

## II) Complexities in data collection:

1. Dynamic nature of data demands collection at uniform intervals of time.
2. Accuracy & Reliability of data collection.
3. Privacy concerns of data.
4. Takes a lot of effort - 50% of effort wasted in data collection.
5. Since the collection process is difficult, collected data is reused often.

## Alternatives

1. Online e-data can be collected.  
The web is vast and diverse. The information from web can be collected and used for the purpose of SNA.

For eg: Information Dynamics lab at HP investigated the presence of ties based off of corporate archive of mail. If enough mails were sent, a tie was assumed.

However, data privacy issues are a major concern.  
~~According to H.~~ since private, sensitive data cannot be easily shared.

Instead, data from ~~not~~ open-source archives can be used. Data Visualisation & diffusion analysis can also be performed.

2.

Data not meant for SNA can also be used:

- archival data (minute of meeting/court records, etc)
- web-data through web-mining.

3.

### Blogs:

They provide a platform for people to interact through comments and links. They're mostly used for marketing trend analysis to predict consumer choices.

- links are used since they're representative of relationships & are authoritative.
- Co-occurrences are also used as an indicator of ties.

Referral Web first used web-mining for referral chaining - to find an expert of a particular niche closest to the user.

Jaccard Coefficient: 
$$\frac{\text{no. of co-occurrences}}{\text{no. of indiv. occurrences}}$$

To extract ties, we need an efficient commercial search engine.

First, the names are extracted using NLP & further used to extract newer names [snowballing]

#### 4. Social media platforms

SSN

5

- Better way of social networking
- Extra features like photo sharing, messaging, etc
- To analyze changes in membership communities  
- e.g. effect of structure
- Guarded data with security measure
- Alternatively, use FOAF networks which store user profiles on their personal blogs, connected by links
- Less support for maintaining & extraction of info.

#### Other disadvantages

- Technological limits to accuracy.
- Can be costly.

18) a) Eccentricity is the longest shortest path  
seen by node 3  
Path length

- (3,1) - 2 (3,8) - 4
- (3,2) - 1 (3,9) - 3
- (3,3) - 0 (3,10) - 4
- (3,4) - 2 (3,11) - 4.
- (3,5) - 3
- (3,6) - 4
- (3,7) - 5

Eccentricity - 5

b) A walk is a ~~pat~~ sequence of nodes & lines in the graph w/ possible repetition.

Walk - n<sub>0</sub>, n<sub>4</sub>, n<sub>9</sub>, n<sub>11</sub>, n<sub>10</sub>, n<sub>9</sub>, n<sub>11</sub>

c) A trial is one walk where the lines are unique and not repeated

Trial - n<sub>1</sub>, n<sub>2</sub>, n<sub>3</sub>, n<sub>0</sub>.

d) A tour is a closed walk in which all lines occur atleast once.

n<sub>0</sub>-n<sub>3</sub>-n<sub>2</sub>-n<sub>1</sub>-n<sub>0</sub> - n<sub>4</sub>-n<sub>5</sub>-n<sub>8</sub>-n<sub>7</sub>-n<sub>6</sub>-n<sub>5</sub>-  
n<sub>4</sub>-n<sub>9</sub>-n<sub>11</sub>-n<sub>10</sub>-n<sub>9</sub>-n<sub>4</sub>-n<sub>0</sub>.

SSN

e) Nodal degree - no. of lines incident

n <sub>0</sub> - 3	n <sub>5</sub> - 3
n <sub>1</sub> - 2	n <sub>6</sub> - 2
n <sub>2</sub> - 2	n <sub>7</sub> - 2.
n <sub>3</sub> - 2.	n <sub>8</sub> - 2.
n <sub>4</sub> - 3.	n <sub>9</sub> - 3.
	n <sub>10</sub> - 2.

f) mean nodal degree.

$$\bar{d} = \frac{\sum_{i=0}^9 d(i)}{9} = \frac{3+2+2+2+3+2+2+2+3}{9} = \frac{28}{12} = 2.33$$

g) Variance.

