



INTRODUCTION TO DATA ANALYTICS

Social Network Analysis

(SNA)

Dr. Rathachai Chawuthai

Department of Computer Engineering

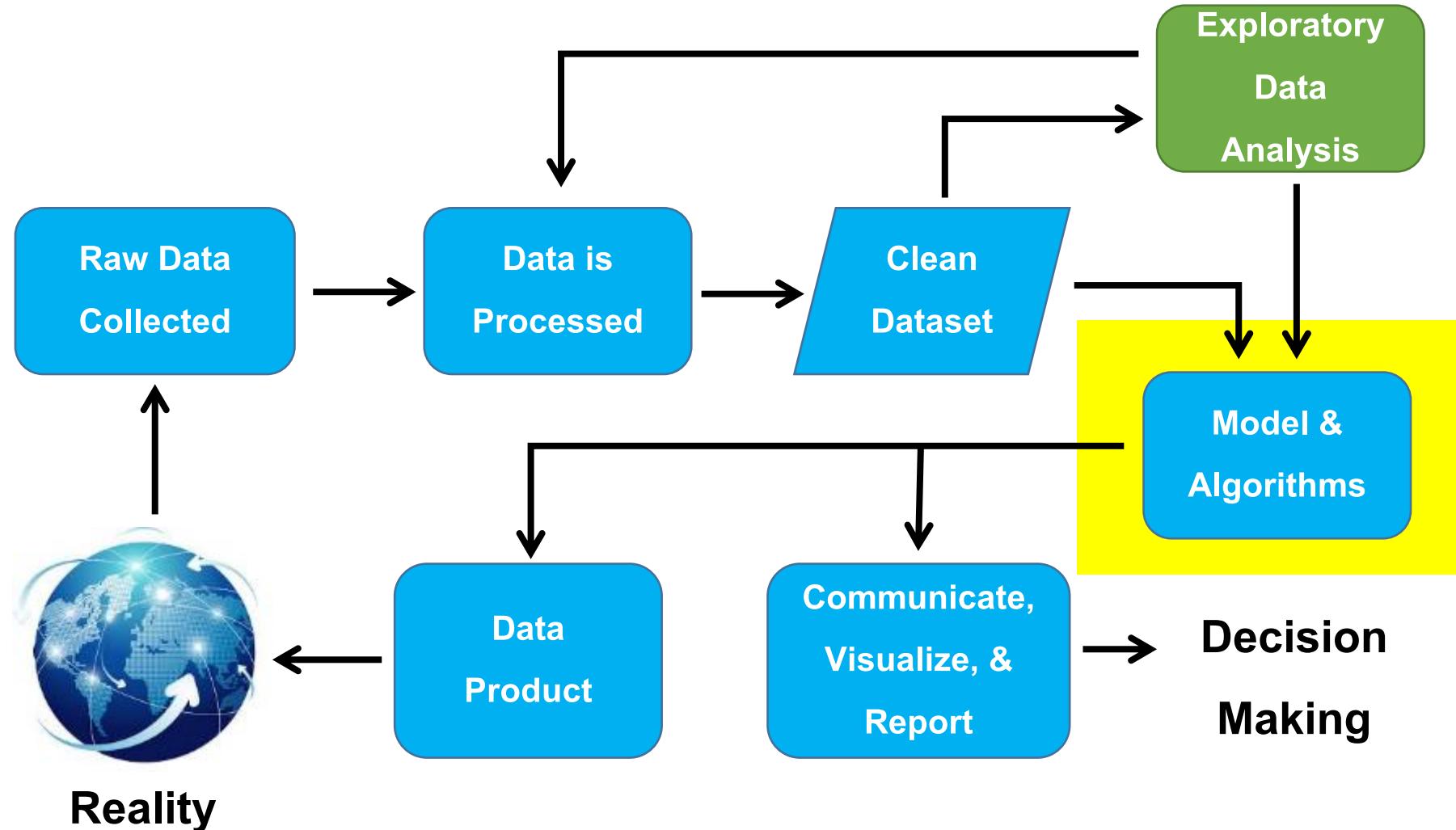
Faculty of Engineering

King Mongkut's Institute of Technology Ladkrabang

Agenda

- Introduction
- SNA Properties
- Community Detection

Data Science Process



Machine Learning



Supervised Learning

Develop predictive model based on both input and output data



Unsupervised Learning

Develop predictive model based on both input and output data



Regression

- Linear Regression
- Polynomial Regression



Classification

- Decision Tree
- Logistic Regression
- Neural Network
- etc.

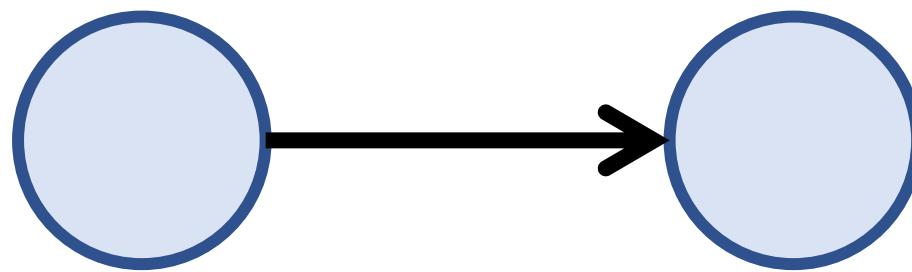


Clustering

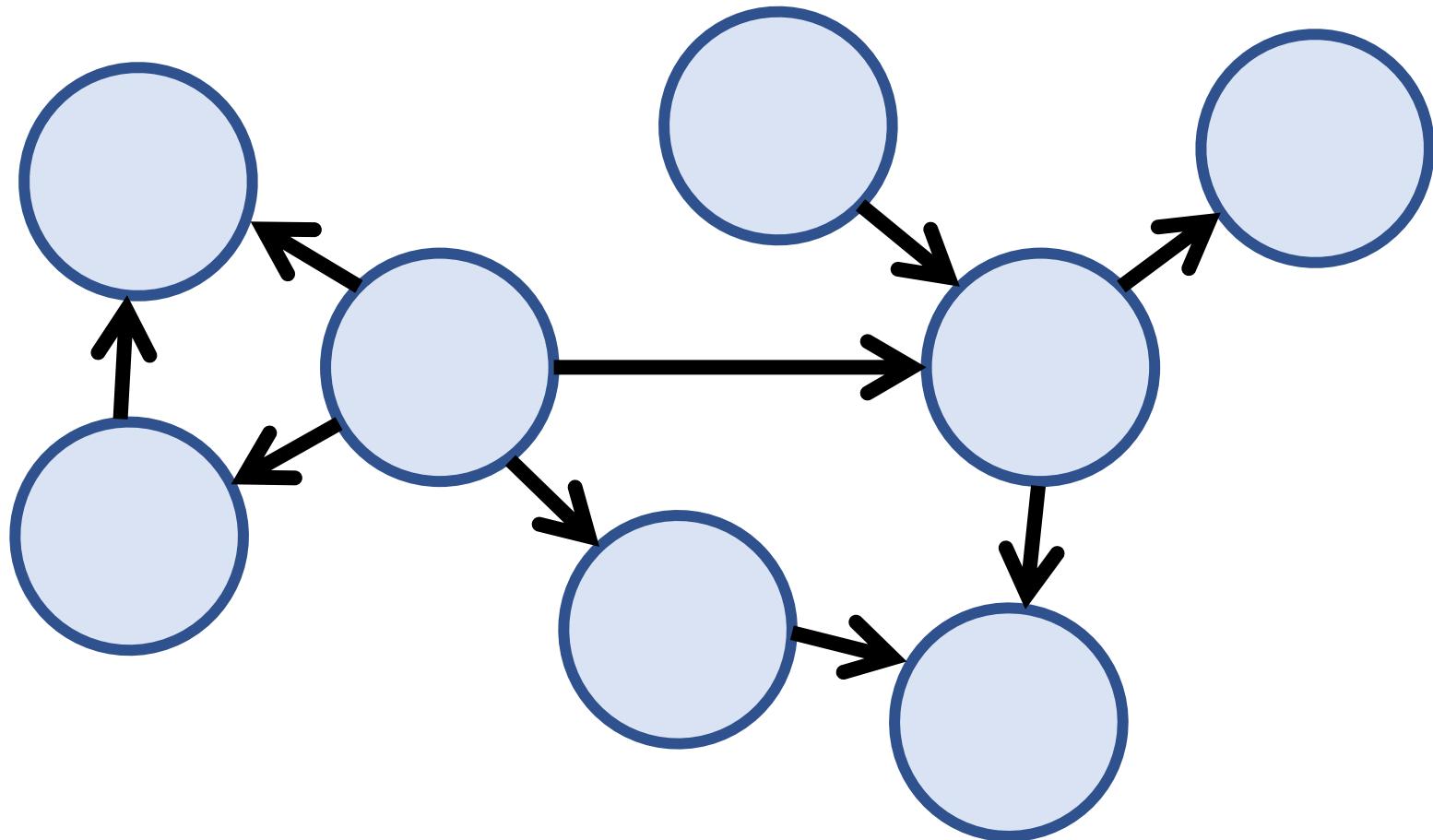
- K-Means
- DB-SCAN
- etc.

