

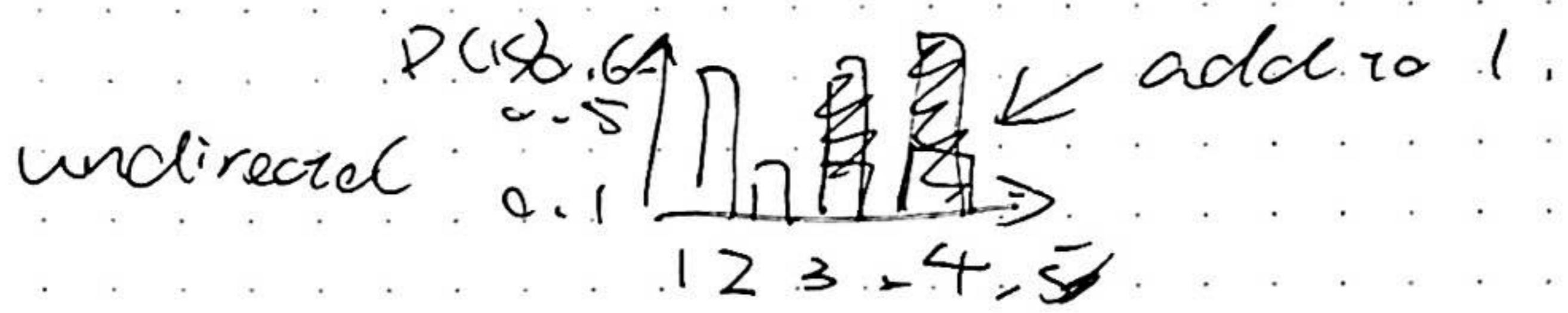
## Properties of graphs.

### (1) Degree distribution $P(k)$ :

Probability that a randomly chosen node has degree  $k$ .

$N_k = \#$  nodes with degree  $k$ .

$$P(k) = N_k / N \rightarrow \text{normalized histogram.}$$



directed: one for in-degree one for out degree.

### (2) Paths in a graph

A path is a sequence of nodes in which each node is linked to the next one.

$$P_n = \{i_0, i_1, i_2, \dots, i_n\} \quad P_n = \{(i_0, i_1), (i_1, i_2), (i_2, i_3), \dots, (i_{n-1}, i_n)\}$$

→ distance (shortest ~~graph~~<sup>path</sup>, geodesic) between a pair of nodes → no. of edges along the shortest path connecting the nodes.

\* if two nodes are unconnected,  $d$  is  $\infty$

for directed graphs → paths need to follow the direction of arrows ⇒ distance is not symmetric

N.B. cf hc-B.

→ network diameter: the maximum (shortest path) distance between any pair of nodes in a graph.  
(but this is a fragile definition as one string of nodes can skew this metric).

an alternative: average path length for a connected graph or a strongly connected direct. graph.

$$\bar{h} = \frac{1}{2E_{\max}} \sum_{i,j \neq i} h_{ij}$$

•  $h_{ij}$  is the distance from node  $i$  to node  $j$ .

•  $E_{\max}$  is the max. no. of edges (total no. of node pairs) =  $n(n-1)/2$ .

→ we compute the average only over the connected pairs of nodes. (ignore infinite length paths)

so -

(3). clustering coefficient (undirected non-directed graph → inspired social network)

↳ how connected are  $i$ 's neighbours to each other?

node  $i$  with degree  $k_i$ :