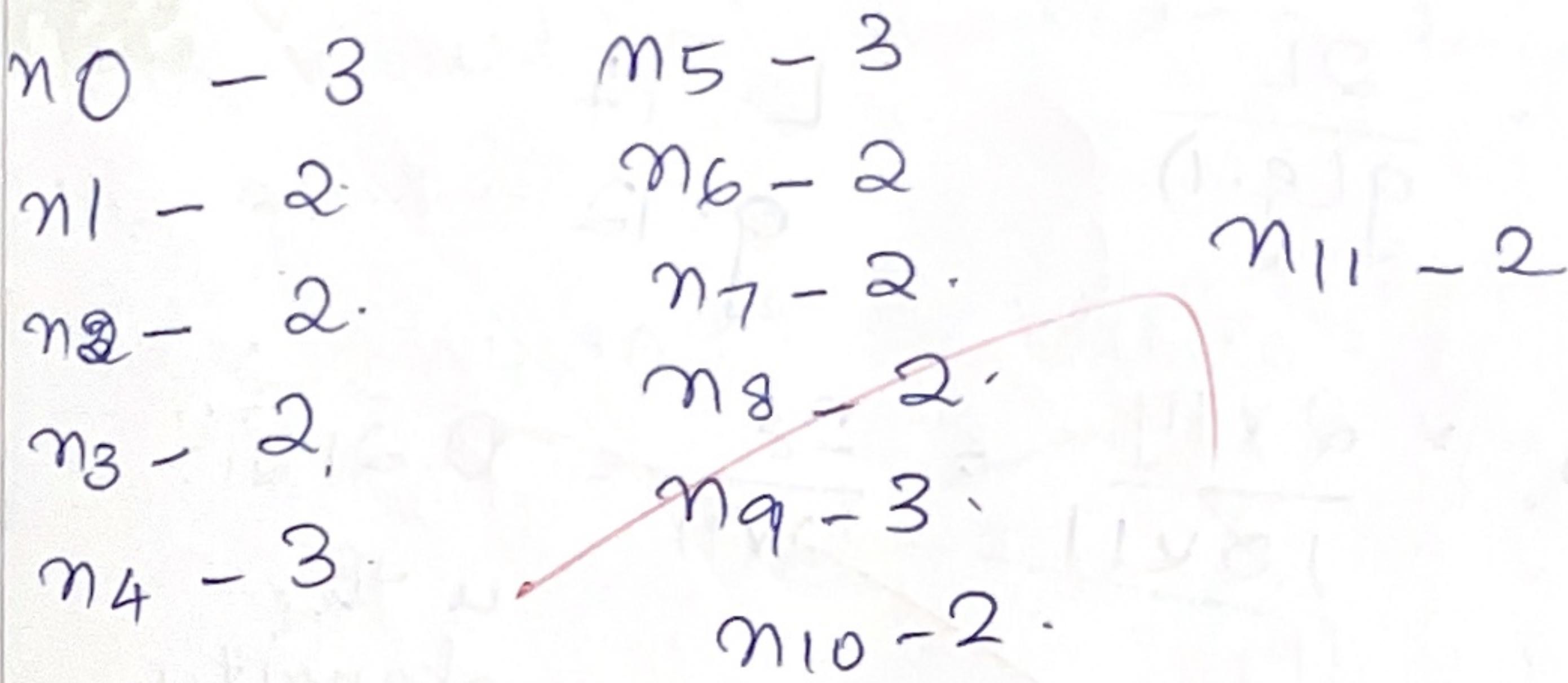
$$S^2 = \frac{\sum (d_i - \bar{d})^2}{9}$$

$$= \frac{(3-2.33)^2 + (2-2.33)^2 + \dots}{12}$$

$$= \frac{4(3-2.33)^2 + 8(2-2.33)^2}{12} = 1.70$$

$$= \frac{2.6668}{12} = 0.222 \text{ is the variance}$$

e) Nodal degree - no. of lines incident to it.



f) mean nodal degree.

$$\bar{d} = \frac{\sum_{i=0}^9 d(i)}{9} = \frac{3+2+2+2+3+3+2+2+2+3+2}{12} = \frac{28}{12} = 2.33 \text{ degrees/node}$$

g) Variance.

$$s^2 = \frac{\sum (d_i - \bar{d})^2}{n}$$

or:

$$= \frac{(3-2.33)^2 + (2-2.33)^2 + \dots}{12}$$

$$= \frac{4(3-2.33)^2 + 8(2-2.33)^2}{12} = \frac{1.7956 + 0.8712}{12}$$

$$= \frac{2.6668}{12} = 0.222 \text{ is the variance}$$

5) Density of graph.