Introduction

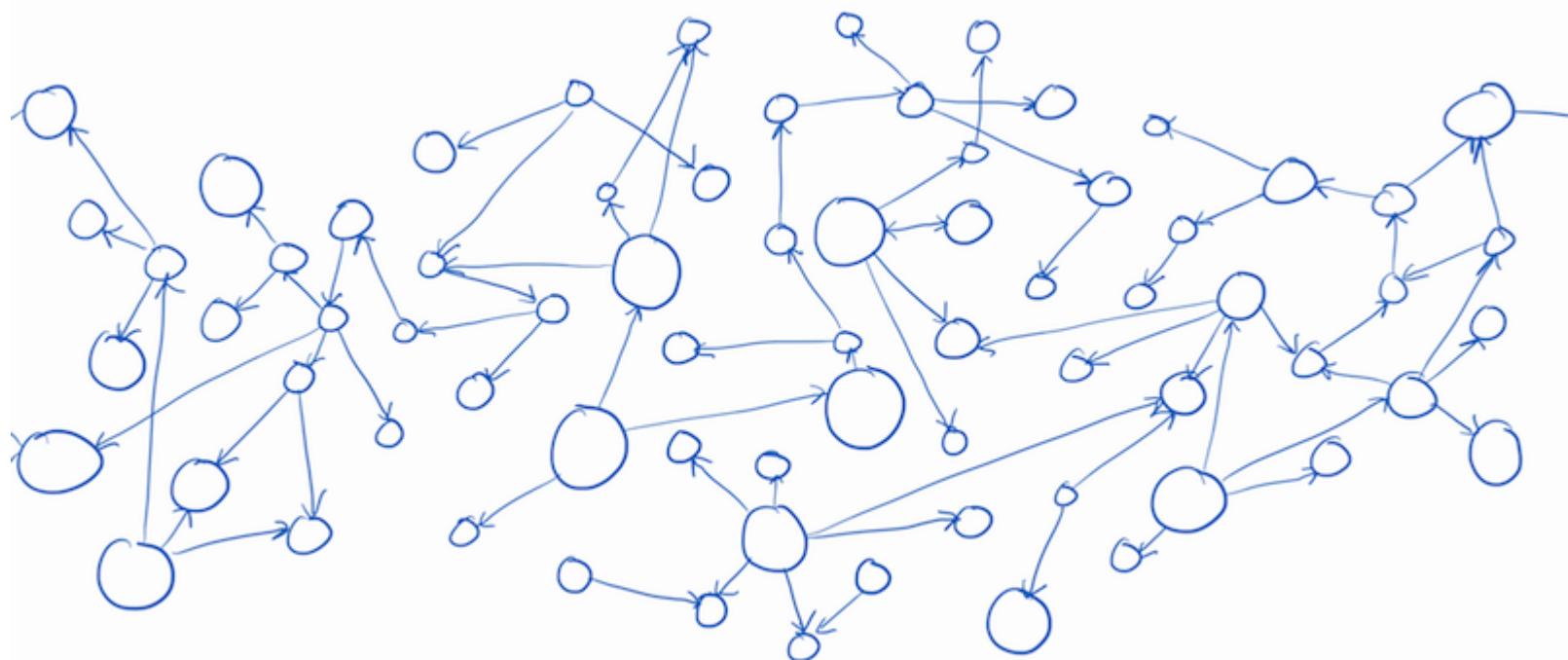
Graph



Graph



Graph



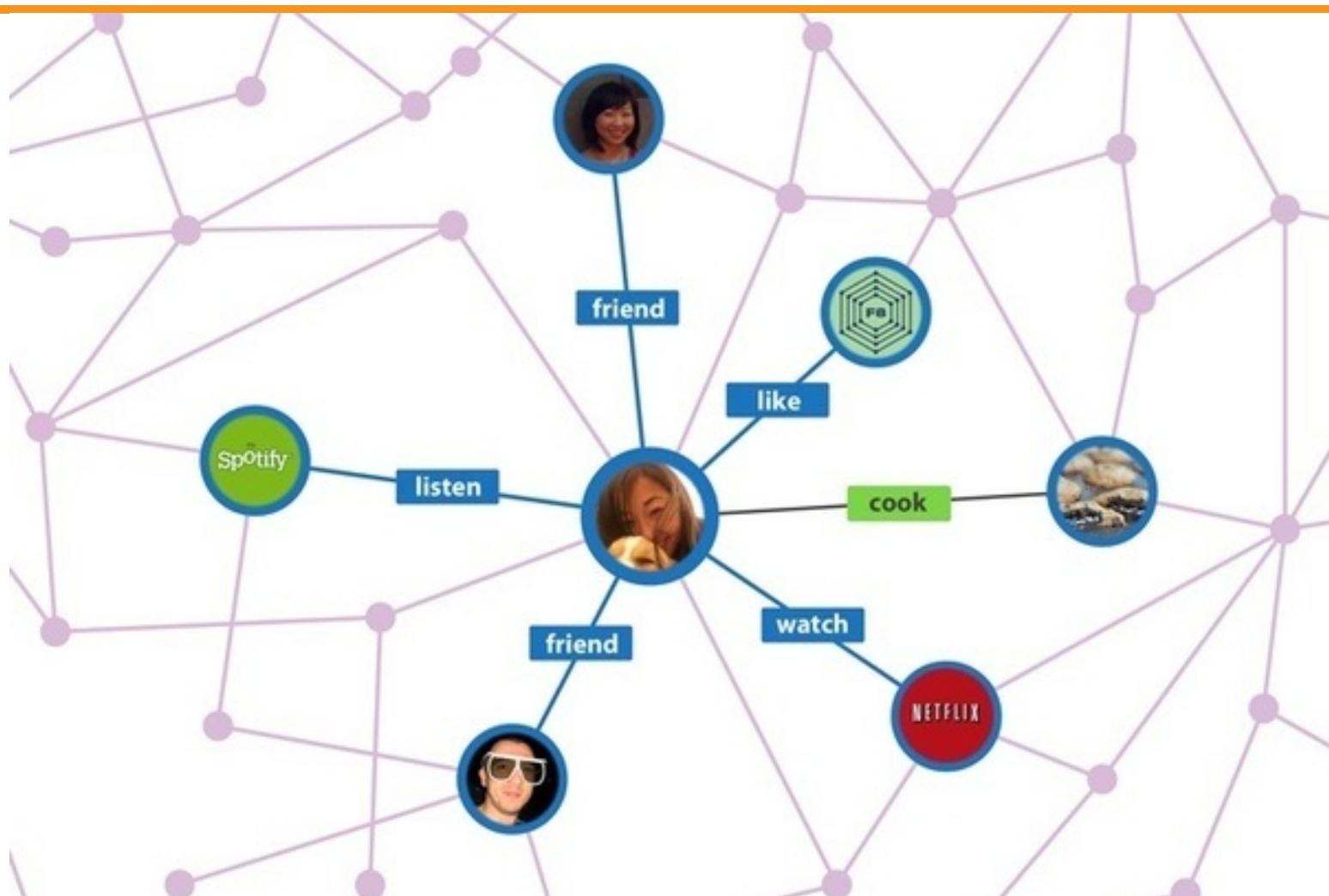
Graph



Graph



Graph



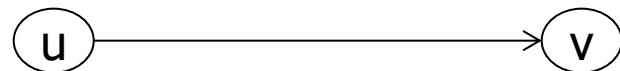
Definitions - Graph

A generalization of the simple concept of a set of dots, links, edges or arcs.

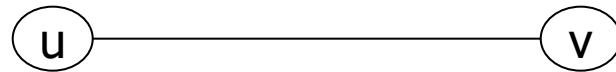
Representation: Graph $G = (V, E)$ consists set of vertices denoted by V , or by $V(G)$ and set of edges E , or $E(G)$

Edge Type

Directed: Ordered pair of vertices. Represented as (u, v) directed from vertex u to v .

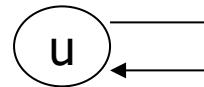


Undirected: Unordered pair of vertices. Represented as $\{u, v\}$. Disregards any sense of direction and treats both end vertices interchangeably.



Edge Type

- **Loop:** A loop is an edge whose endpoints are equal i.e., an edge joining a vertex to it self is called a loop.
Represented as $\{u, u\} = \{u\}$

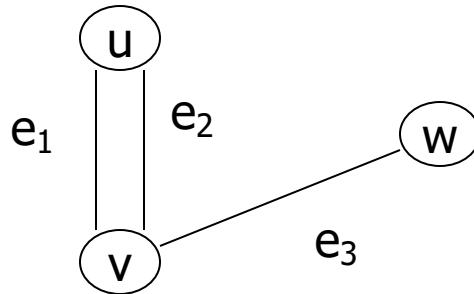


- **Multiple Edges:** Two or more edges joining the same pair of vertices.

Graph Type

Multigraph: $G(V, E)$, consists of set of vertices V , set of Edges E and a function f from E to $\{\{u, v\} \mid u, v \in V, u \neq v\}$. The edges e_1 and e_2 are called multiple or parallel edges if $f(e_1) = f(e_2)$.

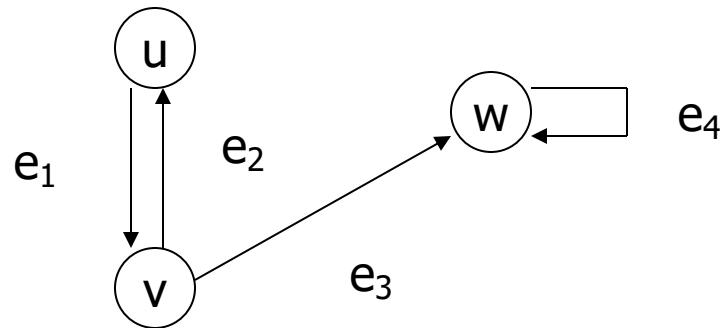
Representation Example: $V = \{u, v, w\}$, $E = \{e_1, e_2, e_3\}$



Graph Type

Directed Multigraph: $G(V, E)$, consists of set of vertices V , set of Edges E and a function f from E to $\{(u, v) \mid u, v \in V\}$. The edges e_1 and e_2 are multiple edges if $f(e_1) = f(e_2)$

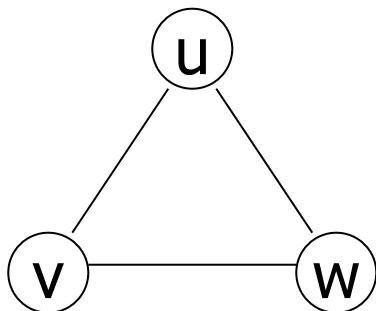
Representation Example: $V = \{u, v, w\}$, $E = \{e_1, e_2, e_3, e_4\}$



Representation- Adjacency Matrix

குமூலரியாக

- Example: Undirected Graph $G(V, E)$

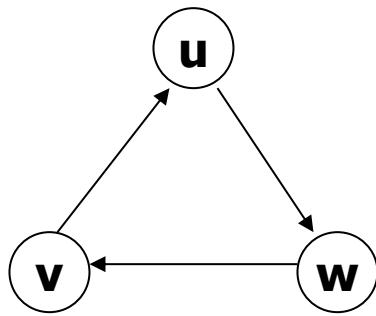


நுதான்சான்

	v	u	w
v	0	1	1
u	1	0	1
w	1	1	0

Representation- Adjacency Matrix

- Example: Directed Graph G (V, E)



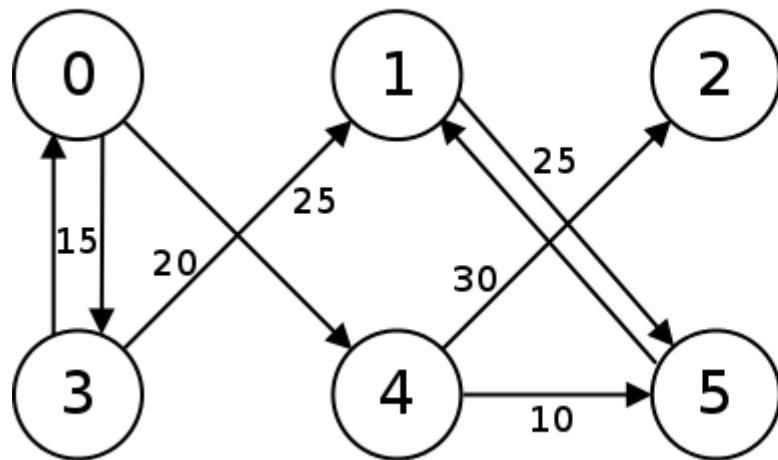
Adjacency Matrix:

	v	u	w
v	0	1	0
u	0	0	1
w	1	0	0

Representation- Adjacency Matrix

- Example: Weighted Graph G (V, E, w)

ແບບນີ້ນຳໃຊ້ກໍດັລ່າຍ່າງ

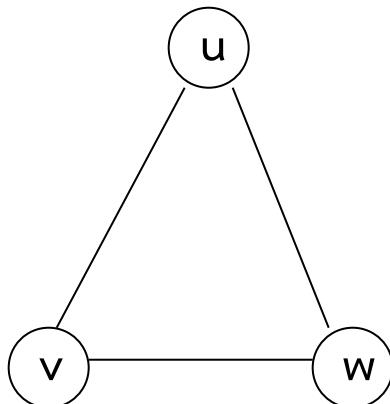


	0	1	2	3	4	5
0				15	20	
1				25		
2						
3	15	25				
4		30			10	
5				25		

Representation- Adjacency List

- Each node (vertex) has a list of which nodes (vertex) it is adjacent

Example: undirected graph $G(V, E)$



node	Adjacency List
u	v , w
v	w, u
w	u , v

Complete Graph

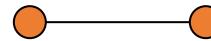
☞ node (ຝຶ່ງກົງເປັນ)

- Complete graph: K_n , is the simple graph that contains exactly one edge between each pair of distinct vertices.

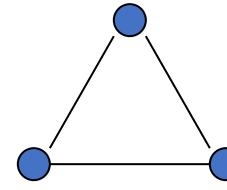
Representation Example: K_1 , K_2 , K_3 , K_4



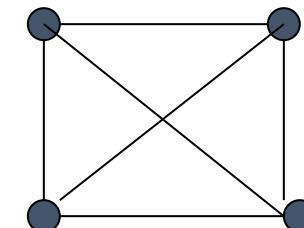
K_1



K_2



K_3



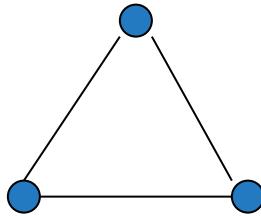
K_4

Cycle

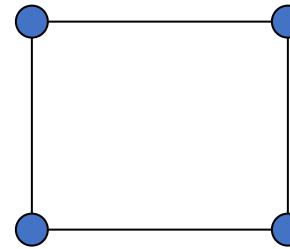
கிருஷ்ணபுரம்

- **Cycle:** C_n , $n \geq 3$ consists of n vertices $v_1, v_2, v_3 \dots v_n$ and edges $\{v_1, v_2\}, \{v_2, v_3\}, \{v_3, v_4\} \dots \{v_{n-1}, v_n\}, \{v_n, v_1\}$

Representation Example: C_3, C_4



C_3



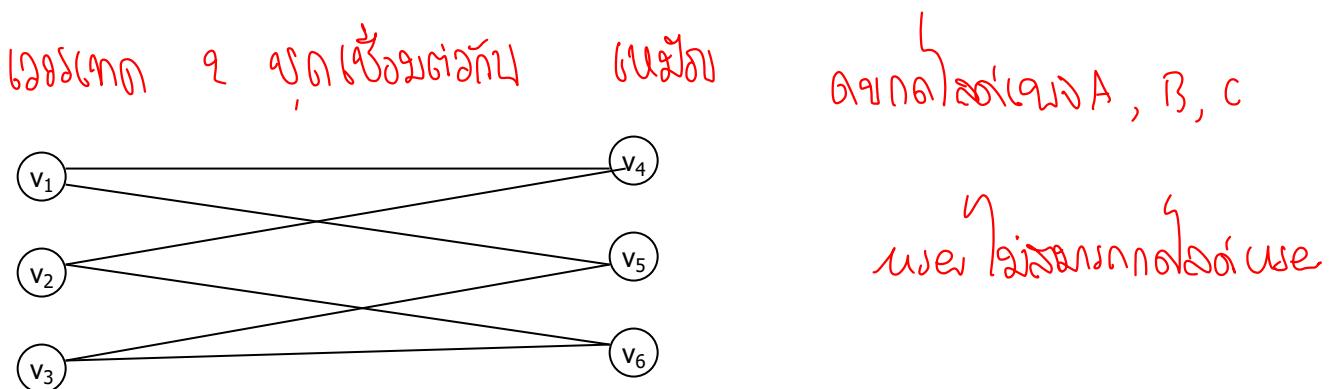
C_4

Bipartite Graphs

- In a simple graph G , if V can be partitioned into two disjoint sets V_1 and V_2 such that every edge in the graph connects a vertex in V_1 and a vertex V_2 (so that no edge in G connects either two vertices in V_1 or two vertices in V_2)

Application example: Representing Relations

Representation example: $V_1 = \{v_1, v_2, v_3\}$ and $V_2 = \{v_4, v_5, v_6\}$,

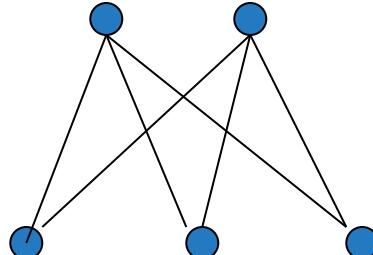


Complete Bipartite Graph

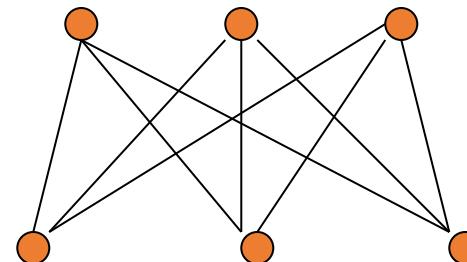
ପୂର୍ଣ୍ଣ ବିଭାଗୀୟ ମୂଳକ ଗ୍ରାଫ୍

- $K_{m,n}$ is the graph that has its vertex set portioned into two subsets of m and n vertices, respectively. There is an edge between two vertices if and only if one vertex is in the first subset and the other vertex is in the second subset.

