

# Network analysis

---

TOPICS OF DATA ENGINEERING

# Networks appear everywhere

---

## ■ Examples

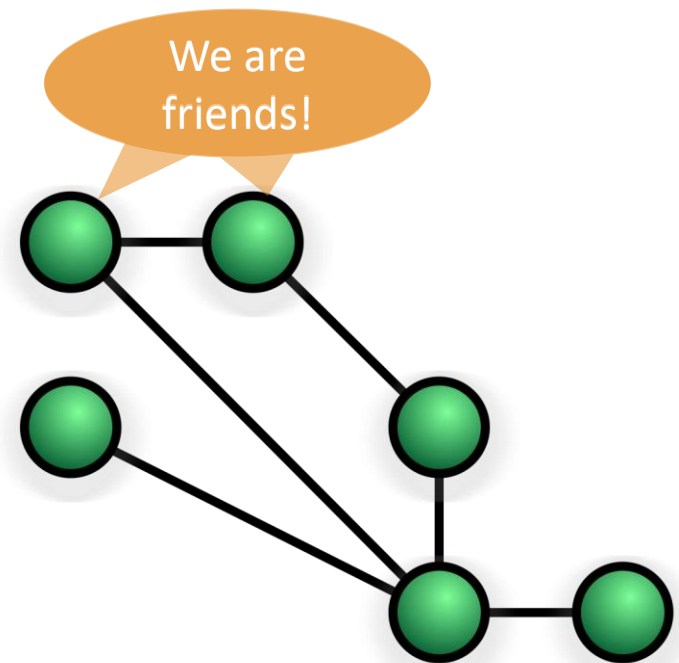
- WWW

- Social networks

  - Friend networks

  - SNS

- Chemical interactions  
in a human body



# Nodes and edges

---

## ■ Node (vertex)

- A fundamental unit of a graph (network).
- Represented as a dot

## ■ Edge

- Connects two nodes

## ■ Each edge connects two nodes

- Undirected edges
- Directed edges
- Weighted edges



# Adjacency matrix

---

- is a representative data structure to describe a network structure.
- The elements are one or zero depending on connectivity.

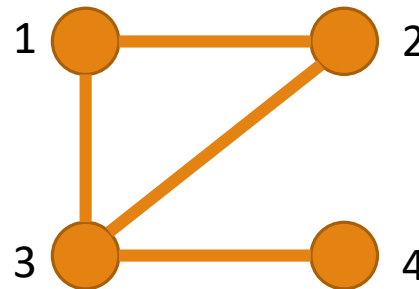
$$A_{ij} = \begin{cases} 1 \\ 0 \end{cases}$$