$$= \frac{d}{g(g-1)} = \frac{2L}{g(g-1)}$$

$$L = 14$$

$$g = 12$$

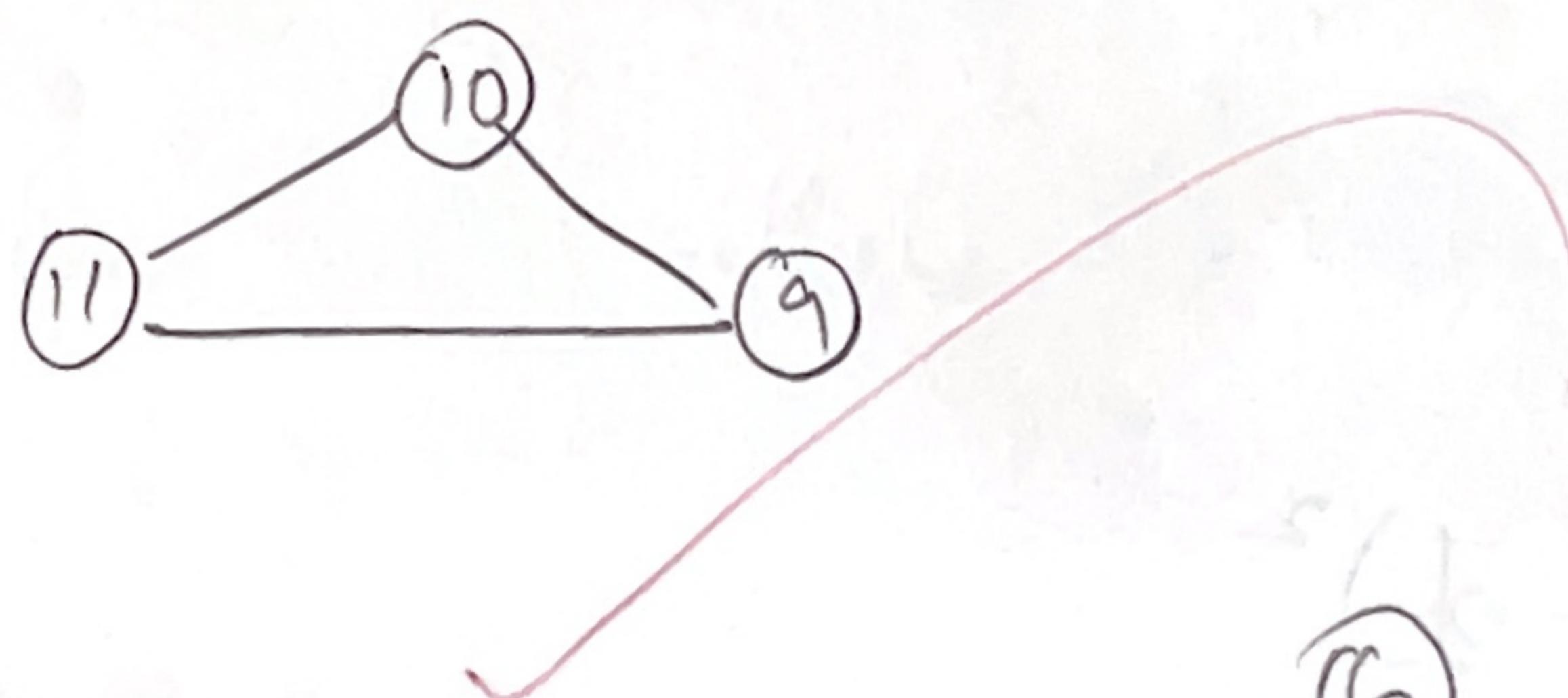
$$\text{Density } D = \frac{2 \times 14}{12 \times 11} = \frac{28}{12 \times 11} = 0.21\bar{2}\dots$$

in the density

(i) cut-point/node connectivity ~~are~~<sup>is</sup> the set of nodes for which the graph will become disconnected if removed.

