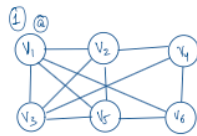


HW #5 - Solution

1. Given the following adjacency matrix

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

- Provide the maximal 2-cliques for this graph.
- Provide the maximal 2-plexes for this graph.
- Are there any structurally equivalent nodes in this graph? If yes, list them.



All nodes are within distance 2 to any other node in the graph.

So, Maximal 2-clique is $\{v_1, v_2, v_3, v_4, v_5, v_6\}$

② For a graph to be k -plex every vertex should have a degree of at least $|V| - k$, where $|V|$ is number of vertices in the graph.

for entire graph the minimum degree is 3 (for v_4, v_6)

So,

$$\Rightarrow |V| - k \leq 3$$

$$\Rightarrow 6 - k \leq 3$$

$$\Rightarrow k \geq 3 \text{ (3-plex graph)}$$

So, the entire graph cannot be a 2-plex

Considering for $|V|=5$

removing any node results in minimum degree of 2 (again for v_4, v_6)

$$\Rightarrow 5 - k \leq 2$$

$$\Rightarrow k \geq 3 \text{ (again 3-plex)}$$

So, node size of 5 cannot be a 2-plex either.

for Node size = 4,

for all the below subgraphs the minimum degree = 2

1. v_1, v_2, v_3, v_4 6. v_1, v_3, v_5, v_6

2. v_2, v_4, v_5, v_6 7. v_1, v_2, v_5, v_6

3. v_2, v_3, v_5, v_6 8. v_3, v_4, v_5, v_6

4. v_1, v_2, v_5, v_6 9. v_1, v_2, v_3, v_4

5. v_1, v_3, v_4, v_6

$$\Rightarrow |V| - k \leq 2$$

$$\Rightarrow k \geq 2 \text{ (2-plex)}$$

So, all the above subgraphs are 2-plex.

Any connected subgraph of node size 3, 4, 1 will also be 2-plex, however only the listed subgraphs of node size 4 are considered "maximal" 2-plex.

③ Yes, $\{v_1, v_2\}$ & $\{v_5, v_6\}$ have same neighbours,

So, they are structurally equivalent.

OR 1(c)

Two nodes are structurally equivalent when they have the exact same neighbors in the graph

- Node 1 neighbors are 2 3 5 6
- Node 2 neighbors are 1 3 4 5
- Node 3 neighbors are 1 2 4 5
- Node 4 neighbors are 2 3 6
- Node 5 neighbors are 1 2 3 6
- Node 6 neighbors are 1 4 5

No two nodes share the same neighbor set

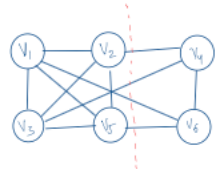
➔ **Therefore no pair of nodes is structurally equivalent**

2. Using the graph from problem 1,

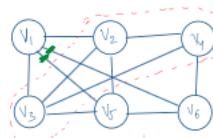
- Find a minimum cut that creates partitions P_1 and P_2 where $|P_1|=4$ and $|P_2|=2$.**
- Find a minimum cut that creates partitions P_1 and P_2 where $|P_1| = |P_2| = 3$.**
- For each of the cuts in a and b above, calculate**
 - Ratio Cut (P)**
 - Normalized Cut (P)**
- Which cut is preferable based on the above metrics?**

2

- (a) for Partitions $|P_1|=4$ and $|P_2|=2$
 Minimum Cut = 4
 partition \Rightarrow
 $P_1: V_1, V_2, V_3, V_5$
 $P_2: V_4, V_6$

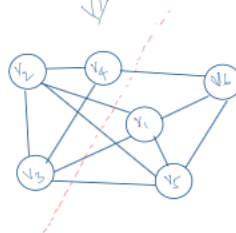


- (b) for Partitions $|P_1|=|P_2|=3$
 Minimum Cut = 5
 partition \Rightarrow
 $P_1: V_2, V_3, V_4$
 $P_2: V_1, V_5, V_6$



\because Edges marked in green will not be cut.

