

Name: Purushothaman, S

Roll No: G24AI1042

Subject: Social Network Analysis.

Subject Code: CSL7390

Exam Type: Trimester II - PGD. Data Engineering - Major Exams.

Date: 19th April 2025 Time: 2.00 PM. - 5.00 PM.

Email ID: g24ai1042@iitj.ac.in

College/University: IIT - Jodhpur.

Course & Batch: PGD. Data Engineering / July 2024.

(1) Convert the incidence Matrix to adjacency Matrix:

(2)

Given Matrix:

$\begin{matrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{matrix}$	<p><u>Column 1:</u> Connects Node 1 and Node 2 \rightarrow edge between 1 and 2</p> <p><u>Column 2:</u> Connects Node 2 and Node 3 Edge between 2 & 3.</p>
--	---

Column 3: Connects Node 2 & Node 4 \Rightarrow Edge between 2 & 4.

So Adjacency Matrix:

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

(1)(b) Which Network Model Assumes that edges are formed between pairs of Nodes with a uniform probability, independent of each other.

Ans: (B) Erdős-Rényi (Random Network)

Reason: This Model Connects Node Pairs with uniform and independent probability.

(1)
(c) In game theory, a situation where no player can improve their outcome by unilaterally changing their strategy, given the strategies of other players is known as:

Ans: c Nash Equilibrium.

Reason: This is when no player benefits by changing strategy unilaterally.

(1) The tendency for individuals in social networks to associate and bond with similar others is defined as.

Ans: Assortative Mixing

Reason: This describes preferences for connecting with similar others.

(2) Why might betweenness centrality be a more relevant measure than degree centrality for identifying critical nodes in a network transmitting information that must follow specific paths.

Ans: Because it quantifies how often a node lies on the shortest paths between other nodes.

Reason: Betweenness centrality highlights nodes critical to information flow.

(1) A key finding about scale-free networks (like those generated by the Barabasi-Albert Model) is their robustness to random node failure but vulnerabilities to targeted attacks on hubs. What underlying property best explains this?

Ans: The presence of many nodes with very high valencies (hubs)

Reason: Hubs keep the network connected; their loss disrupts connectivity.

(2) In community detection, optimizing for high modularity aims to find partitions, where:

Ans: The number of intra-community edges is significantly higher than expected in a random network with the same degree sequences.

Reason: High Modularity reflects strong Community Structure.

(1) Consider the two Nodes X & Y. The Neighbors of X are A, B, C, D. & The Neighbors of Y are C, D, E. What is the Jaccard Co-efficient for link prediction between X & Y?

Ans: $2/5$.

Reason: Jaccard Co-efficient = $|C \cap D| / |C \cup D|$

$$= \{c, d\} / \{a, b, c, d, e\} = 2/5$$

(2) In the Context of Information Cascade Models, how does the Activation Mechanism differ fundamentally between the Independent Cascade Model (ICM) and the Linear Threshold Model (LTM)?

Ans: ICM uses edge probabilities; LTM uses weighted sum (v_i) Threshold.

Reason: ICM is probabilistic per edge;
LTM checks if influence exceeds
threshold.

(1) A Standard graph Convolutional Network (GCN) aggregates information from a Node's immediate Neighbors. Why might this standard message-passing approach be suboptimal for tasks like Node Classification in Networks with high heterophily (where connected Nodes tend to be dissimilar)?

Answer: Because aggregating features from dissimilar Neighbors can blur the Node's own identity representative features, making classification harder.

Reason: Heterophilic graph lead to poor GCN Performance due to Noisy Aggregation.

Qn ②

SolutionProblem Statement:

Vaccinate 5% of the population to minimize infections in a social contact network using Network Analysis concept.

SIR Model:

$S \rightarrow$ Susceptible, $I \rightarrow$ Infected

$R \rightarrow$ Recovered (includes Vaccinated Individuals)

Goal: Select Critical Nodes to vaccinate (Premarktively Move to R)

Strategy: Use two Network Measures.

1) Degree Centrality: Nodes with the highest Number of direct Connections

why?: They can infect Many Neighbors quickly.

② Betweenness Centrality

Measures how often a Node lies on the shortest paths between other Nodes

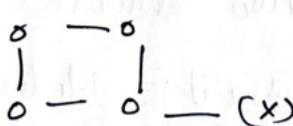
Why?: Nodes acting as bridges between clusters can spread infection across the Network.