cut-point = {n4}

1 → point connectivity

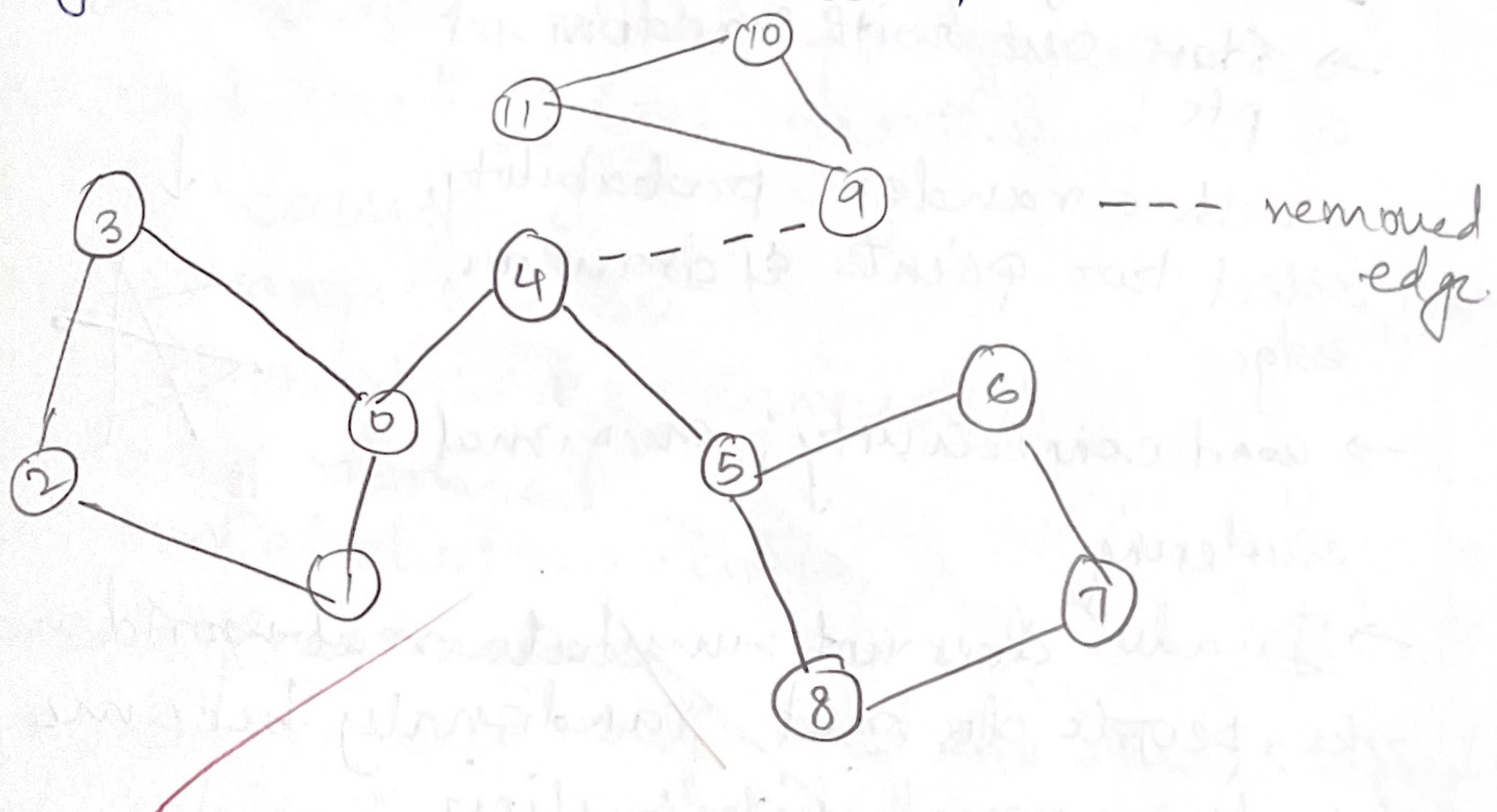


disconnected

(g) Edge connectivity - min no. of edges that should be removed to disconnect the graph.

Edge  $\{n_4-n_9\}$

Value - 1



→) Different models for graph networks

1) Random graph model

→ Start out with random set of pts

→ with a random probability, select two points & draw an edge

→ good connectivity; minimal clustering

→ Disadv: does not simulate real-world networks; people do not randomly become friends; normal distribution

→ can use to simulate complex structures

2) Alpha model

→ proposed in seminal paper.

→ The probability of having a connection between 2 nodes proportional to the no. of shared mutual neighbors.