(c) for Better Visualization & Understanding



- (i) ratio cut & Normalized cut for $|P_1|=4$, $|P_2|=2$

$$\Rightarrow \text{ratio cut} = \frac{1}{2} \left(\frac{4}{4} + \frac{4}{2} \right) = \frac{1}{2} (1+2) = \frac{3}{2}$$

$$\Rightarrow \text{Normalized cut} = \frac{1}{2} \left(\frac{4}{16} + \frac{4}{2} \right) = \frac{1}{2} \left(\frac{1}{4} + \frac{1}{2} \right) = \frac{3}{4}$$

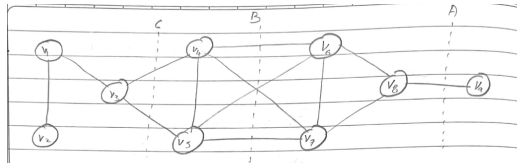
- (ii) ratio cut & Normalized cut for $|P_1|=|P_2|=3$

$$\Rightarrow \text{ratio cut} = \frac{1}{2} \left(\frac{5}{3} + \frac{5}{3} \right) = \frac{5}{3}$$

$$\Rightarrow \text{Normalized cut} = \frac{1}{2} \left(\frac{5}{9} + \frac{5}{9} \right) = \frac{5}{9}$$

- (d) \Rightarrow for ratio cut 4-2 split should be preferred since it is lesser
 \Rightarrow for Normalized cut 3-3 split, since it is lesser than 4-2.

3. Social Media Mining (SMM) Ch 6, problem 7



Now here we have to compute Jaccard and cosine similarity between nodes v_4 and v_8 .

We have to assume that neighbour of a node excludes the node itself.

The node excludes itself :-

$$\therefore N(v_4) = \{v_3, v_5, v_6, v_7\}$$

$$\therefore N(v_8) = \{v_6, v_7, v_9\}$$

We will check for Jaccard

$$J(v_i, v_j) = \frac{|N(v_i) \cap N(v_j)|}{|N(v_i) \cup N(v_j)|}$$

$$J(v_4, v_8) = \frac{|N(v_4) \cap N(v_8)|}{|N(v_4) \cup N(v_8)|}$$

$$= \frac{|\{v_3, v_5, v_6, v_7\} \cap \{v_6, v_7, v_9\}|}{|\{v_3, v_5, v_6, v_7\} \cup \{v_6, v_7, v_9\}|}$$

$$= \frac{|\{v_6, v_7\}|}{|\{v_3, v_5, v_6, v_7, v_9\}|}$$

$$= \frac{2}{5} \approx 0.4$$

Now let's find the cosine

$$C(v_i, v_j) = \frac{|N(v_i) \cap N(v_j)|}{\sqrt{|N(v_i)| \cdot |N(v_j)|}}$$

$$= \frac{|\{v_3, v_5, v_6, v_7\} \cap \{v_6, v_7, v_9\}|}{\sqrt{|\{v_3, v_5, v_6, v_7\}| \cdot |\{v_6, v_7, v_9\}|}}$$

$$= \frac{|\{v_6, v_7\}|}{\sqrt{12}} = \frac{2}{\sqrt{12}} \approx 0.577$$

⇒ The node is including itself.

$$N(v_4) = \{v_3, v_5, v_6, v_7\}$$

$$N(v_8) = \{v_6, v_7, v_9\}$$

We will check for Jaccard

$$J(v_4, v_8) = \frac{|N(v_4) \cap N(v_8)|}{|N(v_4) \cup N(v_8)|}$$

$$= \frac{|\{v_3, v_4, v_5, v_6, v_7\} \cap \{v_6, v_7, v_8, v_9\}|}{|\{v_3, v_4, v_5, v_6, v_7\} \cup \{v_6, v_7, v_8, v_9\}|}$$

$$= \frac{|\{v_6, v_7\}|}{|\{v_3, v_4, v_5, v_6, v_7, v_8, v_9\}|} = \frac{2}{7} \approx 0.286$$

Now let's calculate for cosine

$$C(v_4, v_8) = \frac{|N(v_4) \cap N(v_8)|}{\sqrt{|N(v_4)| \cdot |N(v_8)|}}$$

$$= \frac{|\{v_3, v_4, v_5, v_6, v_7\} \cap \{v_6, v_7, v_8, v_9\}|}{\sqrt{|\{v_3, v_4, v_5, v_6, v_7\}| \cdot |\{v_6, v_7, v_8, v_9\}|}}$$

$$= \frac{|\{v_6, v_7\}|}{\sqrt{5 \cdot 4}}$$

$$= \frac{2}{\sqrt{20}} \approx 0.447$$

4. SMM Ch 6, problem 9

- **What are the maximum and minimum values for NMI? Provide details.**

The NMI is always constrained by the values of 0 and 1.