Representation example: $K_{2,3}$, $K_{3,3}$



$K_{2,3}$



$K_{3,3}$

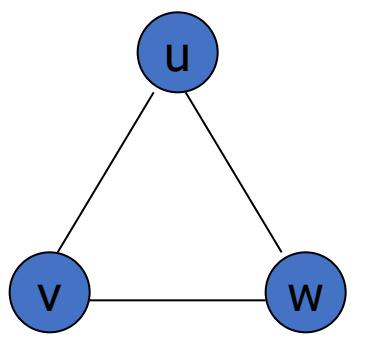
Subgraphs

(អំពីរាយការណ៍..)

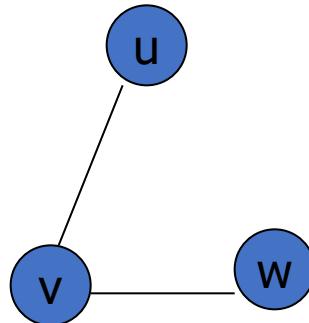
- A subgraph of a graph $G = (V, E)$ is a graph $H = (V', E')$ where V' is a subset of V and E' is a subset of E

Application example: solving sub-problems within a graph

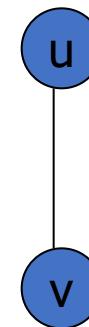
Representation example: $V = \{u, v, w\}$, $E = \{\{u, v\}, \{v, w\}, \{w, u\}\}$, H_1, H_2



G



H_1



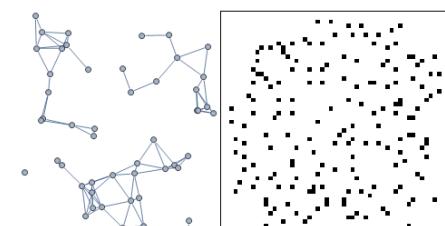
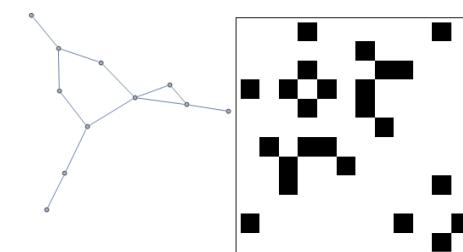
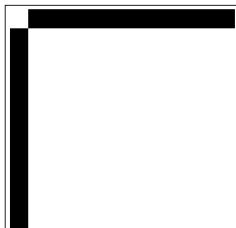
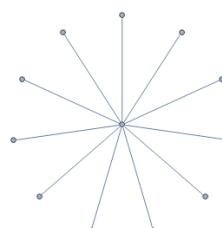
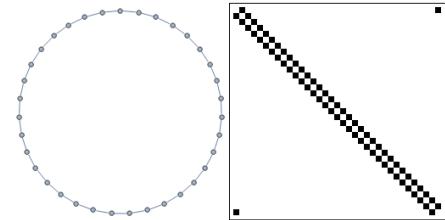
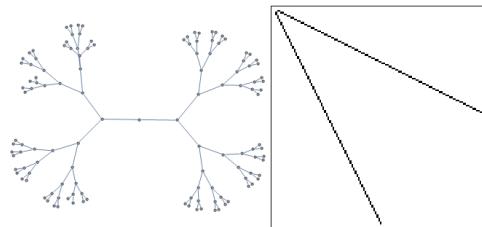
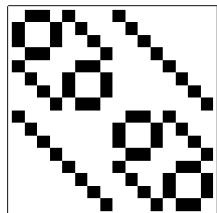
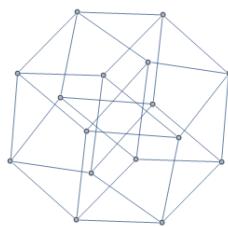
H_2

More Graphs

សំណើរាយការណ៍

- Go to

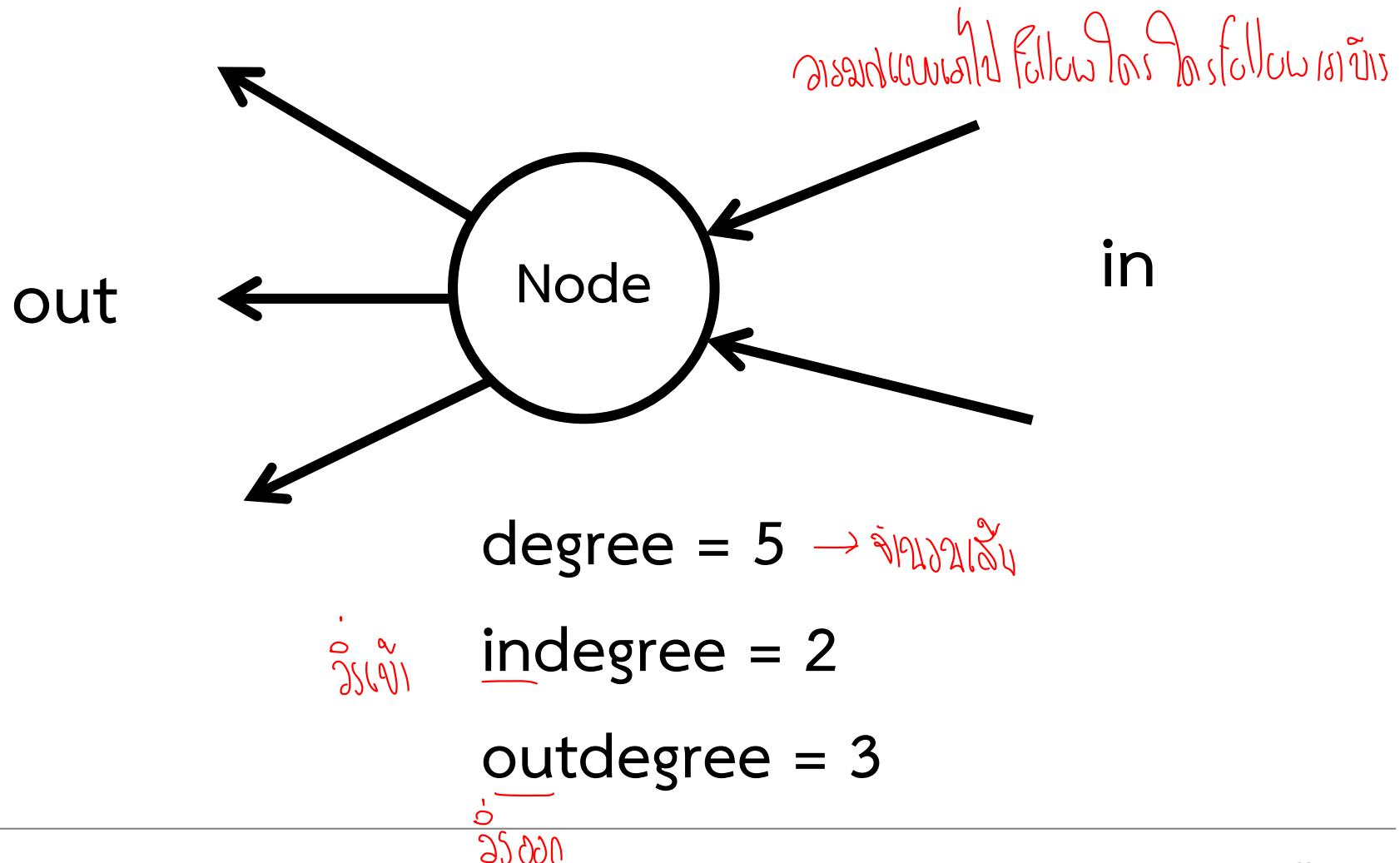
<http://www.orbifold.net/Mathematica/Networks/>



ອົນດະວຸນທະນາຄານ

Properties

Degree

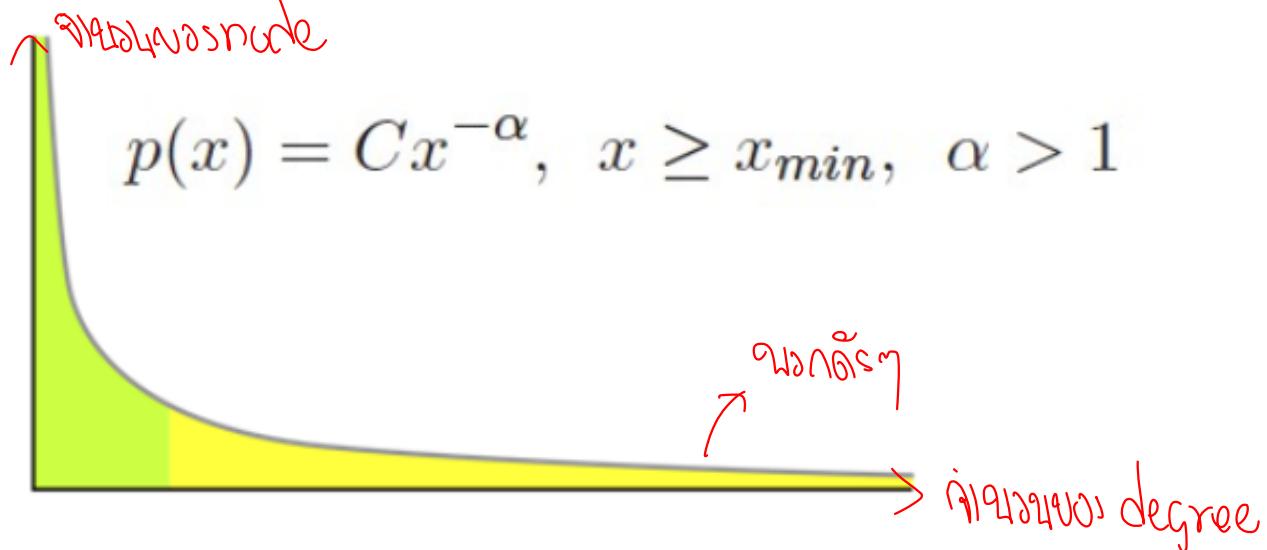


Degree

- Min Degree
- Max Degree
- Average Degree

Degree Distribution

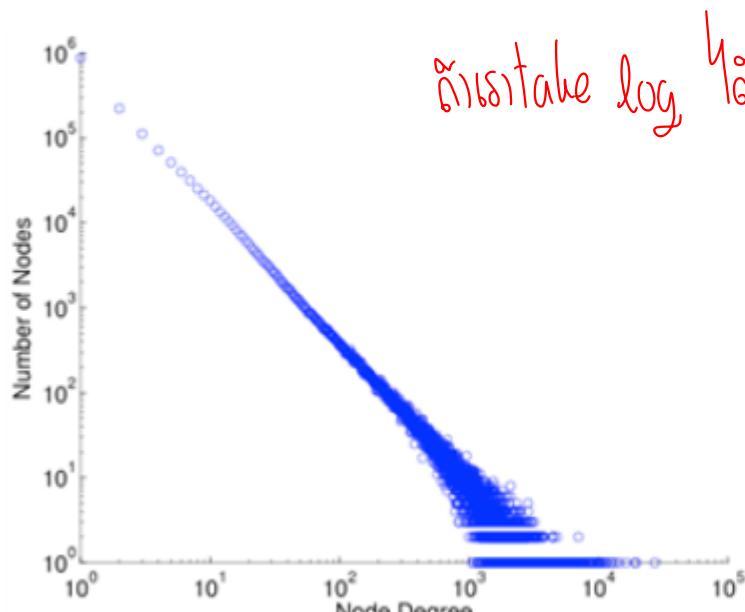
- | Degree distribution in large-scale networks often follows a **power law**, that is, the fraction $p(x)$ of nodes in the network having x connections to other nodes goes for large values of x as:



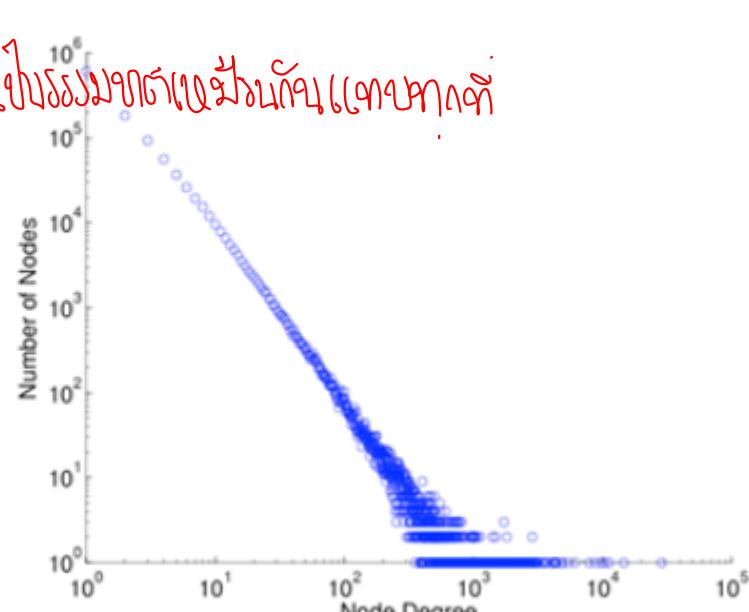
- | A.k.a. **long tail** distribution, **scale-free** distribution

Log-Plot

- | Power law distribution becomes a **straight line** if plotted in a log-log scale



Friendship Network in Flickr

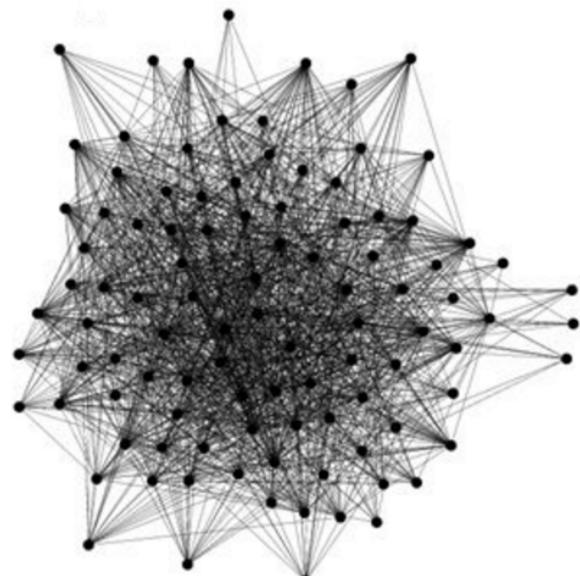


Friendship Network in YouTube

Dense/Sparse

អ្នកគេសរសាងលើមួយ

វេងទឹង



Dense Graph

បញ្ហា



Sparse Graph

Density

- In mathematics, a dense graph is a graph in which the number of edges is close to the maximal number of edges. The opposite, a graph with only a few edges, is a sparse graph. The distinction between sparse and dense graphs is rather vague, and depends on the context.

Undirected Graph

ໜີນກວດຕອບລາຍງາແນະນຳ

$$D = \frac{2|E|}{|V|(|V|-1)}$$

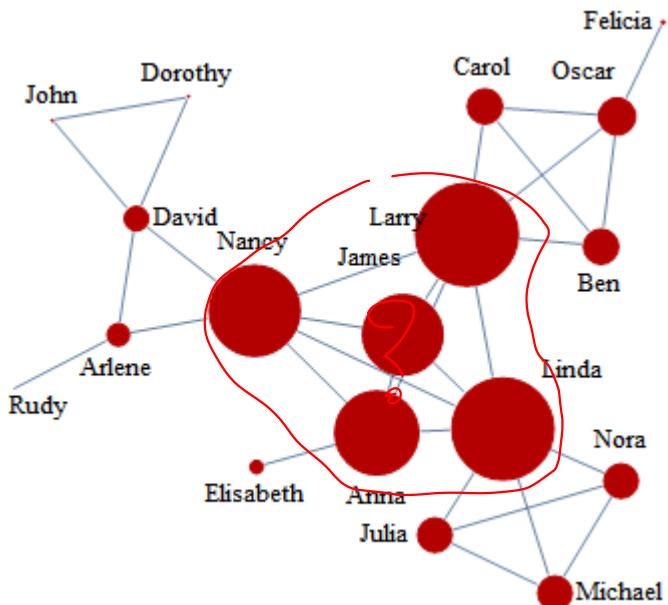
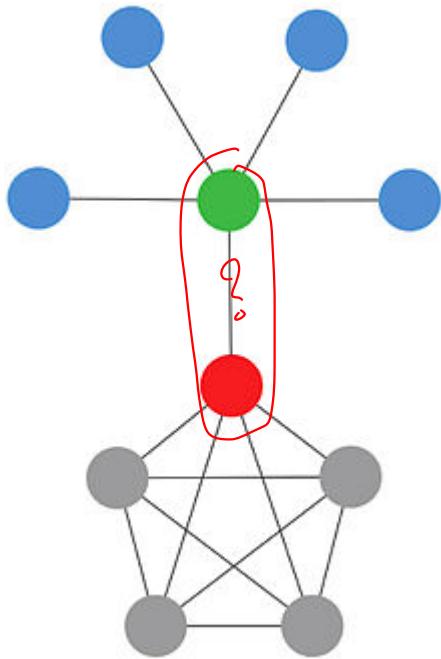
Directed Graph

$$D = \frac{|E|}{|V|(|V|-1)}$$

ຈົດວິທີ
ຈົດວິທີ
ຈົດວິທີ

Centrality

សំណើភាពនៃអ្នកគេ



Closeness Centrality

ສະກອນ

- In a connected graph, the normalized closeness centrality (or closeness) of a node is the average length of the shortest path between the node and all other nodes in the graph. Thus the more central a node is, the closer it is to all other nodes. Closeness was defined by Bavelas (1950) as the reciprocal of the farness, that is:

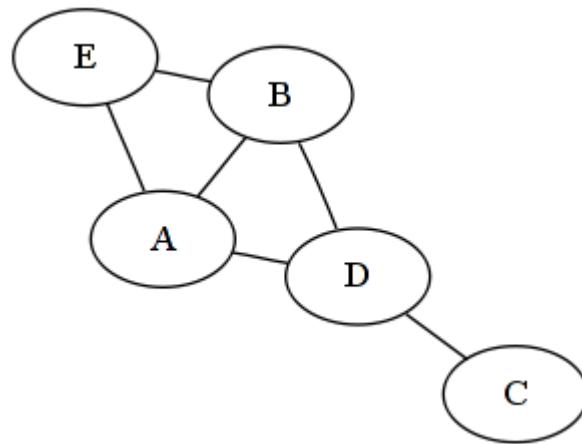
$$C(x) = \frac{N - 1}{\sum_y d(y, x)}$$

▷ *node x*
▷ *node y*
▷ *distance between node x and y*

- where $d(x,y)$ is the distance between vertices x and y . However, when speaking of closeness centrality, people usually refer to its normalized form, generally given by the previous formula multiplied by $N-1$, where N is the number of nodes in the graph. This adjustment allows comparisons between nodes of graphs of different sizes.
- Taking distances from or to all other nodes is irrelevant in undirected graphs, whereas it can produce totally different results in directed graphs (e.g. a website can have a high closeness centrality from outgoing link, but low closeness centrality from incoming links).

Closeness Centrality

$$C(x) = \frac{N - 1}{\sum_y d(y, x)}$$



$$C(A) = \frac{5 - 1}{d(B, A) + d(C, A) + d(D, A) + d(E, A)}$$

ສິນາກີ່ມີປຸງ node ຈາກ A

$$C(A) = \frac{4}{1 + 2 + 1 + 1} = \frac{4}{5} = 0.8$$

ຢືນຢັນວ່າມີເປົ້າໃຫຍ່
ຢືນຢັນວ່າມີເປົ້າໃຫຍ່

Other Centrality

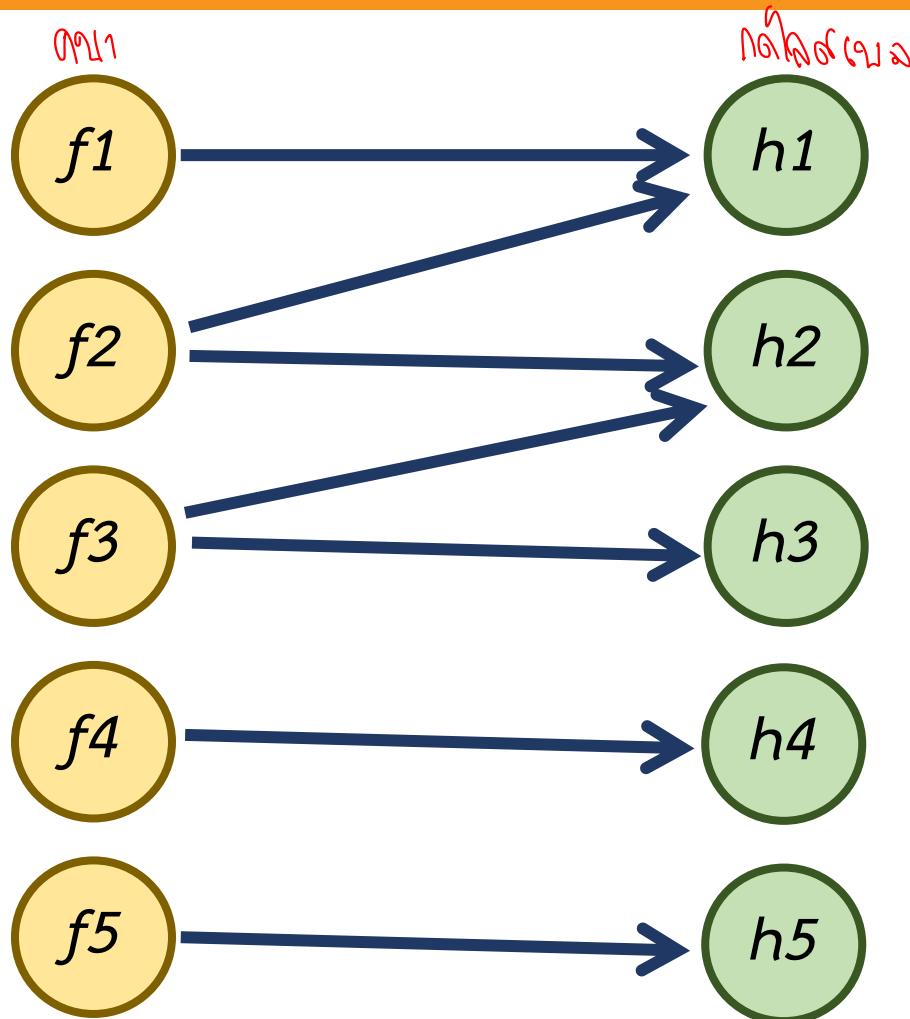
வார்குடி $\rightarrow \frac{x; x = \text{constant}}{\square; \text{distance}}$

இல்லை விடையளவு

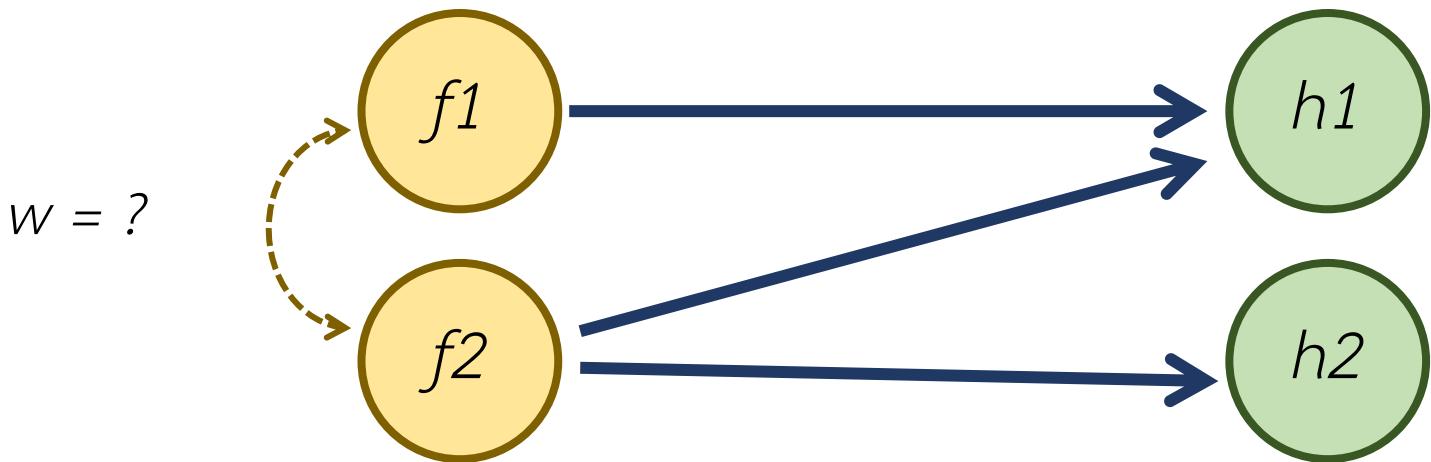
- Degree centrality $C(7) = 13$
- Closeness Centrality $C(1) = 11$
- Harmonic centrality $C(3)$
- Betweenness centrality $C(4)$
- Eigenvector centrality $C(0)$
- Katz centrality
- PageRank centrality
- Percolation centrality
- Cross-clique centrality
- Freeman Centralization