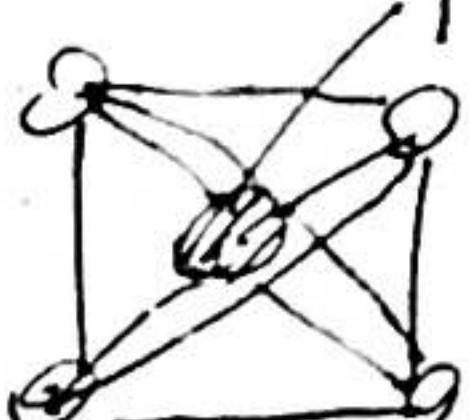
$$C_i \in [0, 1]$$

note:  $\text{represents undirected for nodes with degree 0 or 1.}$

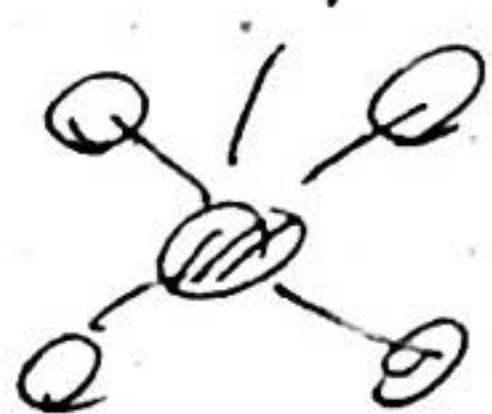
$$C_i = \frac{2e_i}{k_i(k_i-1)}$$

where  $e_i$  is the no. of edges between the neighbours of node  $i$ .

↗ all possible edges between neighbours (max no. of edges between the  $k_i$  neighbours.)



$$C_i = 1$$



$$C_i = C$$

average clustering coefficient:

$$C = \frac{1}{N} \sum_i C_i$$

(friends of mine - they tend to know each other ( $C_i \rightarrow \text{high}$ )).

#### (4) Connecting

Size of the largest connected component.

↳ largest set where any two ~~vertices~~ can be joined by a path.

To find connected components

1. Start from random node and perform breath first search (BFS)
2. Label the nodes that BFS visits.
3. If all nodes are visited, the network is connected.

Otherwise find an unvisited node and repeat BFS.

Once we know these four metrics for a network, what next? Should we be surprised?

→ we need a model!

Simplest model of graphs:

Erdős-Renyi Random Graphs.

Two variants:

- $G_{n,p}$ : undirected graph on  $n$  nodes, where each edge  $(u, v)$  appears iid with probability  $p$

$P$

- $G_{n,m}$ : undirected graph with  $n$  nodes, and  $m$  edges picked uniformly at random.

for the Gnp model we can get:

degree distribution: binomial for Gnp.

Let  $P(k)$  denote fraction of nodes

with degree  $k$

$$P(k) = \binom{n-1}{k} p^k (1-p)^{n-1-k}$$

A bell-shaped curve representing the degree distribution  $P(k)$ . The x-axis is labeled  $k$  and the y-axis is labeled  $P(k)$ . The curve is symmetric and centered around the mean.

$p$   $\rightarrow$  prob. of having  $k$  edges  
 select  $k$  nodes out of the rest  
 of the nodes

$(1-p)$   $\rightarrow$  the rest of the  $n-1-k$  edges

Well curve like Gaussian

mean

$$\text{mean } \bar{k} = p(n-1)$$

$$\frac{\sigma}{\bar{k}} = \left[ \frac{1-p}{p} \frac{1}{n-1} \right]^{1/2}$$

$$\text{Variance } \sigma^2 = p(1-p)(n-1)$$

$$\frac{1}{(n-1)^{1/2}}$$

$\bar{k}$

as the network size increases, the distribution becomes increasingly narrow.  $\Rightarrow$  degree of a node is in the vicinity of  $\bar{k}$ .

clustering coefficient

$c_i = \frac{2e_i}{k_i(k_i-1)}$  no. of edges between  $\downarrow$   
 edges in  $\downarrow$   
 $e_i$  is  $\binom{k_i}{2}$  if  $G_{np}$  appears i.i.d. with  $p$ .  
 neighbours

$$E[e_i] = p^{\binom{k_i}{2}}$$

$$\text{Then } E[c_i] = \frac{p \cdot k_i(k_i-1)}{\binom{k_i}{2}} = p = \frac{k^2}{n-1}$$