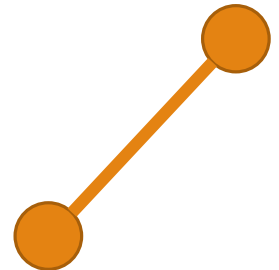
# Examples of adjacency matrices

Undirected graph

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



They are  
equivalent  
under

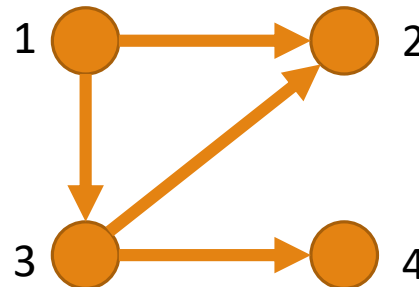


Directed graph

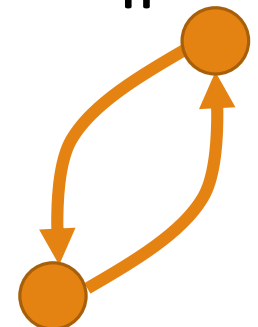
to

From

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



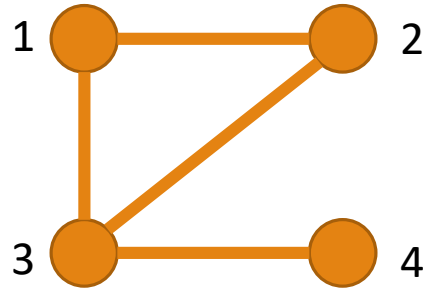
||



# Types of graph structures

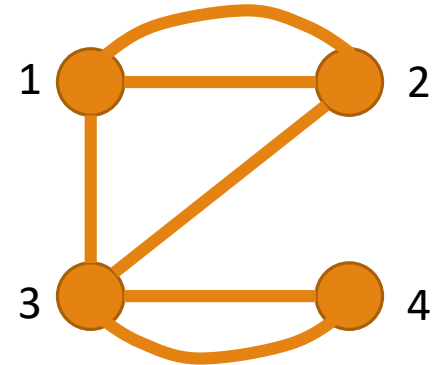
## ■ Simple graph

- No multiple edges
- No self loop



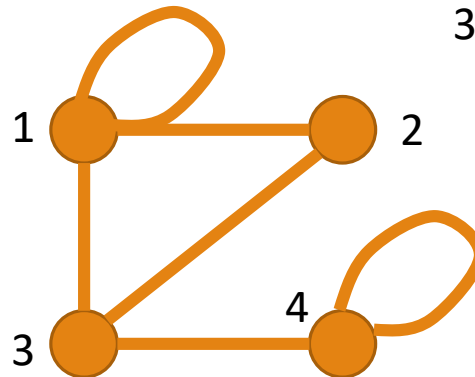
## ■ Multi graph

- Multiple edges exist
- No self loop



## ■ Graph with loop

- Self loops exist



Hereafter, we focus only on a simple graph.

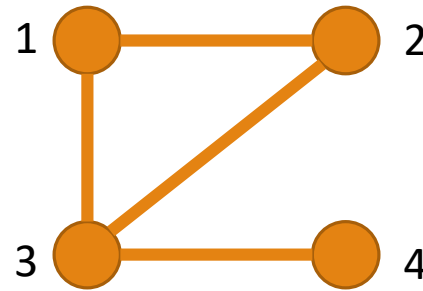
# Degree

---

- The number of edges from a node

$$k_1 = 2 \quad k_2 = 2$$

$$k_3 = 3 \quad k_4 = 1$$



- The sum of degrees are equal to twice the number of nodes  $M$

$$\sum_i k_i = 2 + 2 + 3 + 1 = 8 = 2M$$

# Centrality

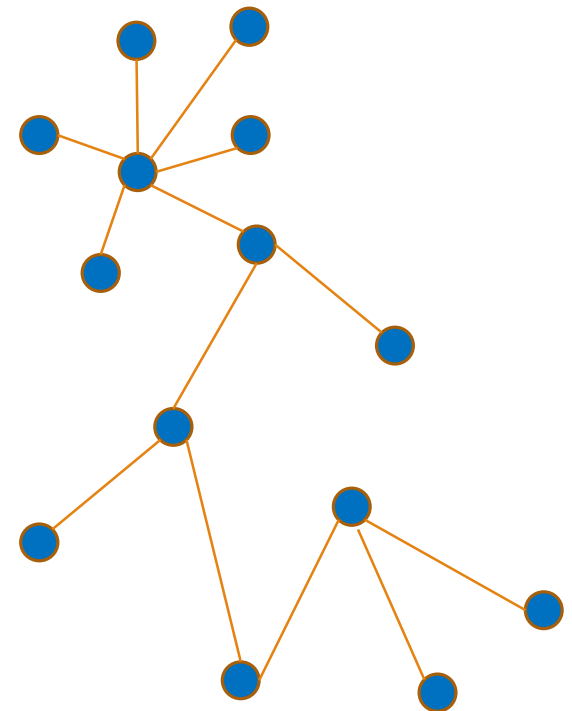
---



# Centrality

---

- Node characteristics to identify importance in a graph
- Many types
  - Degree centrality
  - Betweenness centrality
  - Eigenvector centrality
  - HITS
  - Page rank centrality