Combined Approach:

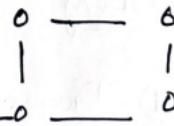
Rank all nodes by a Combined Score
(e.g. weighted sum of degree and betweenness) & then select the top 5%

Diagrams:

Cluster A



Cluster B



(\Rightarrow) high betweenness.

Vaccinating ((x)) cuts off inter-cluster transmissions.

(3)

Suggested Collaborators: link prediction + Embedding.

Task: Use Node2Vec and link prediction to recommend academic collaborators

Solution

(1) Node Embedding. (Node2Vec).

This converts Nodes into Vector representations based on Network structures.

Similar Nodes (eg. co-authors) are closer in vector spaces.

(2) Link predictions

To predict the link we can use Similarity Metrics (like Cosine similarity) between embeddings.

Predict New links (collaborations) between similar researchers.

Cosine Similarity formula:

$$\text{Cosine Similarity} = (\mathbf{A} \cdot \mathbf{B}) / (|\mathbf{A}| * |\mathbf{B}|).$$

(5) Homophily: Researchers in similar fields tend to collaborate more (birds of a feather)

To encourage Cross-Disciplinary Collaboration:

⇒ Use Metadata (e.g.) research field to diversify recommendations.

⇒ Reward link suggestions between structurally similar but topically diverse nodes.

Table Example:

Researcher A	Field	Similarity Score	Cross-Discipline
Dr. X	AI	0.92	No
Dr. Y	Biology	0.85	Yes (Dissimile)

Q: 4: Community Detection: Girvan - Newman
Vs. Louvain. (6)

(a) Girvan - Newman Algorithm

It removes edges with highest betweenness centrality. As edges are removed, graph breaks into communities.

Ex:

Before A — B — C
 |
 D

+ B-D has edge Betweenness

So after removing B-D:

A-B C
 D \Rightarrow Two Communities.

(b) Edge between centrality:

Number of shortest paths between nodes that passes through the edge.

Iteratively remove highest-scoring edges.

(c) Limitations

- i) Expensive computation ($O(n^3)$),
- ii) Not suitable, Scalable for Large graphs.

(d) Louvain Method:

It optimizes modularity (Quality of community partitions). There are two phases:

- i) Local Modularity Optimization,
- ii) Node Aggregation and Iteration

for much faster and scalable.

(n: 5) Page Rank Algorithm:

The intuition behind the Page Rank algorithm:

It is important to link the important pages by other important pages.

It simulates a random web Surfer.

Equation :

$$PR(i) = (1-d)/N + d * \sum [PR(j) / \text{out-degree}(j)]$$

(b) Damping factor:

Damping factor (d) (typically 0.85)
represents the probability of following
a link.

$(1-d)$ represents random Jumps to
any Node (It prevent getting any stuck)
it helps in Maintaining Convergence and
reflecting realistic Surfing Behaviour

(c) Dangling Node Problem:

The Dangling Nodes have No
outgoing edges (No links to follow)

Problem their page rank can leak out of the system, affecting convergence.

Solution: Redistribute their Page Rank uniformly across all Nodes in each iteration

This keeps the transition matrix stochastic (rows sum to 1), ensuring stable and correct result.

Page Rank Formula:

For a Node i , the Page Rank.

$PR(i)$ is calculated as.

$$PR(i) = \frac{1-d}{N} + d \sum_{j \in L(i)} \frac{PR(j)}{\text{outdegree}(j)}$$

Where

d - damping factor (usually 0.85)

$N \rightarrow$ Total Number of Nodes.

$L(i) =$ set of Nodes linking to Node (i) .

(8)

Outdegree (j) = Number of links going out from Node j .

Ex:

$A \rightarrow B$
 $B \rightarrow C$
 $C \rightarrow A$
Each Node linked to one other Node.

Let Assuming the Damping factor
 $d = 0.85$

Initial pageRank of each Node: $PR = 1/3$

Iterating 1 calculated:

(i) $PR(A)$ A has one incoming link from

$$PR(C) = 1/3 \quad \text{Out degree}(C) = 1.$$

$$PR(A) = \frac{1 - 0.85}{3} + 0.85 \times \frac{1}{3} = 0.05 + 0.283 \\ = 0.333$$

$$PR(B) = 0.05 + 0.85 \times \frac{1}{3} = 0.333$$

$$PR(C) = 0.05 + 0.85 \times \frac{1}{3} = 0.333.$$

In this perfectly circular graph, the PageRank remains equal at each step.