Similarity

အသေစိန်များ
မှတ်စွမ်းမှု
is bipartite graph
ပေါင်းပေါင်းလဲ



Similarity



Jaccard Index → නිශ්චත්‍යා නිශ්චත්‍යා සංඛ්‍යා මූලික පෙන්වනුයේ

$$w(f1, f2) = \frac{|\Gamma(f1) \cap \Gamma(f2)|}{|\Gamma(f1) \cup \Gamma(f2)|} = \frac{|\{h1\} \cap \{h1, h2\}|}{|\{h1\} \cup \{h1, h2\}|} = \frac{|\{h1\}|}{|\{h1, h2\}|}$$
$$= \frac{1}{2} \quad = 0.5$$

return set of f_n නොවා තැබ්දු ඇත්තු

Similarity Indices

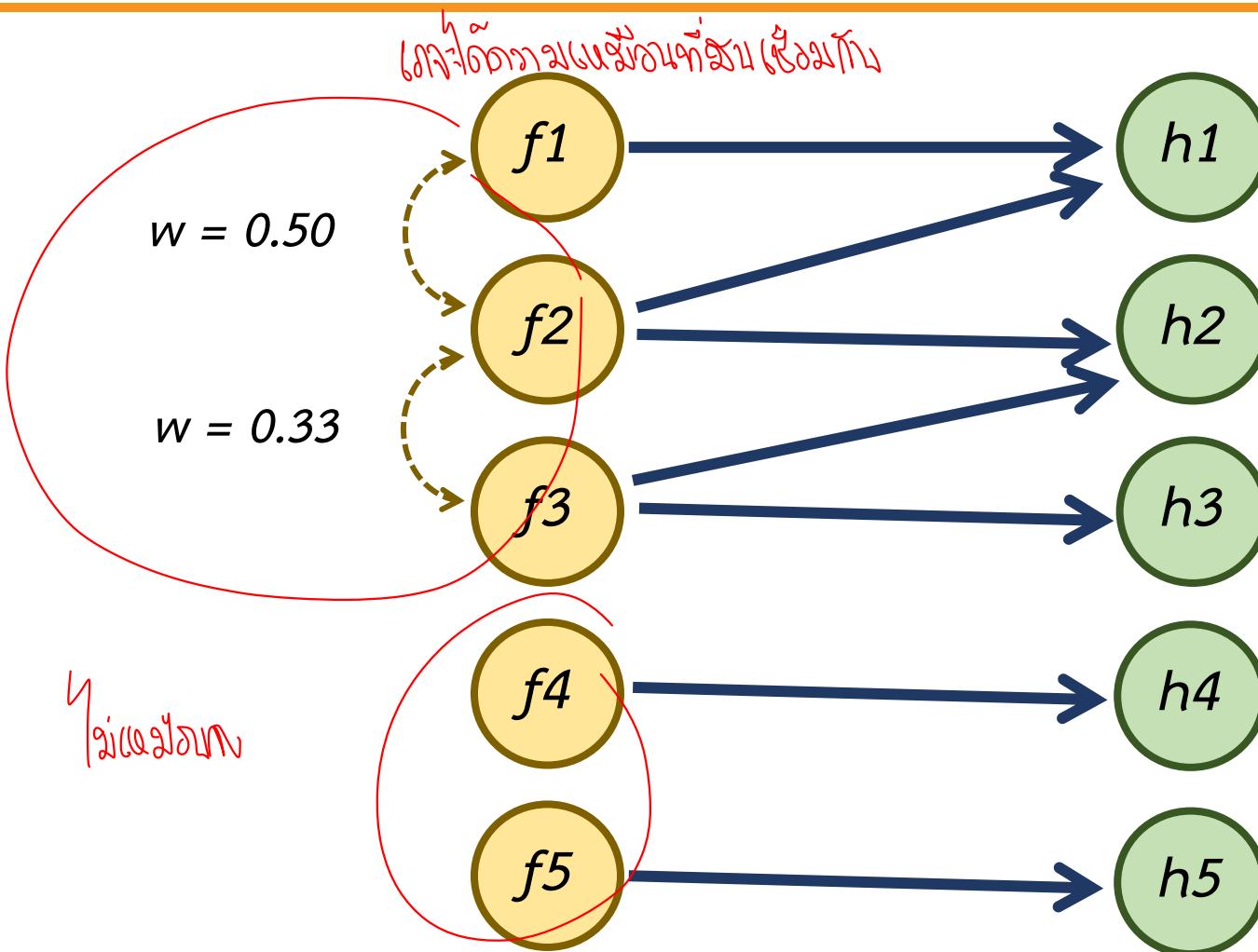
118-2018

$\Gamma(n)$ is a function that returns a set of nodes that interact with the node n .

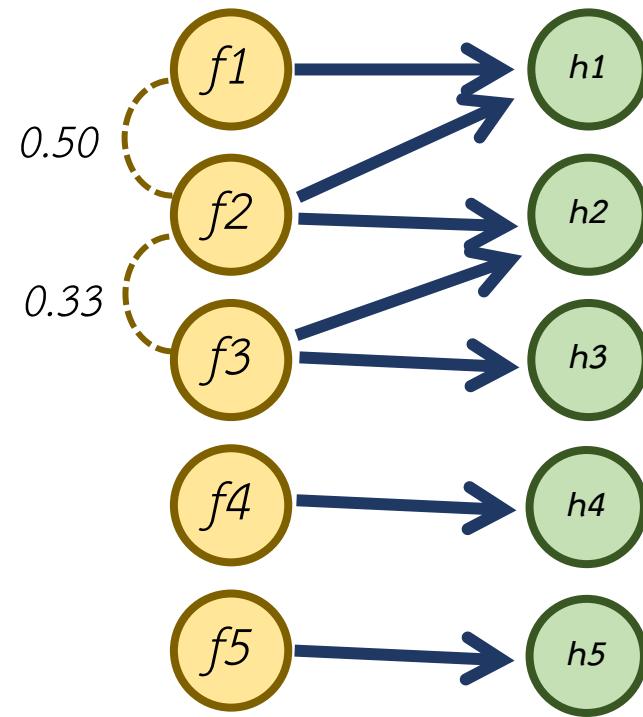
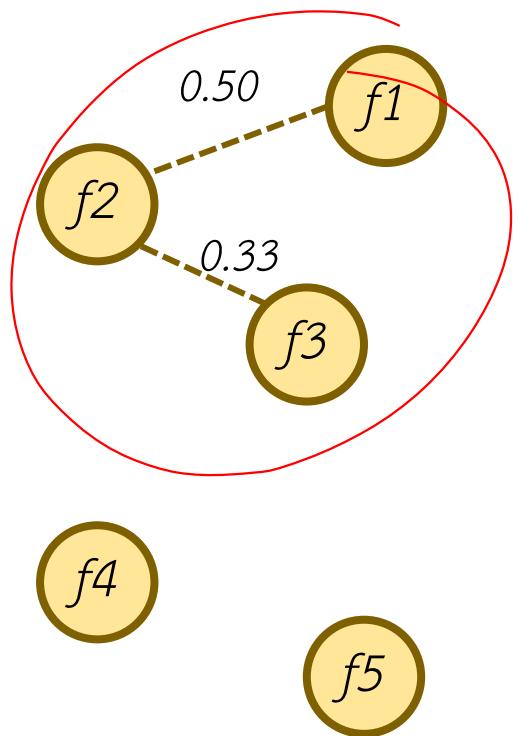
Example: $\Gamma(f2) = \{ h1, h2 \}$

- Common Neighbors (CN) : $|\Gamma(x) \cap \Gamma(y)|$
- Jaccard Index:
$$\frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}$$
- Sørensen index:
$$\frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x)| + |\Gamma(y)|}$$
- Hub Depressed Index (HDI):
$$\frac{|\Gamma(x) \cap \Gamma(y)|}{\max(\Gamma(x), \Gamma(y))}$$
- Resource Allocation Index (RA) :
$$\sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{|z|}$$