→  $\alpha$  parameter used to control the influence of mutual friends on connection

→ highly connected; Great clustering

→ However, normally distributed

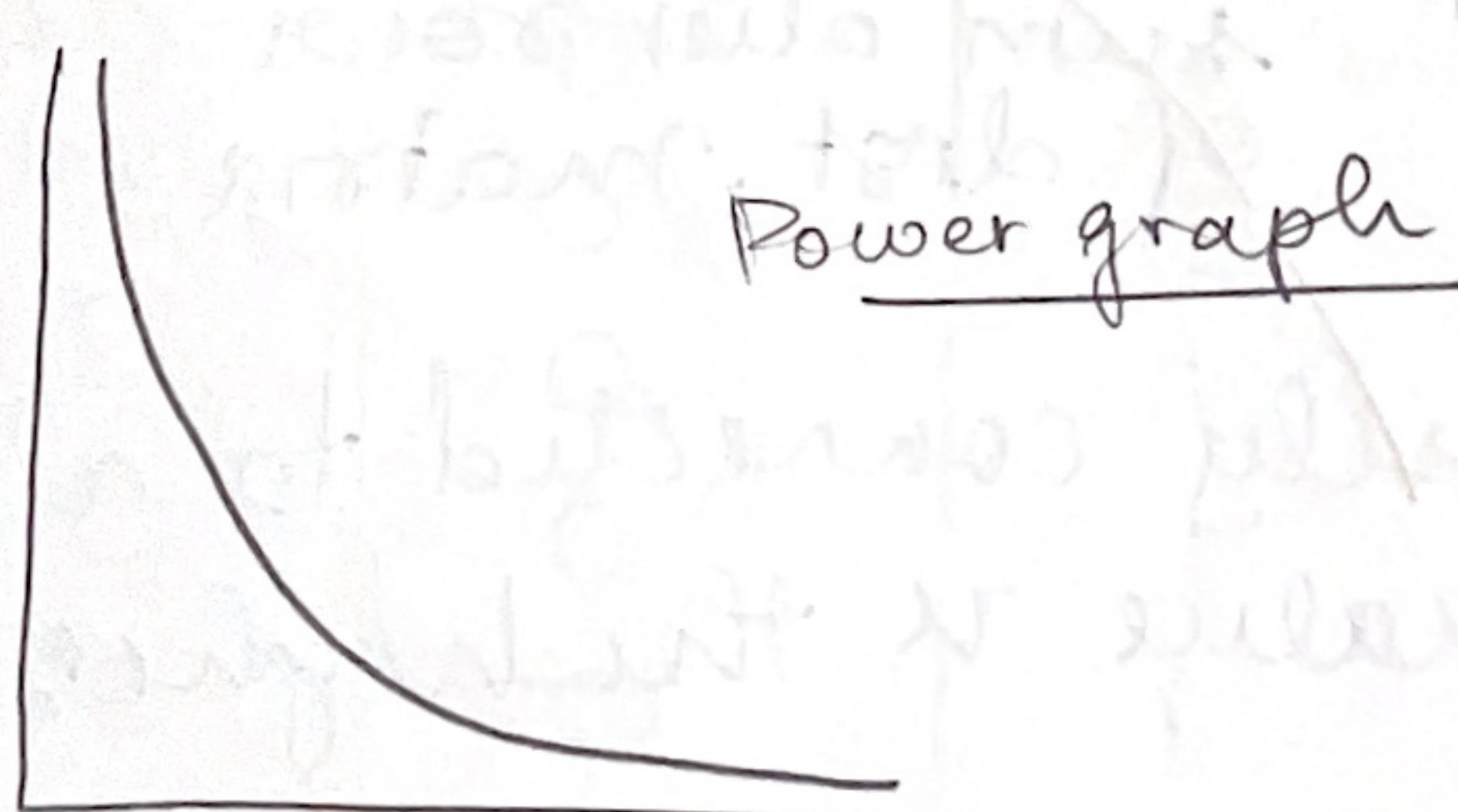


### 3) Beta model

- Starts with a toroidal lattice,
- Nodes connected to neighbours & neighbors of neighbors.
- with random probability  $\beta$ , keeping one end fixed, select another end & reassign it.
- Rearrange like so.
- Well clustered & connected networks formed
- Normal distribution