When the minimum value of 0 occurs it means that the detected communities and the true labels have no relationship whatsoever. In such cases, knowing the community to which a given node belongs does not help you guess its true label at all.

When the maximum value of 1 occurs, there is a perfect match between the communities and true labels (although it may only use different names). In this case, once you know the community, you know the label for sure so it is a perfect match.

Thus, $NMI \in [0, 1]$. Values close to 1 indicate very strong similarity between communities found and communities corresponding to the true labels and values close to 0 indicate there is essentially no relationship.

- **Explain how NMI works (describe the intuition behind it)**

NMI takes the idea of mutual information and transforms it into a score from 0 to 1.

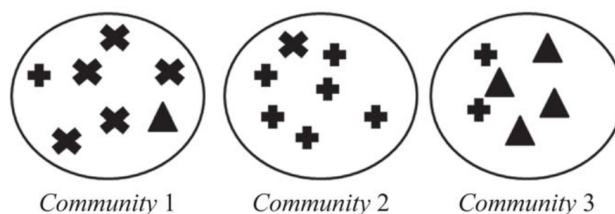
Mutual information tells us how much knowing one of the thing gives you information to use in learning about the other thing. Here, the two things are:

- the community assigned to a given node, and
- the true label of that node.

When clustering is good, knowing the community provides us with a lot of information about the true label. In this case, NMI is high. When clustering is bad, knowing the community does not help us guess the label of the node and NMI is low.

Since we normalize this information we are assured NMI will stay between 0 and 1 allowing us to compare different clusterings fairly even if they have different numbers of communities or different label distributions.

5. SMM Ch 6, problem 10



n_h
 $h=1$ 7 ($7 \div 1 \Rightarrow \times$)
 $h=2$ 7 ($7 \div 2 \Rightarrow +$)
 $h=3$ 6 ($6 \div 3 \Rightarrow \Delta$)

$So, n_{h1}=6, n_{h2}=9, n_{h3}=5$

$n_{h,j}$	$j=1$	$j=2$	$j=3$
$h=1$	5	1	1
$h=2$	1	6	0
$h=3$	0	2	4

$$NMJ = \frac{MI}{\sqrt{H(L)H(H)}} = \frac{\sum_{h \in H} \sum_{j \in L} n_{h,j} \log \frac{n_{h,j}}{n_{h \cdot j}}}{\sqrt{\sum_{j \in L} n_{j \cdot} \log \frac{n_{j \cdot}}{n} \cdot \sum_{h \in H} n_{h \cdot} \log \frac{n_{h \cdot}}{n}}}$$

$$MI = 5 \log \left(\frac{20 \times 5}{7 \times 6} \right) + \log \left(\frac{20 \times 1}{7 \times 9} \right) + \log \left(\frac{20 \times 1}{7 \times 5} \right) + \log \left(\frac{20 \times 1}{7 \times 6} \right) + 6 \log \left(\frac{20 \times 6}{7 \times 9} \right) + 2 \log \left(\frac{20 \times 2}{6 \times 9} \right) + 4 \log \left(\frac{20 \times 4}{6 \times 5} \right)$$

$$= 3.9424$$

$$H(L) = 6 \log \left(\frac{6}{20} \right) + 9 \log \left(\frac{9}{20} \right) + 5 \log \left(\frac{5}{20} \right)$$

$$= -9.2686$$

$$H(H) = 7 \log \left(\frac{7}{20} \right) + 7 \log \left(\frac{7}{20} \right) + 6 \log \left(\frac{6}{20} \right)$$

$$= -9.5203$$

$$\therefore NMJ = \frac{3.9424}{\sqrt{(-9.2686) \cdot (-9.5203)}} = 0.419$$

6. SMM Ch 6, problem 11

High Precision alone is insufficient for community evaluation because it ignores completeness. Precision focuses only on the percentage of correctly identified elements in clusters.

A model with high precision may still fail to identify many true positives, leading to low recall and poor overall performances.

(1) Precision :- Emphasizes the accuracy of the detected clusters.

(2) Recall :- Highlights the completeness of the detected clusters.

Combining precision and recall ensures a balanced and comprehensive evaluation of clustering performances.

Let us consider an example of assessing community detection in a network containing both discovered clusters and real communities.