Similarity



Projection



Projection of F

Bipartite Graph

Community Detection

diagramm cluster von network

Modularity Maximization

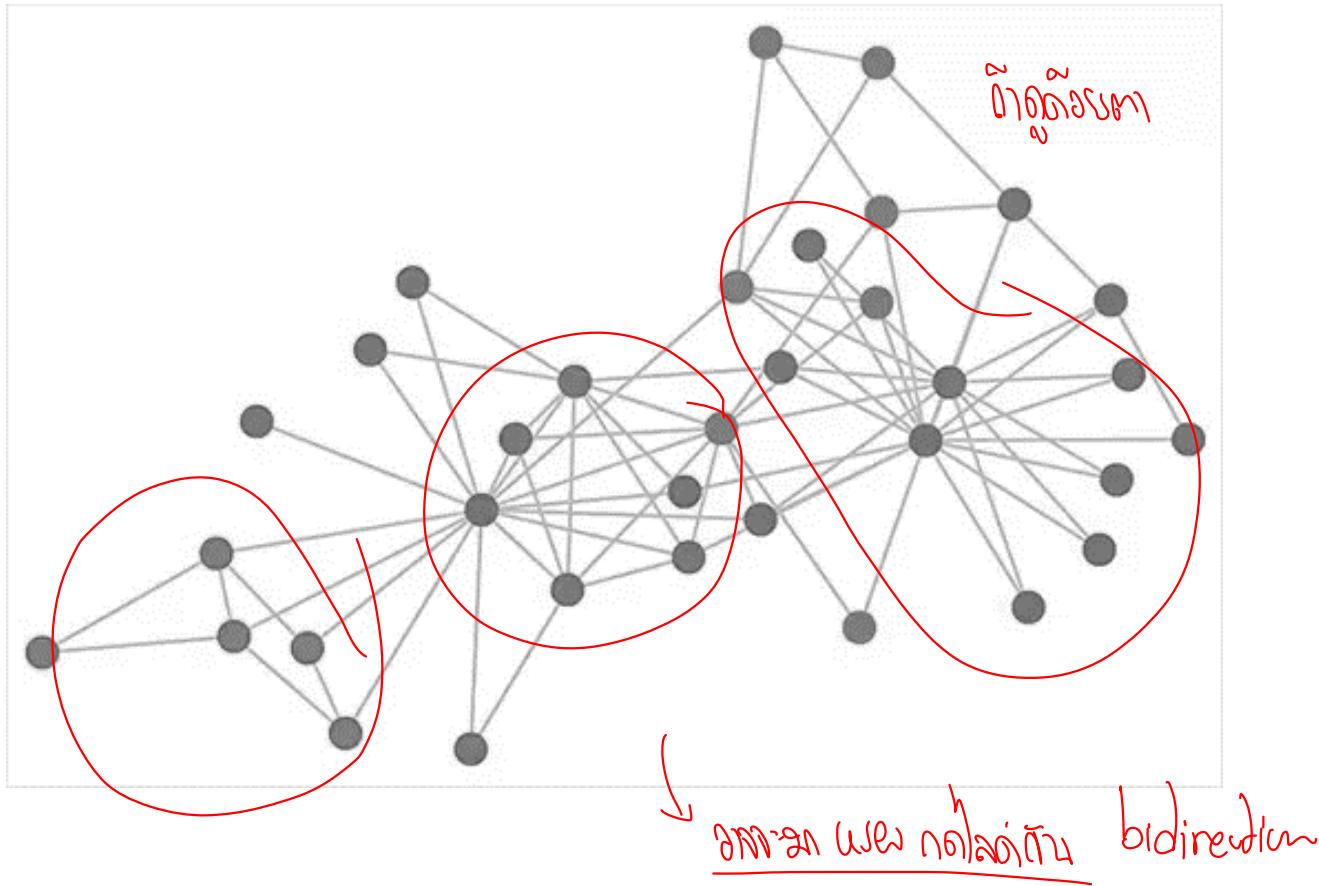
ការអនុវត្តន៍វាង network

នូវការអនុវត្តន៍

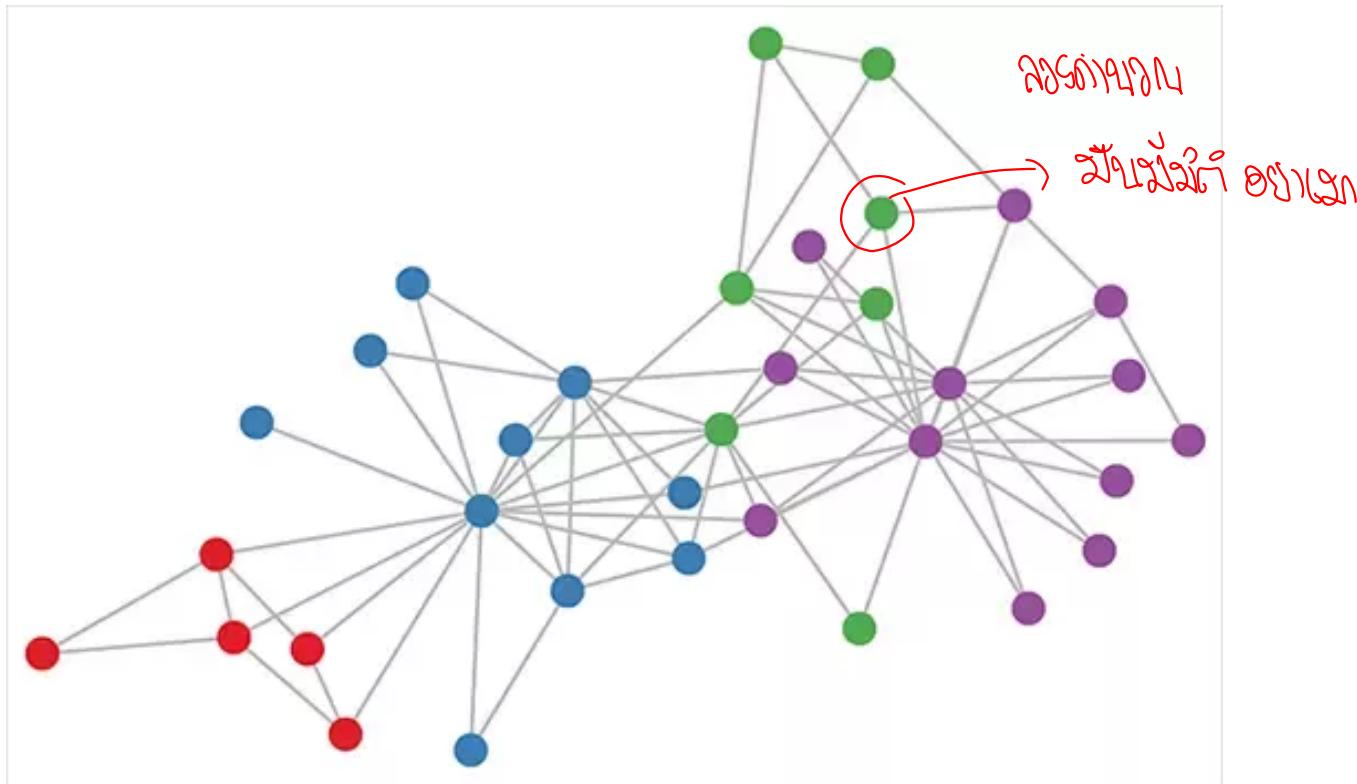
រួចរាល់, រួចរាល់

- Modularity is a **measure** of the network structure.
- It was designed to **evaluate the strength of division of a network** into modules (also called groups, clusters or communities).
- Networks with high modularity have **dense connections** between the nodes within modules but sparse connections between nodes in different modules.
- Modularity is often used in **optimization methods** for detecting community structure in networks
- A simple calculation is good for **unweighted** and **undirected** graphs

Community Structure in SNA



Community Structure in SNA



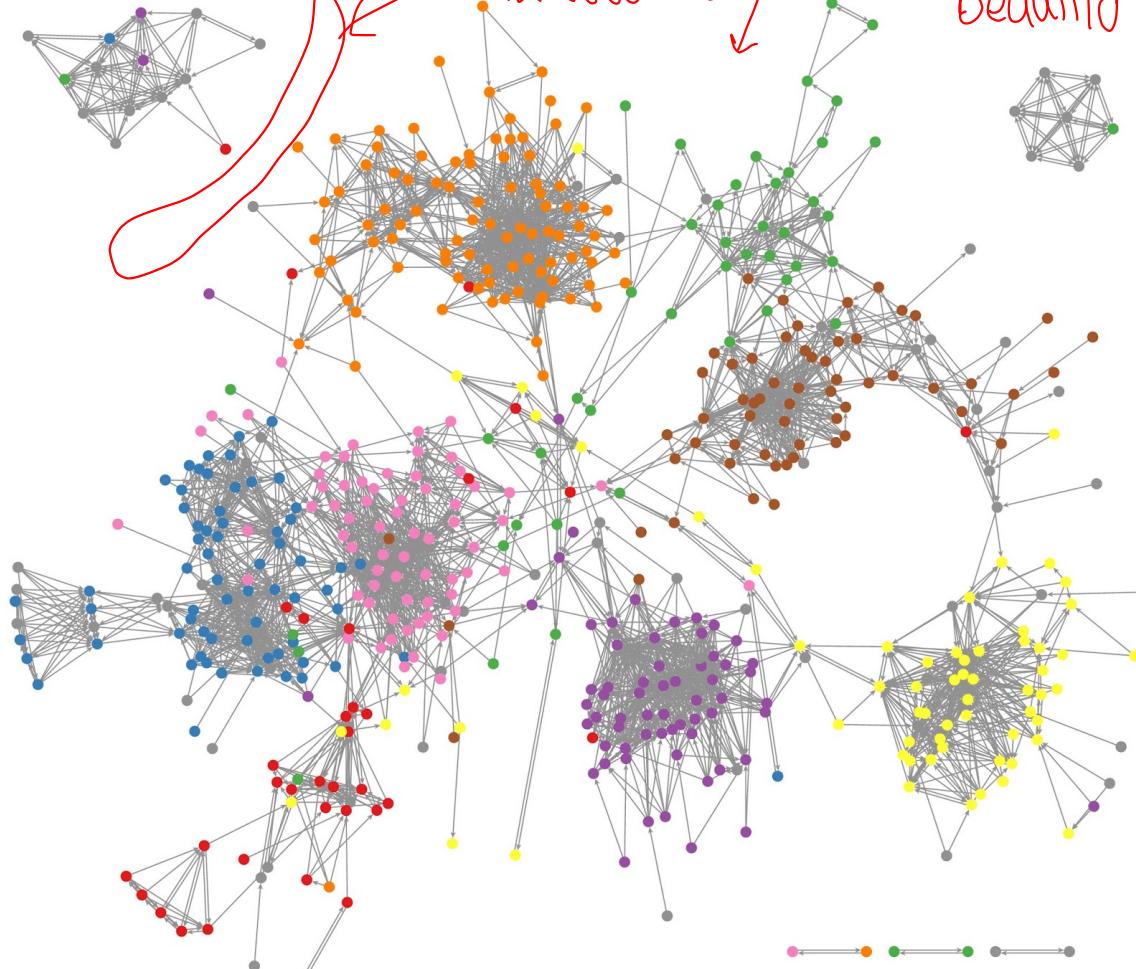
Community Structure in SNA

շաբաթը՝
կամ bi part

շաբաթ = պինդաց տարր

կամ եղանակ

beautiful



Equation

សមត្ថធម៌

anomaly

វិធាននៅលើខាងក្រោម

$$Q = \frac{1}{2m} \sum_{i,j \in C} \left(A_{i,j} - \frac{d_i d_j}{2m} \right)$$

Equation

$$Q = \frac{1}{2m} \sum_{i,j \in C} \left(A_{i,j} - \frac{d_i d_j}{2m} \right)$$

i กับ j อยู่ในกลุ่มเดียวกัน ดังนั้น $A_{i,j}$ จะเท่ากับ 1

Every pairs of nodes i and j being in the same community

Equation

$$Q = \frac{1}{2m} \sum_{i,j \in C} \left(A_{i,j} - \frac{d_i d_j}{2m} \right)$$

m is the number of links *અપાનોએડજ નુંથું*

Equation

$$Q = \frac{1}{2m} \sum_{i,j \in C} \left(A_{i,j} - \frac{d_i d_j}{2m} \right)$$



Adjacency matrix

If having link, $A_{i,j} = 1$

If having no link, $A_{i,j} = 0$

Equation

$$Q = \frac{1}{2m} \sum_{i,j \in C} \left(A_{i,j} - \frac{d_i d_j}{2m} \right)$$



probability a random edge would go between i and j

ຕາງໝາຍຂີບທີ່ຈະປັ້ງສະກັນໄວ້

Equation

d is a degree of a node == number of edges connected

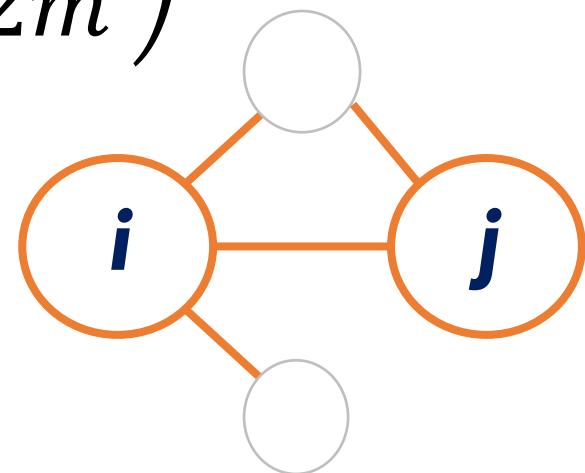
$$Q = \frac{1}{2m} \sum_{i,j \in C} \left(A_{i,j} - \frac{d_i d_j}{2m} \right)$$

↑↑↑↑ degree varnished

In this case,

$$d_i = 3$$

$$d_j = 2$$

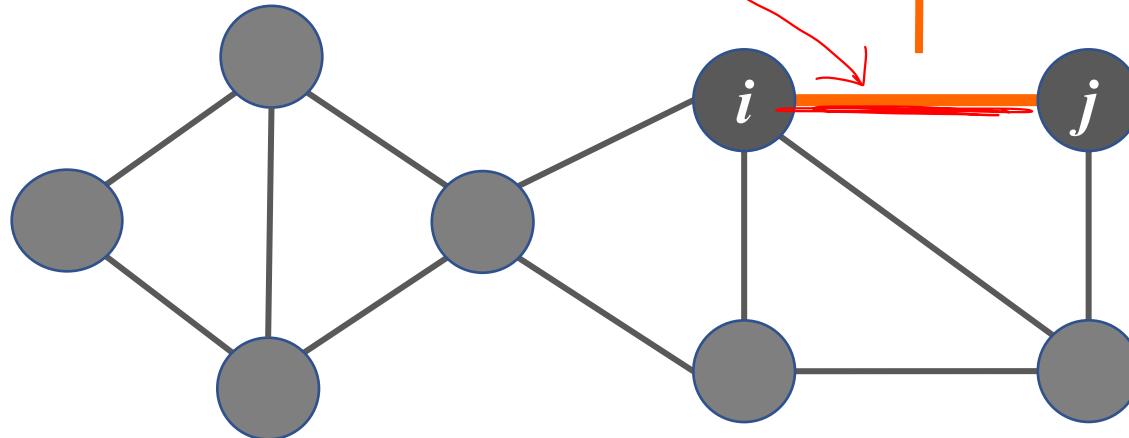


How To Calculate

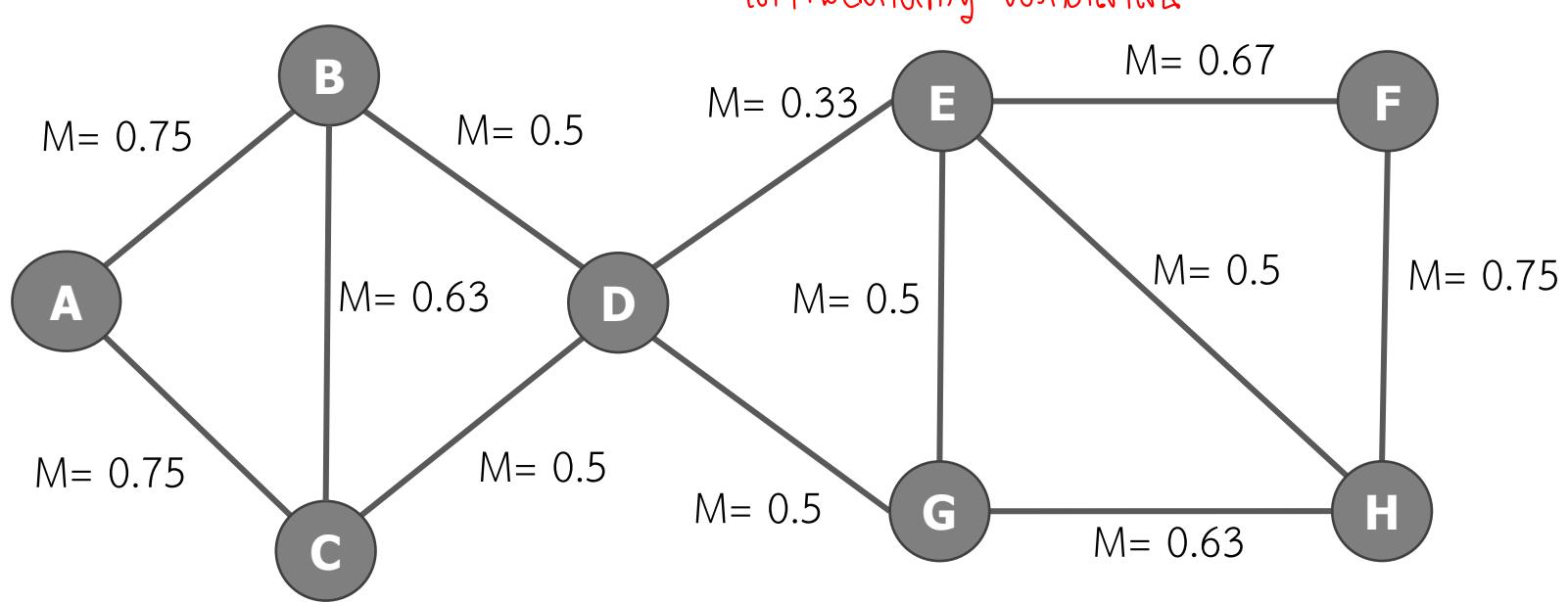
$$A_{i,j} - \frac{d_i d_j}{2m} = 1 - \frac{4 \times 2}{2 \times 12} = 0.67$$