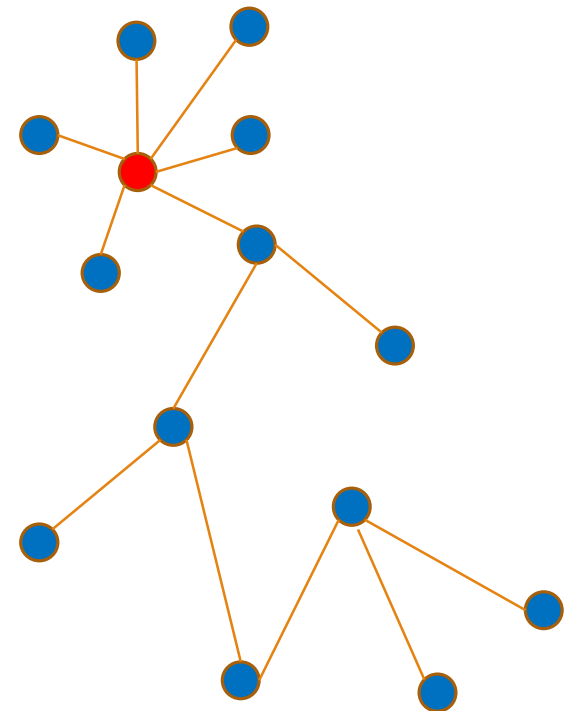
# Degree centrality

---

■ Based on the idea “the person who has the most friends is important”.

■ Definition

●  $C_d(i) = k_i$



# Betweenness centrality

---

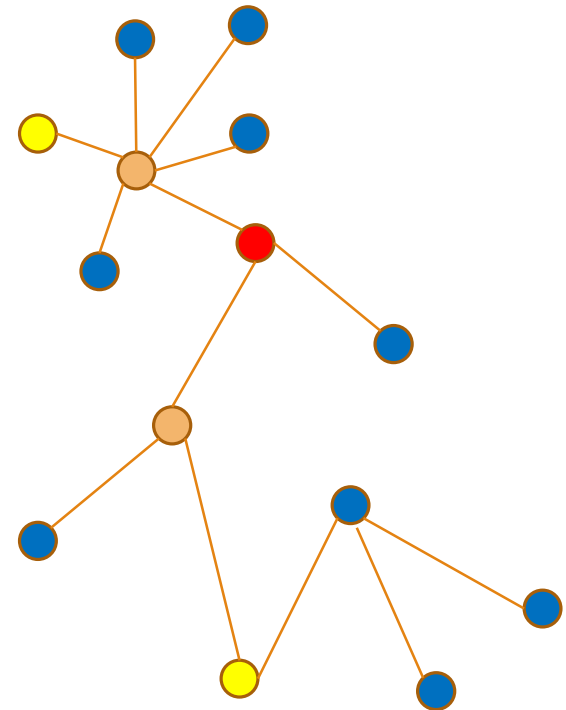
■ Based on the idea “the person who can mediate between the most pair of friends is important”.

■ Definition

$$C_b(i) = \frac{\sum_{j,k} sp_{j,k}(i)}{n(n-1)/2}$$

where  $sp_{j,k}(i) = \begin{cases} 1 & \text{if Node } i \text{ is on the shortest path} \\ & \text{between Node } j \text{ and Node } k. \\ 0 & \text{Otherwise} \end{cases}$

