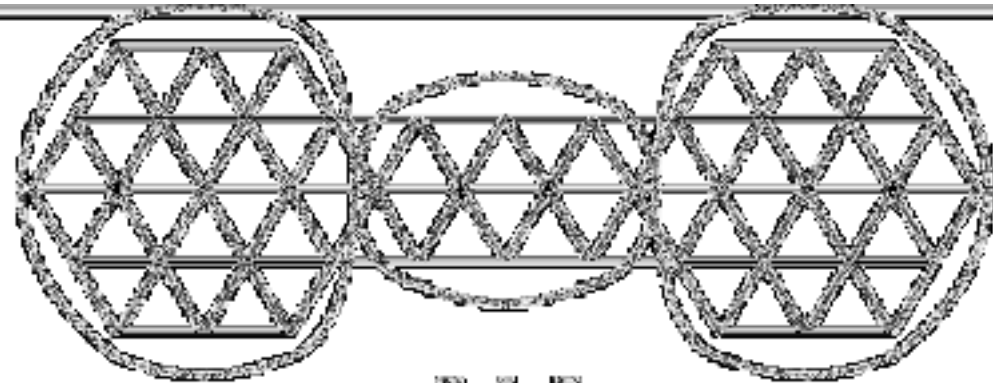
normalization

$$\begin{bmatrix} 0.47 & 0.33 & 0.45 & 0.33 \\ 0.20 & 0.33 & 0.05 & 0.33 \\ 0.13 & 0.02 & 0.45 & 0.02 \\ 0.20 & 0.33 & 0.05 & 0.33 \end{bmatrix}$$

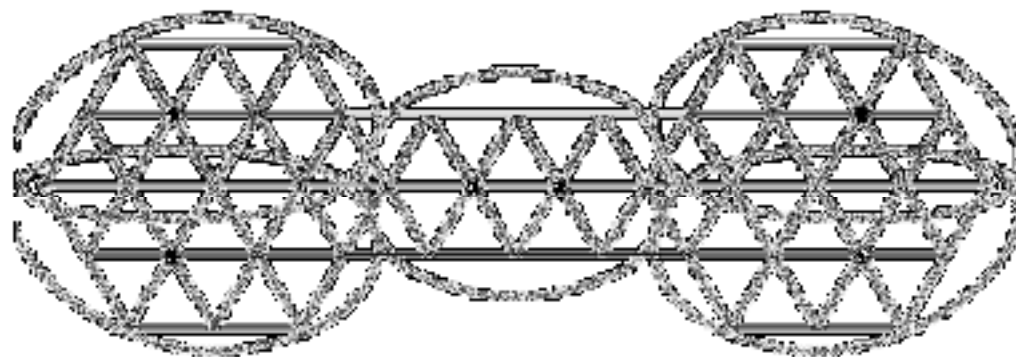
Các Inflation



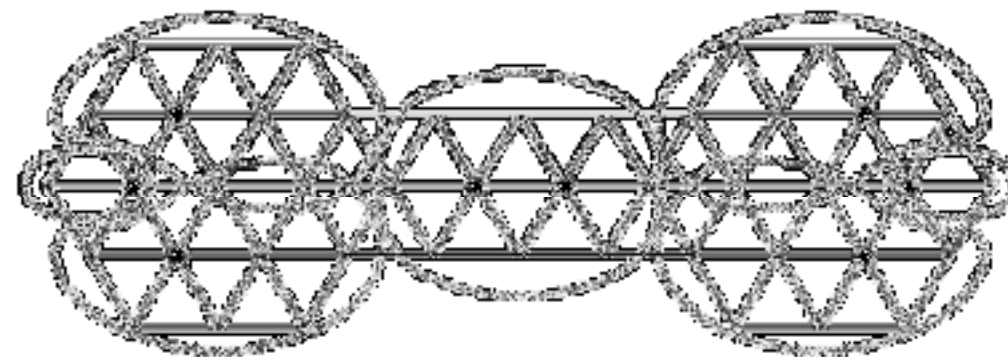
R 1.4



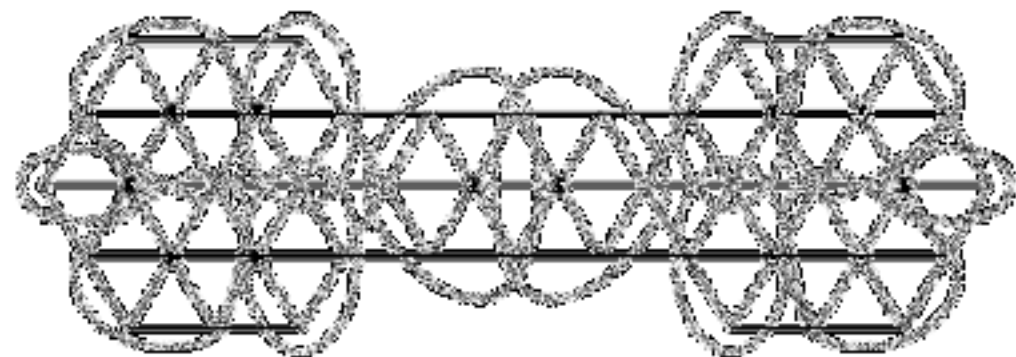
R 1.5



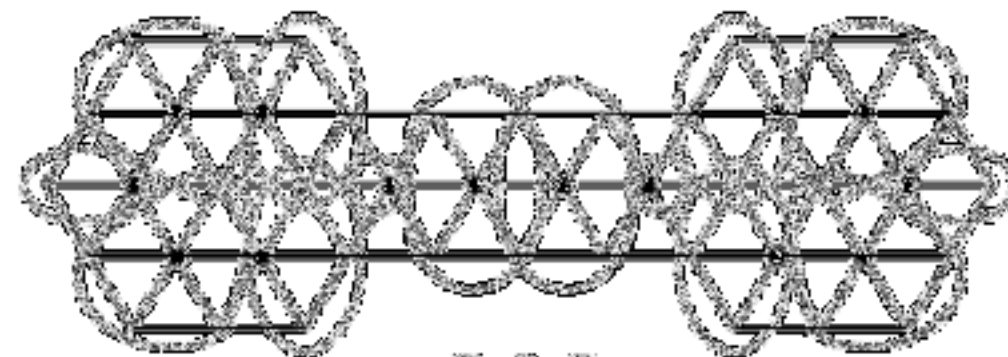
R 1.7



R 2.0



R 2.1



R 2.5

Tài liệu tham khảo

- Mihailcevic, Rados, and Dragomir Radev. *Graph-based natural language processing and information retrieval*. Cambridge university press, 2011.