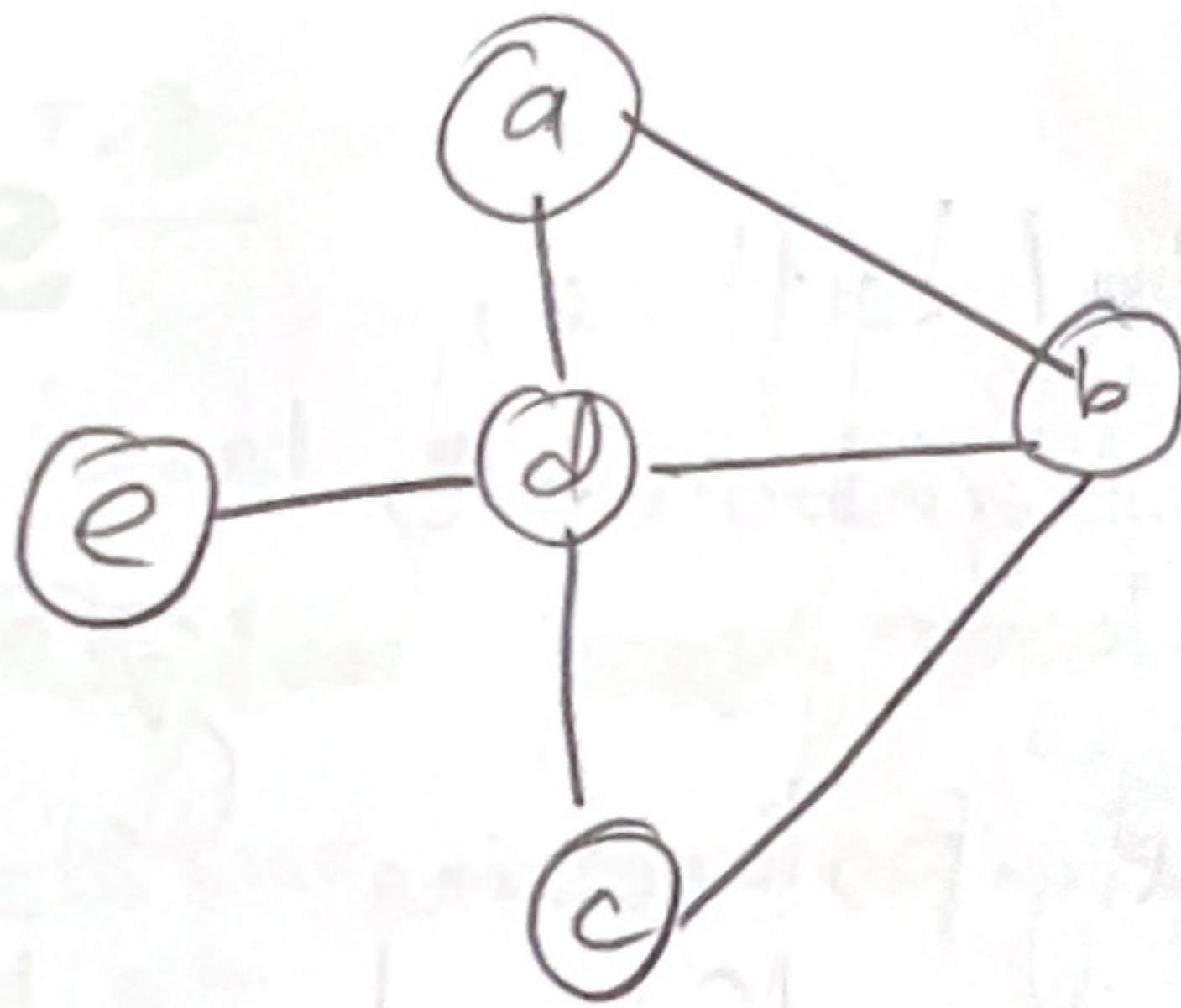
### 4) Scale-free models

- prev. models only cared abt small world prop not scale-free characteristic.
- nodes with large no. of degrees lesser than ones with less no. of degrees
- Rich get richer idea.
- Distribution follows power law  $d^{-k}$   
for  $k > 0$



$$d^{-k} \quad [k > 0]$$

3)



Whom do S.S.N.

SSN

Closeness centrality.

Distance matrix.