Ground truth communities :-

➔ Community 1 :- { A, B, C, D, E }

➔ Community 2 :- { F, G, H, I, J }

Detected Clusters:-

- Cluster 1 : {A, B}
- Cluster 2 : {F, G, H, I, J}

Calculating Precision and Recall for Cluster:-

(1) True Positive :- {A, B} \Rightarrow 2

(2) False Positive :- \emptyset None \Rightarrow 0

(3) False Negative :- {C, D, E} \Rightarrow 3

$$\begin{aligned}\Rightarrow \text{Precision} &= \text{TP}/(\text{TP} + \text{FP}) \\ &= 2/(2 + 0) \\ &= 1\end{aligned}$$

$$\begin{aligned}\Rightarrow \text{Recall} &= \text{TP}/(\text{TP} + \text{FN}) \\ &= 2/(2 + 3) \\ &= 0.4\end{aligned}$$

This demonstrates that while precision is perfect for Cluster 1, the recall is low.

In conclusion, a high precision alone can be misleading as a good measure of performance, especially when the data are imbalanced; only by taking into account precision and recall can we meaningfully assess the performance of a community detection method.

7. SMM Ch 6, problem 12

Purity is a metric that measures the homogeneity of clusters with respect to the ground truth labels.

Although it is easy and intuitive to compute, purity is either inappropriate or misleading in some cases as a method of evaluating community detection performances.

Situation where Purity is Misleading :-

(1) Lack of Balance between Clusters :-

The value of purity can be driven upwards when the algorithm obtains multiple small "pure" clusters. These will not capture the global community structure well.

Let take an example to solve it.

- Ground Truth communities:-
Community 1 : { A, B, C, D, E }
Community 2 :- { F, G, H, I, J }
- Detected Cluster:-
Cluster1 – { A }
Cluster2 – { B }
Cluster3– { C }
Cluster4 – { F, G, H, I, J }

If each detected cluster includes nodes exclusively from the same ground truth community, purity will be 1.

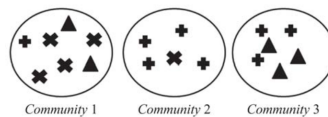
(2) Ignoring Completeness :-

Purity focuses only on how homogenous a cluster is. It overlooks whether all relevant nodes from a community are assigned to the same cluster. As a result cluster is ignored.

(3) Imbalanced Data :-

When ground truth communities are imbalanced purity becomes deceptive. Assigning nodes from smaller communities to larger ones artificially inflates the score, favoring larger community unfairly.

8. SMM Ch 6, problem 13



• So let us solve the Precision and recall

(1) True Positive :- similar members are assigned to the same community.

$$TP = \left(\binom{4}{2} + \binom{2}{2} \right) + \binom{4}{2} + \left(\binom{3}{2} + \binom{3}{2} \right)$$

Community 1 Community 2 Community 3

$$= (6+1) + 6 + (3+3)$$

$$= 7 + 6 + 6$$

$$= 19$$

(2) False Positive :- Disimilar members are assigned to the different community.

$$FP = (4 \times 1 + 4 \times 2 + 2 \times 1) + (4 \times 1) + (3 \times 3)$$

Community 1 Community 2 Community 3

$$= (4+8+2) + 4 + 9$$

$$= 27$$

(3) False negative :- Similar members are assigned to different communities.

$$FN = \underbrace{(4 \times 1)}_{(x)} + \underbrace{(1 \times 4 + 6 \times 3 + 8 \times 1)}_{(+1)} + \underbrace{(2 \times 3)}_{(\Delta)}$$

$$= 4 + 19 + 6$$

$$= 29$$

(4) True negative :- dissimilar members are assigned to the different communities.

$$TN = \underbrace{(4 \times 4 + 1 \times 1 + 2 \times 4 + 2 \times 1)}_{\text{Community 1+2}} + \underbrace{(4 \times 3 + 1 \times 3 + 1 \times 3)}_{\text{Community 2+3}}$$

$$+ \underbrace{(4 \times 3 + 4 \times 3 + 1 \times 3 + 2 \times 3)}_{\text{Community 1+3}}$$