$n = \# \text{ of nodes}$



# Eigenvector centrality

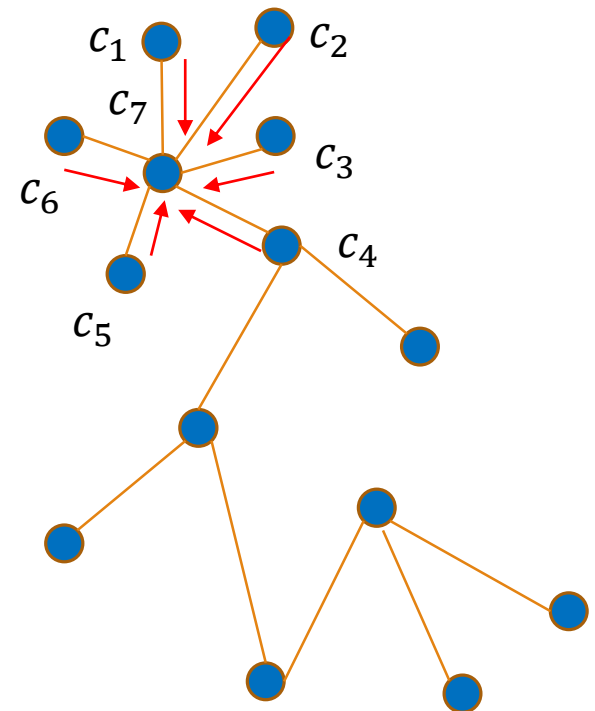
■ Based on the idea “if important persons think a friend is important, he/she is important”.

■ Definition

●  $c_i^{(n+1)} = \frac{1}{\lambda} \sum_j A_{ij} c_j^{(n)}$

●  $\mathbf{c}^{(n+1)} = \frac{1}{\lambda} A \mathbf{c}^{(n)}$

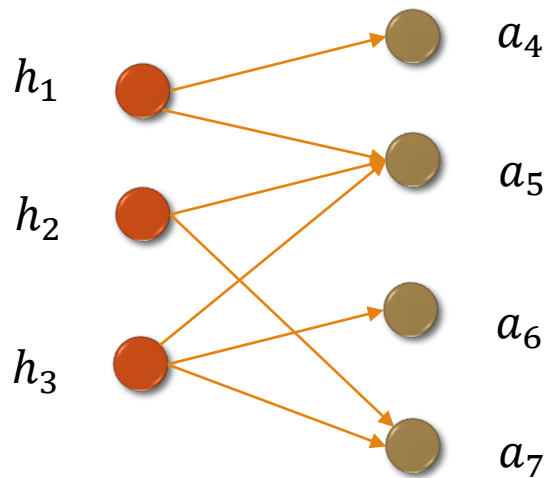
● If  $\lim_{n \rightarrow \infty} \mathbf{c}^{(n)} = \mathbf{c}$ ,  $A\mathbf{c} = \lambda\mathbf{c}$



# HITS

---

- An extension of eigenvector centrality
- It assumes hubs and authorities



$$B = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 1 \end{pmatrix}$$

$$\mathbf{h}^{(n+1)} = \frac{1}{\lambda_1} B \mathbf{a}^{(n)}$$

$$\mathbf{a}^{(n+1)} = \frac{1}{\lambda_2} B^T \mathbf{h}^{(n)}$$

# PageRank

- The Google algorithm of web page importance algorithm

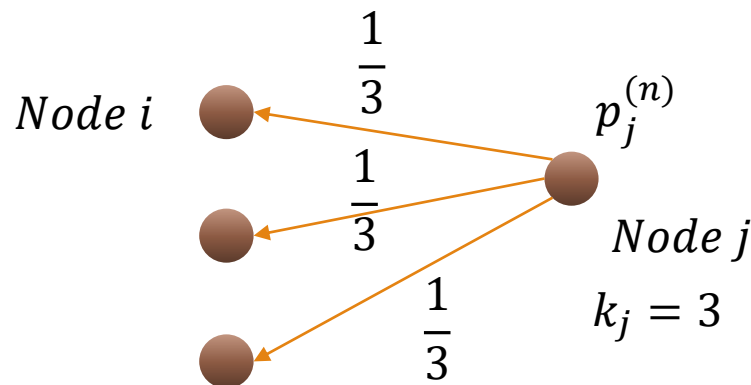
$n = \# \text{ of nodes}$

- $$p_i^{(n+1)} = \alpha \sum_j \frac{A_{ij}}{k_j} p_j^{(n)} + (1 - \alpha) \frac{1}{n}$$