	A	B	C	D	E	Sum	closeness centrality
A	0	1	2	1	2	6	0.1666
B	1	0	1	1	2	5	0.2
C	2	1	0	1	2	6	0.1666
D	1	1	1	0	1	4	0.25
E	2	2	2	1	0	7	0.143

where  $D_{ij}$  = shortest path between node i and node j

closeness centrality =  $\frac{1}{\text{sum over rows of dist. matrix}}$

since D is most closely connected to all other nodes, its value is the highest, with 0.25

Betweenness centrality =  $\frac{\sigma_{uv}(w)}{\sigma_{uv}}$

$\sigma_{uv}$  = shortest path between  $u \in V$

$\sigma_{uv}(w)$  = shortest path b/w  $u \in V$ , including  $w$ .

$w=A$

Path	$\sigma_{uv}(w)$	$\sigma_{uv}$	BC
BC	0	1	0
BD	1	2	0.5
BE	1	2	0.5
CD	0	1	0
CE	0	1	0
DE	0	1	0

$w=B$

Path	$\sigma_{uv}(w)$	$\sigma_{uv}$	BC
AC	1	2	0.5
AD	0	1	0
AE	0	1	0
CD	0	1	0
CE	0	1	0
DE	0	1	0

$\omega = C$ 

Path	$\sigma_{ij}$	$\sigma_{ij}(\omega)$	BC
AB	1	0	0
AD	1	0	0
AE	1	0	0
BD	2	1	0.5
BE	2	1	0.5
DE	0	1	0
			<u>1</u>

 $\omega = D$ 

Path	$\sigma_{ij}$	$\sigma_{ij}(\omega)$	BC
AB	1	0	0
AC	2	1	0.5
AE	1	1	1
BC	1	0	0
BE	2	2	1
CE	1	1	1
			<u>3.5</u>

highest for  
node D = 3.5

 $\omega = E$ 

Path	$\sigma_{ij}$	$\sigma_{ij}(\omega)$	BC
AB	1	0	0
AC	2	0	0
AD	1	0	0
BE	1	0	0
BD	2	0	0
CD	1	0	0
			<u>0</u>

a) 2 mode matrix.

let people be

Allison : P<sub>1</sub>

Drew : P<sub>2</sub>

Elliot : P<sub>3</sub>

Keith : P<sub>4</sub>

Ross : P<sub>5</sub>

Sarah : P<sub>6</sub>

Events / Affiliation

SSN

Party 1 : a<sub>1</sub>

Party 2 : a<sub>2</sub>

Party 3 : a<sub>3</sub>

	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>	P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>
P <sub>1</sub>	0	0	0	0	0	0	1	0	1
P <sub>2</sub>	0	0	0	0	0	0	0	1	0
P <sub>3</sub>	0	0	0	0	0	0	0	1	0
P <sub>4</sub>	0	0	0	0	0	0	0	0	1
P <sub>5</sub>	0	0	0	0	0	0	1	1	1
P <sub>6</sub>	0	0	0	0	0	0	1	1	0
a <sub>1</sub>	1	0	0	0	1	1	0	0	0
a <sub>2</sub>	0	1	1	0	1	1	0	0	0
a <sub>3</sub>	1	0	1	1	1	0	0	0	0

# One-mode

SSN

	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>	P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>
P <sub>1</sub>	2	0	1	1	2	1
P <sub>2</sub>	0	1	0	0	1	1
P <sub>3</sub>	1	2	1	2	1	
P <sub>4</sub>	1	0	1	1	1	0
P <sub>5</sub>	2	1	2	1	3	2
P <sub>6</sub>	1	1	1	0	2	2

## PART-A

SSN

19

1) Social network analysis is the conceptualisation of social networks & their representation in various forms like graphs / sociomatrices. It also involves the analysis of these structures to form general inferences on the nature of social groups.

- 2) → Highly connected subgroups with sparse intergroup connection.
- Two mode-affiliation network rep by bipartite graph.
- core-periphery networks: Core has highly connected nodes; peripheral nodes are not connected directly, but with core; minimises error.

## Global Views of social structure

powerful, but not central:

~~Actor~~: has a lot of connections.  
A french actor, in a web of actors might be powerful within her own community, but not central, since they do not act as a connection to the overall web.

SSN<sup>21</sup>

140

Similarly, in a film network, a director might be central, in knowing/connecting other people. However, they might not be powerful themselves.

b) Three schemes:

→ Sociometric Analysis:

~~the~~ equivalence of blocks.

→ ~~the~~ Relational Algebra:

Role & Positional Analysis

→ Graph Theoretic

for dyad/triad/closeness analysis.

Total Marks

1/2023