d_i d_j
m = number of edges

12 edges
m = 12

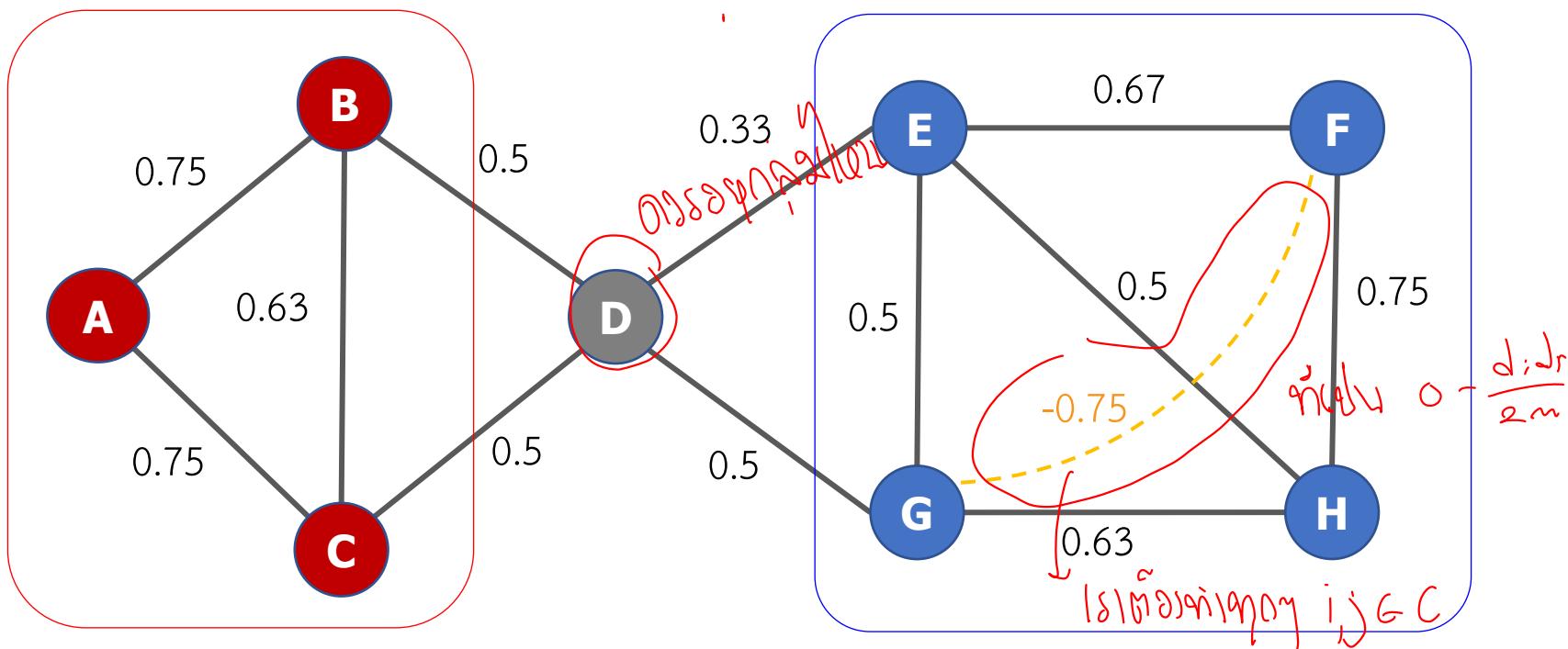


After Calculate



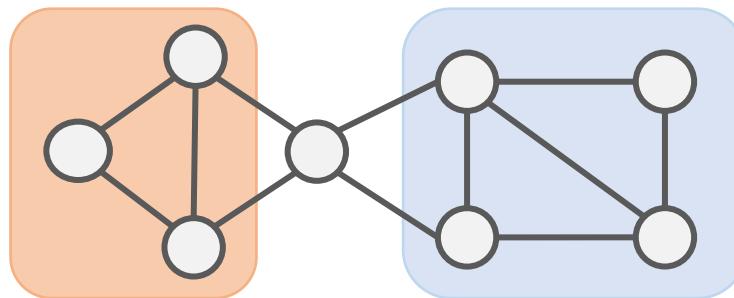
All edges have own Modularity Score

After Calculate



$$\begin{aligned}
 Q &= ((0.75 + 0.75 + 0.63) + (0.5 + 0.5 + 0.63 + 0.67 + 0.75 - 0.75)) / 2m \\
 &= 0.184
 \end{aligned}$$

Next Iteration

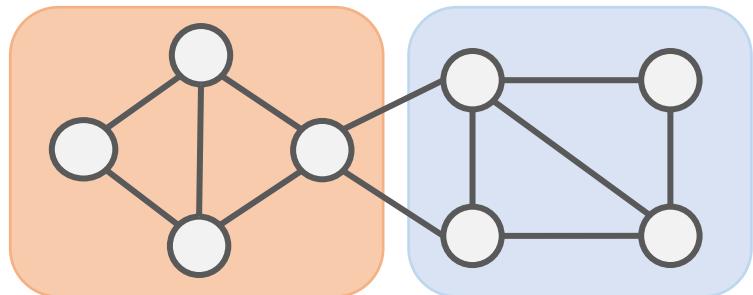
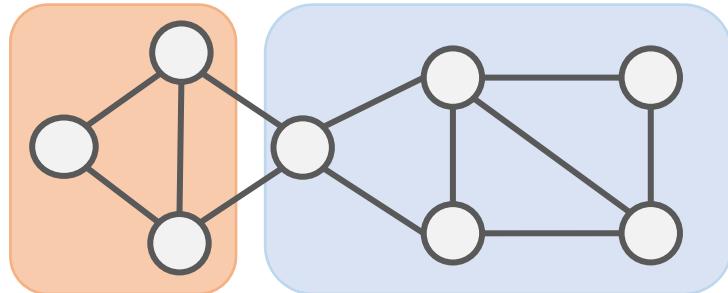


(A)

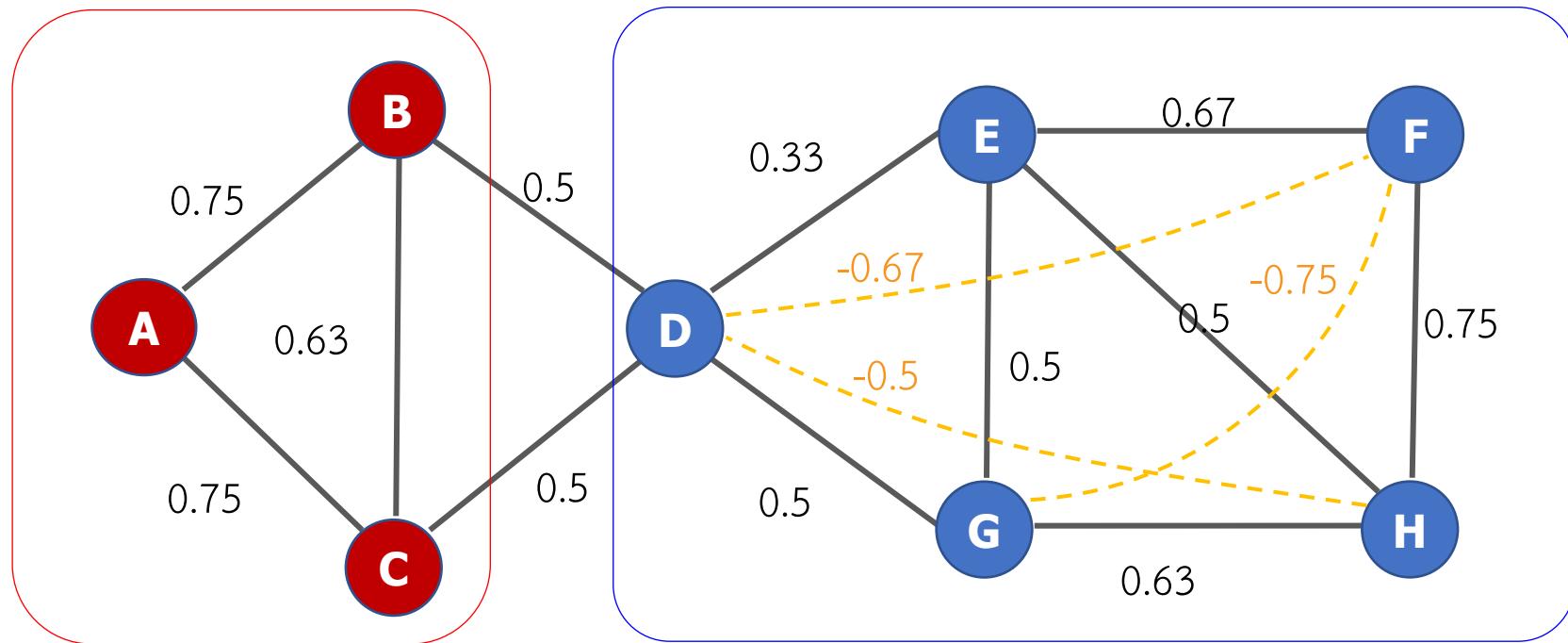


↓ ↓ 1110

(B)



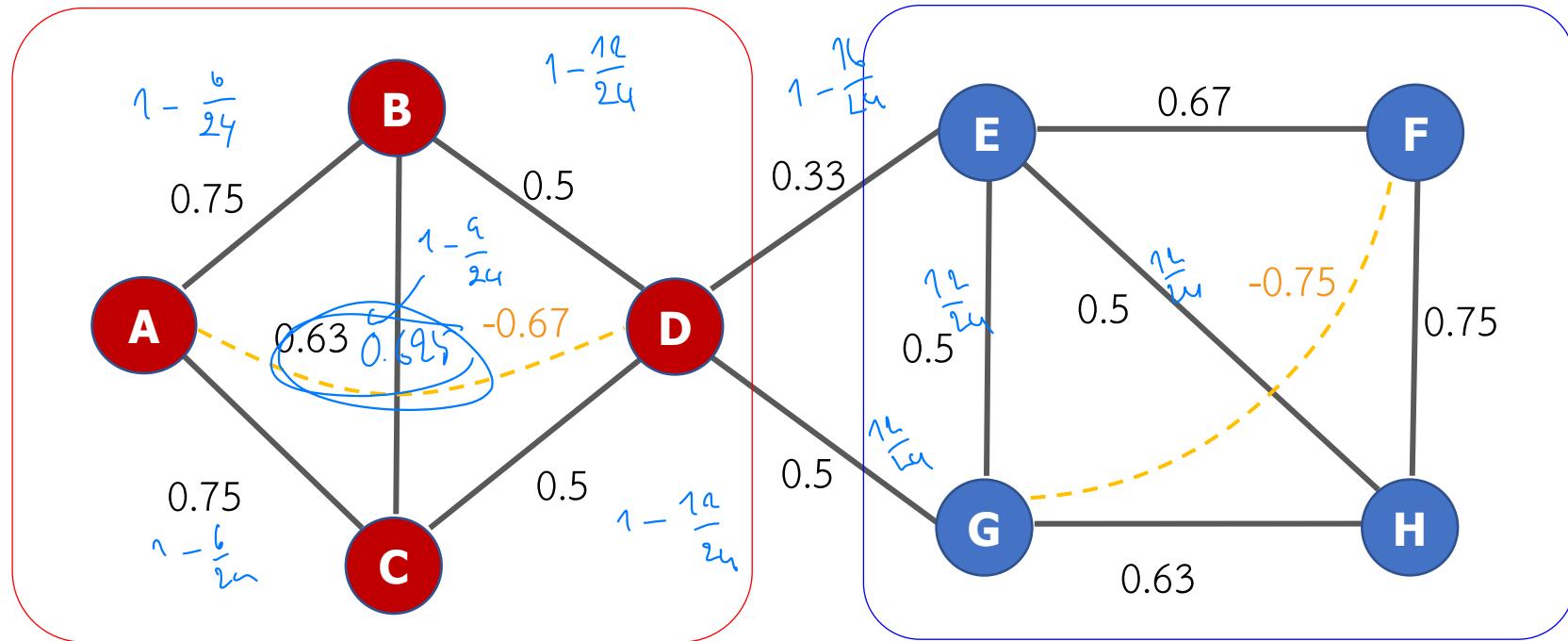
Candidate - A



$$\begin{aligned} Q &= ((0.75 + 0.75 + 0.63) + (0.5 + 0.5 + 0.63 + 0.67 + 0.75 + 0.33 + 0.5 - 0.5 - 0.67 - 0.75)) / 2m \\ &= 0.170 \end{aligned}$$

(ব্যবস্থাপনার মাধ্যমে)