The probability  
to see Node i

The probability  
after seeing Node j

The probability  
to randomly see Node i



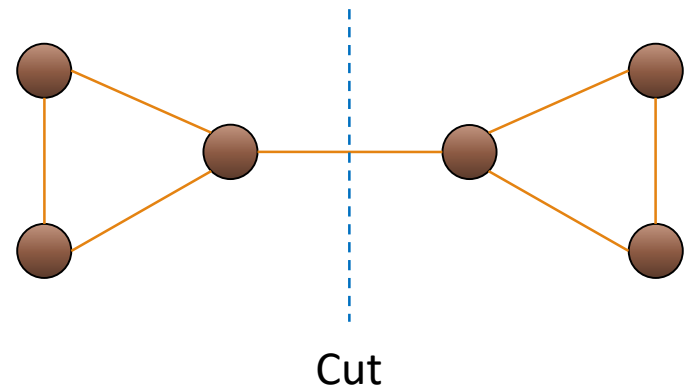
# Network clustering

---

# Cluster

---

- A part of the network which is densely connected to each other
- Cut the edges which connect clusters (MinCut)
- Then, how should we define “densely connected” or “MinCut”?





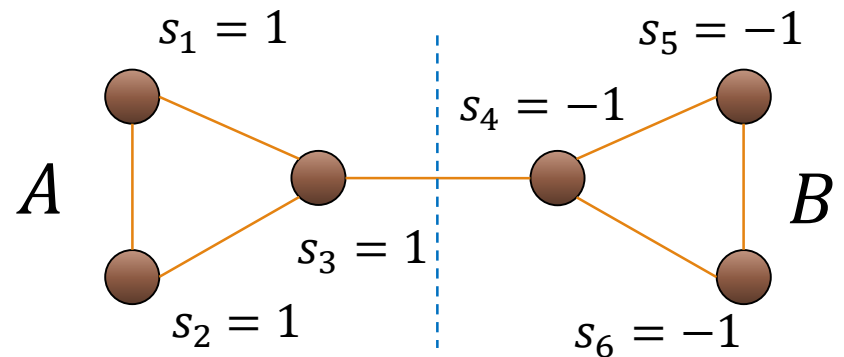
# Spectral clustering

---

# Spectral clustering

---

- Separate a network into two clusters
  - easily extendable to more clusters
- Assign +1 or -1 as a cluster label to each node
- Two clusters should have a minimum bridge between them



# Normalized cut – (1)

---

- $cut(A, B)$ : the number of edges bridging Cluster A and Cluster B

$$\begin{aligned} cut(A, B) &= \sum_{i \in A} \sum_{j \in B} A_{ij} \\ &= \frac{1}{2} \cdot \frac{1}{4} \sum_{i,j} A_{ij} (s_i - s_j)^2 \end{aligned}$$

# Normalized cut – (2)

---

- $vol(A)$ : the number of edges connected to nodes in Cluster A

$$vol(A) = \sum_{i \in A} k_i = \frac{1}{2} \sum_i k_i (1 - s_i)$$
$$vol(B) = \sum_{i \in B} k_i = \frac{1}{2} \sum_j k_j (1 + s_j)$$

# Normalized cut – (3)

---

■  $ncut(A, B)$ : the ratio of the number of bridges to nodes in clusters

$$\begin{aligned} ncut(A, B) &= \frac{cut(A, B)}{vol(A)} + \frac{cut(A, B)}{vol(B)} \\ &= \underbrace{cut(A, B)}_{\text{Part 1}} \underbrace{\left( \frac{1}{vol(A)} + \frac{1}{vol(B)} \right)}_{\text{Part 2}} \end{aligned}$$