$$= (16 + 1 + 8 + 2) + (12 + 3 + 3) + (12 + 12 + 3 + 6)$$

$$= 27 + 33 + 18$$

$$= 78$$

So Now Calculate Precision and recall

$$\therefore P = \frac{TP}{TP + FP} = \frac{19}{19 + 27} = 0.4130$$

$$\therefore R = \frac{TP}{TP + FN} = \frac{19}{19 + 29} = 0.3958$$

• F-measure :- $2 \times \frac{P \times R}{P + R}$

$$= 2 \times \frac{(0.4130 \times 0.3958)}{(0.4130 + 0.3958)}$$

$$= 2 \times \frac{0.1634}{0.8088}$$

$$= 0.4040$$

• Purity :- $\frac{1}{N} \sum_{i=1}^K \max_j |C_i \cap L_j|$

Here, $N = 18$
 $K = 3$

Majority labels in Com-1 is X, which has 4 nodes
Majority label in Com-2 is 7, which has 4 nodes
Majority labels in Com-3 is Δ, which has 3 nodes

$$\therefore \text{Purity} = \frac{1}{18} \sum_{i=1}^3 \max_j |C_i \cap L_j|$$

$$= \frac{1}{18} (4 + 4 + 3)$$

$$= \frac{11}{18}$$

$$= 0.6111$$

$$n_h$$

$h=1$	7	$(d=1 \Rightarrow X)$
$h=2$	5	$(d=2 \Rightarrow +)$
$h=3$	6	$(d=3 \Rightarrow \Delta)$

$$S_e, n_{d=1} = 5, n_{d=2} = 8, n_{d=3} = 5$$

n_{h-1}	d_1	d_2	d_3
$h=1$	4	1	2
$h=2$	1	4	0
$h=3$	0	3	3

$$\therefore NMI = \frac{MI}{\sqrt{H(L) \cdot H(H)}} = \frac{\sum_{h \in H} \sum_{l \in L} n_{h,l} \log \left(\frac{n_{h,l}}{n_l \cdot n_h} \right)}{\sqrt{\left(\sum_{l \in L} n_l \log \frac{n_l}{n} \right) \cdot \left(\sum_{h \in H} n_h \log \frac{n_h}{n} \right)}}$$

$$\rightarrow MI = 4 \log \left(\frac{18 \times 4}{7 \times 5} \right) + \log \left(\frac{18 \times 1}{7 \times 8} \right) + 2 \log \left(\frac{18 \times 2}{7 \times 5} \right) + \log \left(\frac{18 \times 1}{5 \times 5} \right)$$

$$+ 4 \log \left(\frac{18 \times 4}{5 \times 8} \right) + 3 \log \left(\frac{18 \times 3}{6 \times 8} \right) + 3 \log \left(\frac{18 \times 3}{6 \times 5} \right)$$

$$= 2.5823$$

$$H(L) = 5 \log \left(\frac{5}{18} \right) + 8 \log \left(\frac{8}{18} \right) + 5 \log \left(\frac{5}{18} \right)$$

$$= -8.3804$$

$$H(H) = 7 \log \left(\frac{7}{18} \right) + 5 \log \left(\frac{5}{18} \right) + 6 \log \left(\frac{6}{18} \right)$$

$$= -8.5154$$

$$\therefore NMI = \frac{2.5823}{\sqrt{(-8.3804) \cdot (-8.5154)}} = 0.305$$

9. SMM Ch 7, problem 2

We need to establish an experiment that addresses herd behavior, where individuals can make decisions consecutively, each being able to see what the previous individuals have decided, but still have some private information of their own. There will need to be two clear choices for example let's take Community Area X vs Community Area Y. Each individual is assigned a private signal and individuals are instructed not to share this signal with anyone or talk about the signal itself. Individuals are instructed only to view what others have decided. Under this procedure we will be able to see if people continue to use their own information or they start to use the majority behavior they see in all the previous people.

For a concrete design, consider an example concerning a study of the best locations for new housing: Community Area X or Community Area Y. The city knows in secret that one area has slightly better conditions than the other, but the participants do not. The participants come in one at a time. Each participant is given a private brochure that weakly recommends one of the two community areas X or Y that has been simulated on data related to safety, schools, or services. After reading the brochure, the participant must publicly state I would recommend Area X or I would recommend Area Y, and their recommendation is written on a board that all of the subsequent participants will see. The next participant enters the room after all of the previous participants have made their public recommendation. The subsequent participant sees all of the previous recommendations on the board, reads the brochure they have received, and makes their public recommendation thereafter. Over time, if many of the earlier participants pick the same community area to recommend, then the late participants may ignore

their own brochures and simply follow the apparent majority recommendation. This is an example of ensemble behavior in a realistic community area setting.