Dangling Node case:

Let's say Node C has No Outgoing links (a dangling Node):

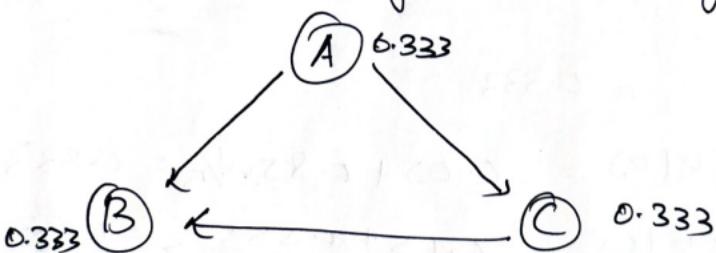
A $\xrightarrow{?}$ B
B $\xrightarrow{?}$ C
C $\xrightarrow{?}$ \emptyset

} To handle dangling Node C

Its Rank Redistributed uniformly to A, B, & C.

so $PR(C) = 0.3$ this 0.3 is divided equally. Each Node gets $0.3/3 = 0.1$.

This ensures the page Rank remains Normalized & algorithm converges.



Q/N(6)

Solution:

Game between two players are here shown in table.

<u>Payoff Matrix:</u>		Strategy A	Strategy B
U (Player 1)	(3, 2)	(0, 1)	
L (Player 2)	(2, 0)	(2, 3)	

here first value = player 1's payoff
 second value = player 2's payoff

(a) Pure Strategy Nash Equilibrium

A Nash equilibrium is where no player wants to change strategy unilaterally.

check each cell:

$\rightarrow (U, A)$: Player 2 prefers B since $(\overline{U} > 2)$.

(U, B) player 1 prefers L (since $2 > 0$)

(L, A) player 2 prefers B (since $3 > 0$)

(L, B) Both players are best-responding
to each other \rightarrow Nash Equilibrium.

Answer: only (L, B) is pure strategy
Nash equilibrium.

(b) Expected payoffs for players above.

Let player 1 play:

U with probability P.

L with probability $(1-P)$.

Payoff of player 2:

If they choose A:

$$= 2P + 0(1-P) = 2P.$$

If they choose B

$$1P + 3(1-P) = P + 3 - 3P = 3 - 2P.$$

(C) For ($P=0.7$)

If $P=0.7$ then. $A : 20 \cdot 0.7 = 1.4$ &

$$B = 3 - 20 \cdot 0.7 = 1.6$$

Here player A should be choose strategy B, as it gives higher expected pay off.

Here the Table:

		Strategy A.	
		(3, 2)	(0, 1)
		4	(2, 0)
Player 1	L	(2, 0)	
	U	(3, 2)	(2, 1)

Qn. ④

GNN Layer update for Node B.

Problem Statement: The given

directed graph where Nodes A, C and D point to Node B. Here we need to calculate the updated feature

Vector $h_B^{(1)}$ using simple GINN Layer

Solns

GINN Layer Operations

- ⇒ Aggregate Neighbor Features (Average)
- ⇒ Transform using a weight Matrix
- ⇒ Activate using ReLU

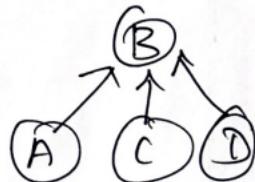
Given Data:

Initial Feature Vectors.

$$\cdot h_A^{(0)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\cdot h_C^{(0)} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

$$\cdot h_D^{(0)} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$



Neighbors of A, C, D.

Weight Matrix W:

$$W = \begin{bmatrix} 0.5 & 0 \\ 0.1 & 0.2 \end{bmatrix}$$

Aggregate - Neighbor Feature (Average):

$$h_N^{(0)}(B) = \frac{1}{3} (h_A^{(0)} + h_C^{(0)} + h_D^{(0)}) = \frac{1}{3} \left[\frac{1+0+2}{1+3+0} \right] = \frac{1}{3} \left[\frac{3}{6} \right]$$

Linear Transformation

$$W \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0.5 \times 1 + 0.2 \\ 0.1 \times 1 + 0.2 \times 2 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.1 + 0.4 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

Apply ReLU Activation

$\text{ReLU}(x) = \max(0, x)$ applied element wise

$$h_B^{(1)} = \text{ReLU} \left(\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \right) = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

So the final Answer:

$$h_B^{(1)} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

This is the updated feature vector for Node B after one GNN Layer.