# Part 2

---

$$\begin{aligned}
 \frac{2}{\sum_i k_i(1-s_i)} + \frac{2}{\sum_j k_j(1+s_j)} &= 2\left(\frac{1}{2\sum_i k_i - \sum_i k_i(1+s_i)} + \frac{1}{\sum_j k_j(1+s_j)}\right) \\
 &= \frac{1}{\sum_i k_i} \left( \frac{1}{1 - \frac{\sum_i k_i(1+s_i)}{2\sum_i k_i}} + \frac{1}{\frac{\sum_i k_i(1+s_i)}{2\sum_i k_i}} \right) \\
 &= \frac{1}{\sum_i k_i} \left( \frac{1}{1-r} + \frac{1}{r} \right) \\
 &= \frac{1}{\sum_i k_i} \frac{1}{r(1-r)} = \frac{1}{\sum_i k_i \left( \frac{1-r}{2}(1+s_i) - \frac{r}{2}(1-s_i) \right)^2}
 \end{aligned}$$

# The detail of the calculation

---

$$\begin{aligned}& \sum_i k_i \left( \frac{(1-r)(1+s_i) - r(1-s_i)}{2} \right)^2 \\&= \frac{1}{4} \sum_i k_i (1-r)^2 (1+s_i)^2 - \frac{2}{4} r(1-r) \sum_i k_i (1+s_i)(1-s_i) + \frac{r^2}{4} \sum_i k_i (1-s_i)^2 \\&= \frac{1}{2} (1-r)^2 \sum_i k_i (1+s_i) + \frac{1}{2} r^2 \sum_i k_i (1-s_i) \\&= (1-r)^2 \sum_i k_i r + r^2 \sum_i k_i (1-r) \\&= \sum_i k_i r(1-r)\end{aligned}$$

# Part 1

---

$$\begin{aligned}\frac{1}{8} \sum_{i,j} A_{ij} (s_i^2 - 2s_i s_j + s_j^2) &= \frac{1}{4} \sum_{i,j} (A_{ij} s_i^2 - A_{ij} s_i s_j) \\&= \frac{1}{4} \sum_{i,j} (k_i \delta_{ij} s_i s_j - A_{ij} s_i s_j) \\&= \frac{1}{4} \sum_{i,j} (k_i \delta_{ij} - A_{ij}) s_i s_j \\&= \frac{1}{2} \sum_{i,j} (k_i \delta_{ij} - A_{ij}) \frac{(1-r)(1+s_i) - r(1-s_i)}{2} \frac{(1-r)(1+s_j) - r(1-s_j)}{2}\end{aligned}$$



# The detail of the calculation

---

$$\begin{aligned} & \sum_{i,j} (k_i \delta_{ij} - A_{ij}) \frac{(1-r)(1+s_i) - r(1-s_i)}{2} \frac{(1-r)(1+s_j) - r(1-s_j)}{2} \\ &= \frac{1}{4} \sum_{i,j} (k_i \delta_{ij} - A_{ij}) ((1-r)^2(1+s_i)(1+s_j) + r^2(1-s_i)(1-s_j) - 2r(1-r)(1-s_i)(1+s_j)) \\ &= \frac{1}{2} \sum_{i,j} (k_i \delta_{ij} - A_{ij}) (r^2 s_i s_j + (1-r)^2 s_i s_j + 2r(1-r)s_i s_j) \\ &= \frac{1}{2} \sum_{i,j} (k_i \delta_{ij} - A_{ij}) (1 - 2r + 2r^2 + 2r(1-r)) s_i s_j \\ &= \frac{1}{2} \sum_{i,j} (k_i \delta_{ij} - A_{ij}) s_i s_j \end{aligned}$$