Candidate - B



$$Q = ((0.75 + 0.75 + 0.63 + 0.5 + 0.5 - 0.67) + (0.5 + 0.5 + 0.63 + 0.67 + 0.75 - 0.75)) / 2m$$

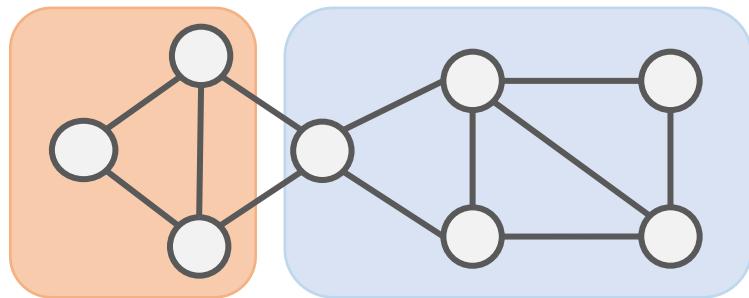
$$= 0.198$$

பொதுமான ப
யார்யாவுக்கு

(B)

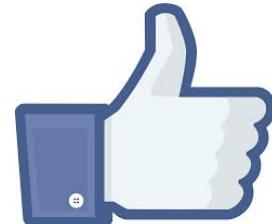
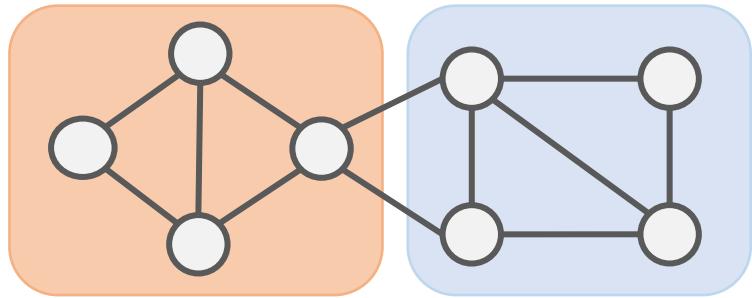
We Choose

(A)

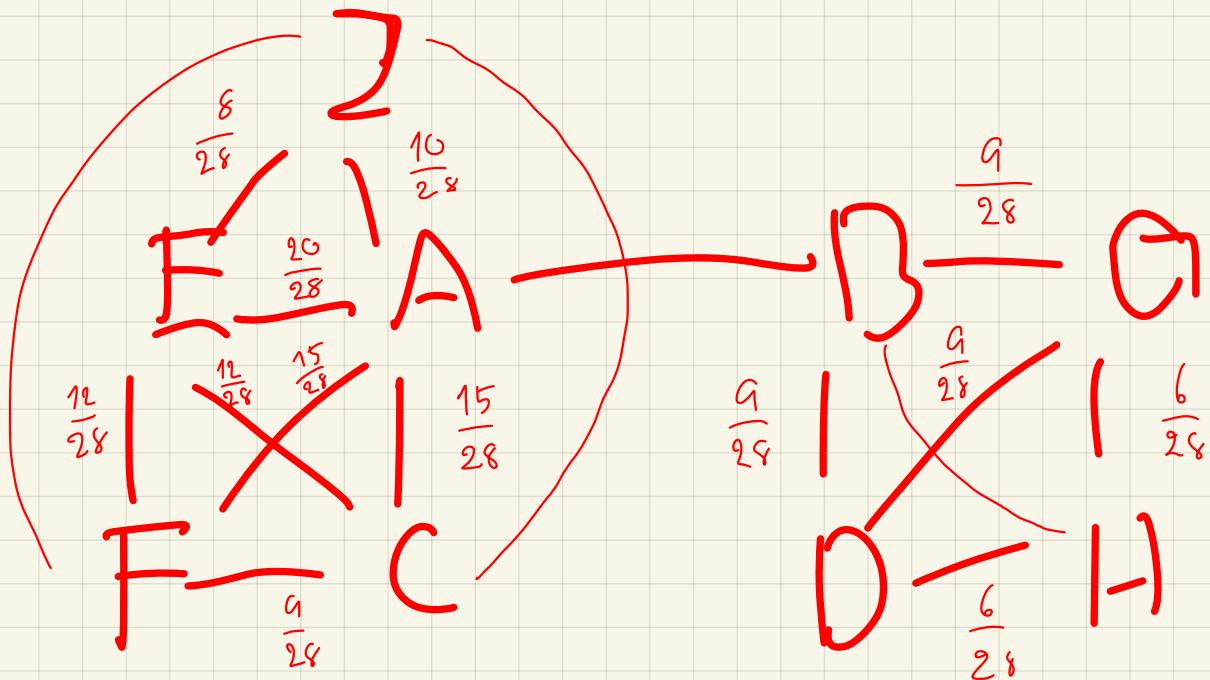


$$Q = 0.170$$

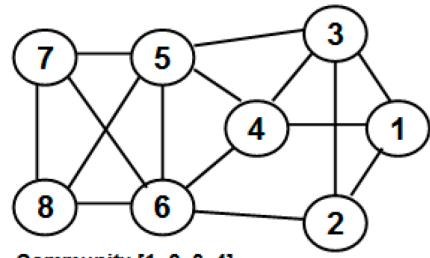
(B)



$$Q = 0.198$$



Other Example



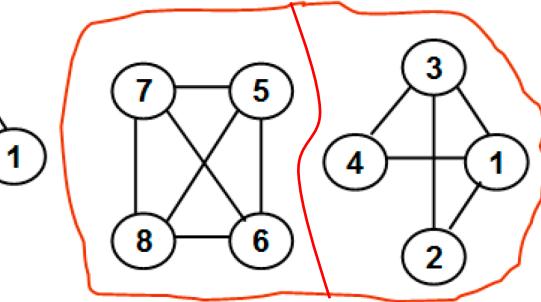
Community [1, 2, 3, 4]

Edges with $A_{ij} = 1$ Modularity

1 - 2	$1 - (3)(3)/(2^*15) = 0.70$
1 - 3	$1 - (3)(4)/(2^*15) = 0.60$
1 - 4	$1 - (3)(4)/(2^*15) = 0.60$
2 - 3	$1 - (3)(3)/(2^*15) = 0.70$
3 - 4	$1 - (4)(4)/(2^*15) = 0.47$
Edges with $A_{ij} = 0$	
2 - 4	$0 - (3)(4)/(2^*15) = -0.40$

Total Modularity Score for

Community [1, 2, 3, 4]



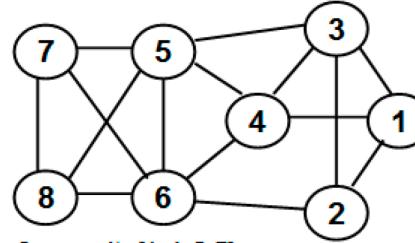
Community [5, 6, 7, 8]

Edges with $A_{ij} = 1$ Modularity

5 - 6	$1 - (5)(5)/(2^*15) = 0.17$
5 - 7	$1 - (3)(5)/(2^*15) = 0.50$
5 - 8	$1 - (3)(5)/(2^*15) = 0.50$
6 - 7	$1 - (3)(5)/(2^*15) = 0.50$
6 - 8	$1 - (3)(5)/(2^*15) = 0.50$
7 - 8	$1 - (3)(3)/(2^*15) = 0.70$

Total Modularity Score for
Community [2, 3, 6, 8]

ເຕີມກຳນົດໄວ້

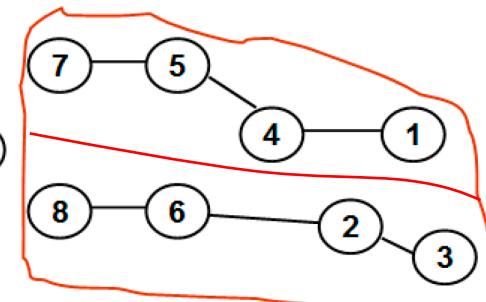


Community [1, 4, 5, 7]

Edges with $A_{ij} = 1$ Modularity

1 - 4	$1 - (3)(4)/(2^*15) = 0.60$
4 - 5	$1 - (4)(5)/(2^*15) = 0.33$
5 - 7	$1 - (3)(5)/(2^*15) = 0.50$
Edges with $A_{ij} = 0$	
1 - 5	$0 - (3)(5)/(2^*15) = -0.50$
1 - 7	$0 - (3)(3)/(2^*15) = -0.30$
4 - 7	$0 - (4)(3)/(2^*15) = -0.40$

Modularity Score for
Community [1, 4, 5, 7]



Community [2, 3, 6, 8]

Edges with $A_{ij} = 1$ Modularity

2 - 3	$1 - (3)(4)/(2^*15) = 0.60$
2 - 6	$1 - (3)(5)/(2^*15) = 0.50$
6 - 8	$1 - (3)(5)/(2^*15) = 0.50$
Edges with $A_{ij} = 0$	
2 - 8	$0 - (3)(3)/(2^*15) = -0.30$
3 - 6	$0 - (4)(5)/(2^*15) = -0.67$
3 - 8	$0 - (4)(3)/(2^*15) = -0.40$

Modularity Score for
Community [2, 3, 6, 8]

Cumulative Modularity Score for the two Communities: $2.67 + 2.87 = 5.54$

(a)

$$Q = 5.54/2m$$

(b)

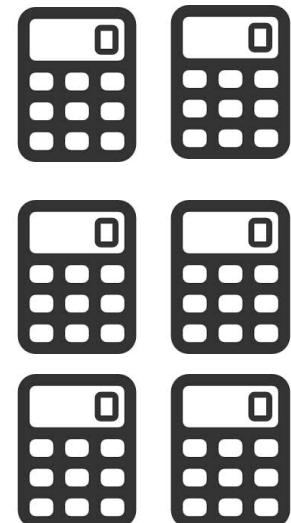
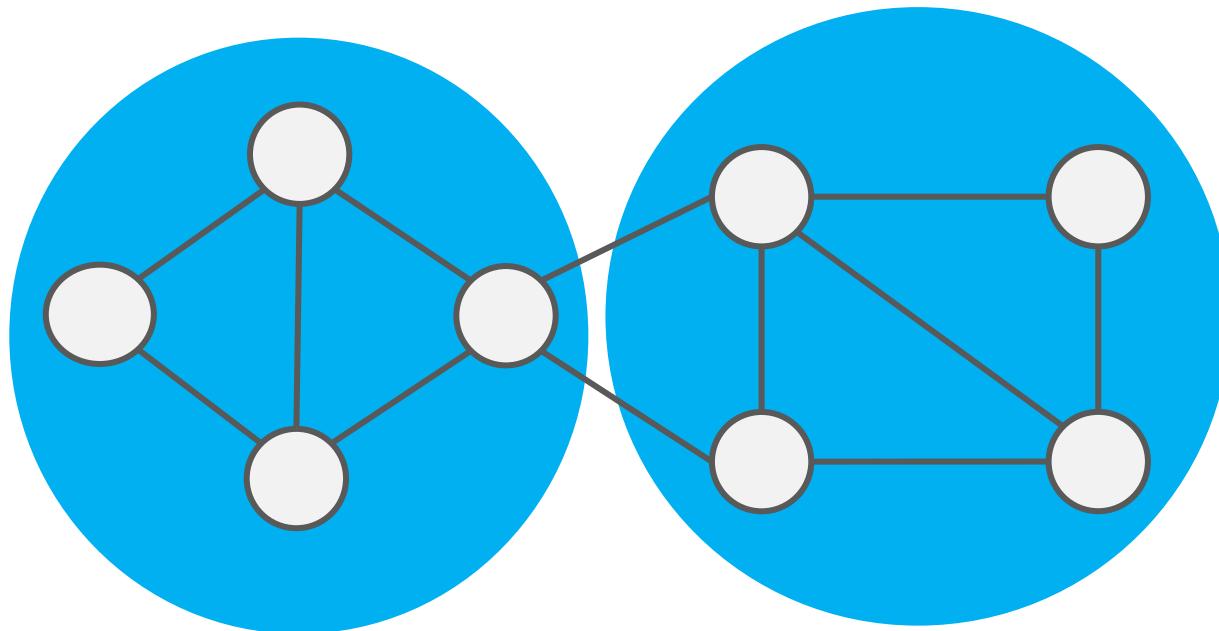
$$Q = 0.44/2m$$

Walk Trap

ແລ້ວ
ມີສົດສະບັບ

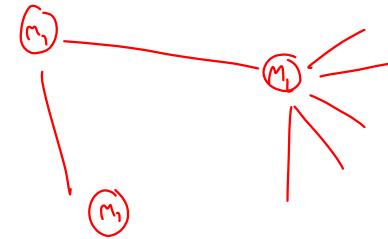
Random Walk

- Calculate **Whole Modularity** in every iteration



“

ឧច្ចាស់ការងារ (បន្ទីមបង្ហាញ)
ទំនាក់ទំនង



Every once in a while, a new technology, an old problem, and a big idea turn into an innovation.

វឌ្ឍន៍ដែលការពារ
Max mxm នៃ

(បន្ទីមបង្ហាញ) visualize graph នូវ Gephi → Java & Google & Free

Dean Kamen
បិន្តុ

”

つづく