10. SMM Ch 7, problem 3

The dissemination of a novel idea or product is modelled by

$$dA(t)/dt = i(t)[P-A(t)]$$

where $A(t)$ refers to the number of the adopters at time t , P signifies the total number of individuals who could adopt, and $i(t)$ denotes the strength of influence at time t . The phrase $[P-A(t)]$ relays the number of individuals left to adopt the idea or product.

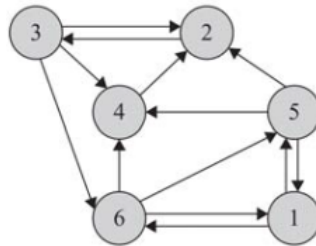
In the **external-influence model**, a person adopts primarily due to external factors, particularly the mass media and advertising, so adoption is instantaneous and grows at an even rate over time. In the **internal-influence model**, a person adopts due to their social network, and we will see the adoption of new ideas follow an S-curve with initial slow adoption, rapid adoption, and finally, a decreasing gain as adoption nears completion. In the **mixed-influence model**, you see both influences, so we will see early adoption and rapid gain in adoption due to external influence but eventually rapid gain, especially when the majority adopts because of internal influence.

11. SMM Ch 7, problem 4

The innovation is a community safety app introduced in an urban neighborhood to address local safety needs. The population is all residents who could potentially use the app. At first, only a few early adopters download it and report issues like broken streetlights or suspicious activity, influencing their neighbors by using and talking about it. At the same time the city promotes the app through ads, emails, signs, adding external influence. Together these internal and external influences drive the diffusion process: adoption is slow in the beginning then speeds up as more people hear about and use the app and finally levels off once most interested residents have adopted it.

An intervention is any deliberate action that changes this diffusion process. To speed up diffusion the city can increase external influence by running a larger awareness campaign or sending direct invitations and increase internal influence by rewarding people for inviting neighbors. To slow down a harmful innovation officials can warn people reduce promotion and discourage sharing so the idea spreads much more slowly.

12. SMM Ch 7, problem 7



The process follows the Independent Cascade Model. Each node gets only one chance, in the time step after it becomes active, to activate each of its outgoing neighbors. In this problem a node i activates a neighbor j only if there is a directed edge $i \rightarrow j$ and $i - j \equiv 1 \pmod{3}$.

Node 5 is given as active at time 0. Let $A(i, t) = 1$ if node i is active at time t .

Time 0:-

Only node 5 is active.

$$A(5, 0) = 1.$$

Active set = {5}.

Time 1:-

Node 5 now tries to activate its out-neighbors 4, 1 and 2.

$$5 - 4 = 1 \equiv 1 \pmod{3} \Rightarrow \text{node 4 becomes active, } A(4, 1) = 1.$$

$$5 - 1 = 4 \equiv 1 \pmod{3} \Rightarrow \text{node 1 becomes active, } A(1, 1) = 1.$$

$$5 - 2 = 3 \equiv 0 \pmod{3} \Rightarrow \text{node 2 is not activated.}$$

New active nodes at time 1: {4, 1}.

Cumulative active set after time 1: {5, 4, 1}.

Time 2:-

Only the nodes that became active at time 1 (nodes 4 and 1) get a chance now.

Node 4 has out-neighbor 2:

$$4 - 2 = 2 \equiv 2 \pmod{3} \Rightarrow \text{node 2 is not activated.}$$

Node 1 has out-neighbors 5 and 6:

$$1 - 5 = -4 \equiv 2 \pmod{3} \Rightarrow \text{node 5 is not newly activated.}$$

$$1 - 6 = -5 \equiv 1 \pmod{3} \Rightarrow \text{node 6 becomes active, so } A(6, 2) = 1.$$

New active node at time 2: {6}.

Cumulative active set after time 2: {5, 4, 1, 6}.

Time 3:-

Only node 6 is newly active, so it now gets its one chance to activate its neighbors 1, 4 and 5.

$6 - 1 = 5 \equiv 2 \pmod{3} \Rightarrow$ node 1 is not newly activated.

$6 - 4 = 2 \equiv 2 \pmod{3} \Rightarrow$ node 4 is not newly activated.

$6 - 5 = 1 \equiv 1 \pmod{3}$ but node 5 is already active, so there will be no change.

No new nodes are activated after time 3, so the process has converged.

The Independent Cascade Model converges with the active nodes {1, 4, 5, 6}.