# Normalized cut – (3)

---

■  $ncut(A, B)$ : the ratio of the number of bridges to nodes in clusters

$$\begin{aligned} & 2 \cdot ncut(A, B) \\ &= \frac{\sum_{i,j} (k_i \delta_{ij} - A_{ij}) y_i y_j}{\sum_i k_i y_i^2} \\ &= \frac{\mathbf{y}^T (D - A) \mathbf{y}}{\mathbf{y}^T D \mathbf{y}} \end{aligned}$$

where  $y_i = \frac{(1-r)(1+s_i)-r(1-s_i)}{2}$

# Relaxation of condition

---

- We want to find  $\mathbf{y}$  that minimizes  $\frac{\mathbf{y}^T (D-A) \mathbf{y}}{\mathbf{y}^T D \mathbf{y}}$
- Let elements in  $\mathbf{y}$  take arbitrary real values
- Since there is a scaling symmetry,  $\mathbf{y} \rightarrow a\mathbf{y}$  gives the same ncut value.
  - Take  $\mathbf{y}$  so as to  $\mathbf{y}^T D \mathbf{y} = 1$
  - Then, the problem is to minimize  $\mathbf{y}^T (D - A) \mathbf{y}$  under the condition  $\mathbf{y}^T D \mathbf{y} = 1$

# Generalized eigenvalue problem

---

- Applying Lagrange multiplier, we get

$$L = \mathbf{y}^T (D - A) \mathbf{y} - \lambda (\mathbf{y}^T D \mathbf{y} - 1)$$

$$\frac{1}{2} \frac{\partial L}{\partial \mathbf{y}} = (D - A) \mathbf{y} - \lambda D \mathbf{y} = 0$$

$$(D - A) \mathbf{y} = \lambda D \mathbf{y}$$

# After solving the generalized eigenvalue problem ...

---

- Apply K-means problem

# Modularity maximization

---

# Modularity $Q$

---

- $e_{aa} = \sum_{ij} \frac{A_{ij}}{2M} \delta(c_i, a) \delta(c_j, a)$ 
  - the probability to find edges connecting nodes inside Cluster  $a$
- $a_a = \sum_b e_{ab} = \sum_i \frac{k_i}{2M} \delta(c_i, a)$ 
  - the probability of edges connecting to Cluster  $a$
  - $a_a^2$  is a probability to find edges under the assumption that edges are randomly connected to nodes
- $Q = \sum_a (e_{aa} - a_a^2) = \frac{1}{2M} \sum_{ij} (A_{ij} - \frac{k_i k_j}{2M}) \delta(c_i, c_j)$

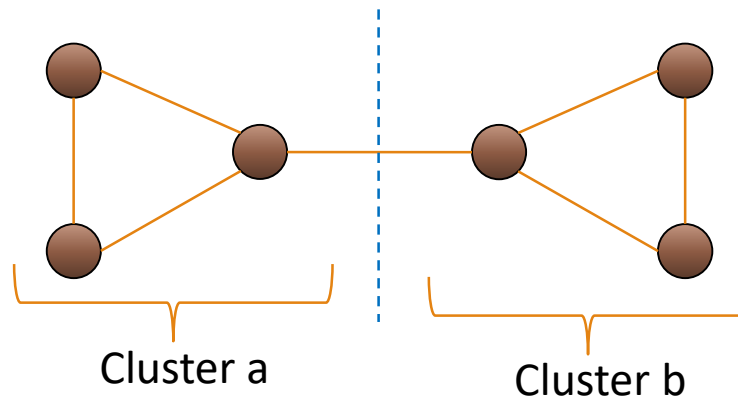
# An example

■ The number of edges  $M = 7$

■  $e_{aa} = \frac{3}{7}, e_{ab} = \frac{1}{7}, e_{bb} = \frac{3}{7}$

■  $a_a = e_{aa} + e_{ab} = \frac{4}{7}, a_a^2 = \frac{16}{49} = a_b^2$

■  $Q = e_{aa} - a_a^2 + e_{bb} - a_b^2 = 2 \left( \frac{3}{7} - \frac{16}{49} \right) = \frac{10}{49}$





# Resolution limit

---

- Modularity maximization tends to give larger clusters than expected.