clustering coefficient  $\approx \frac{k}{n}$

of a random graph is small

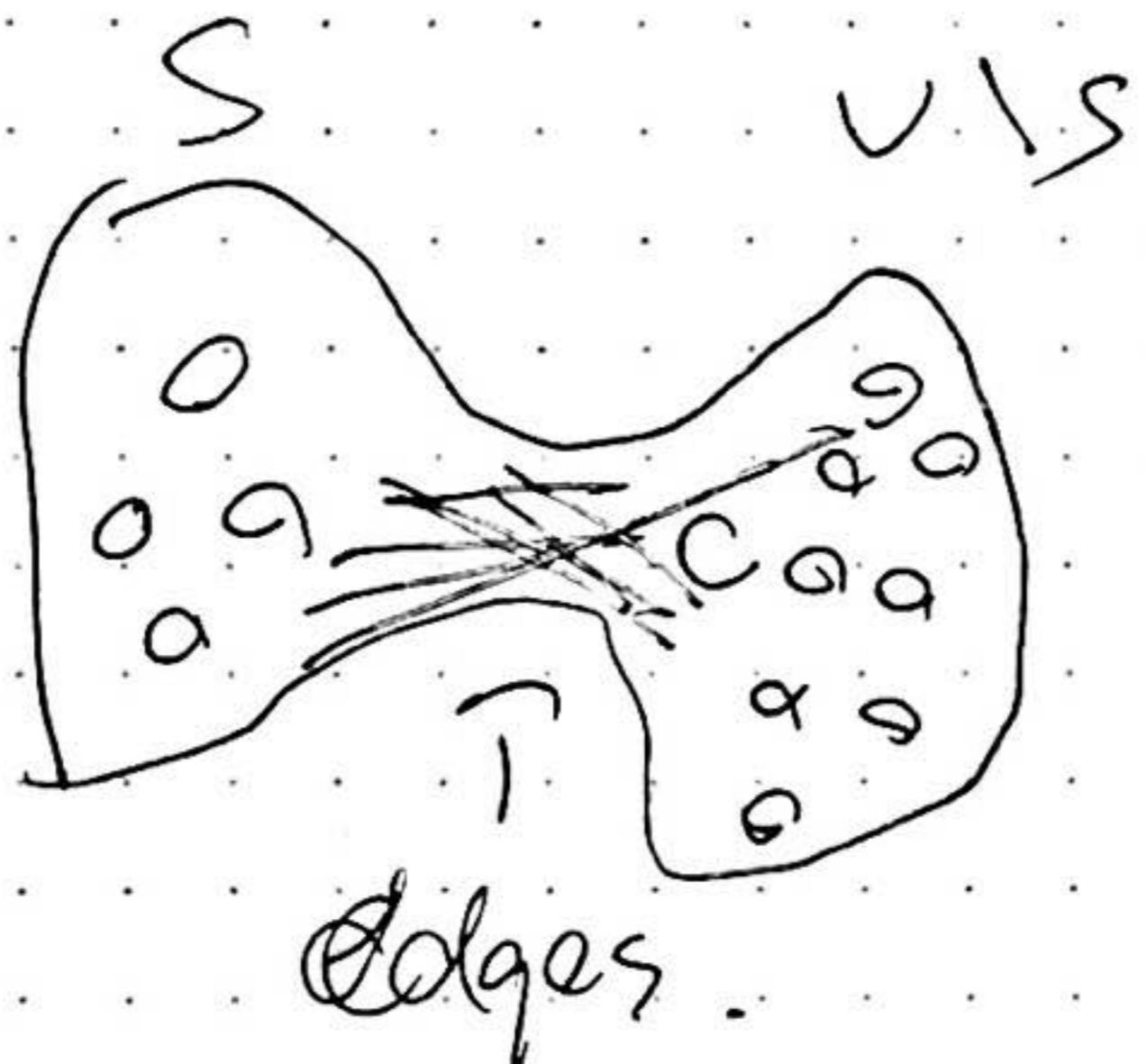
$n \rightarrow \text{large } E[c_i] \rightarrow c$

Def Expansion

Graph  $G(V, E)$  has expansion  $\alpha$  if

$\forall S \subseteq V$  # of edges leaving  $S \geq \alpha \cdot \min(|S|, |V \setminus S|)$

$$\hookrightarrow \alpha = \min_{S \subseteq V} \frac{\# \text{ edges leaving } S}{\min(|S|, |V \setminus S|)}.$$



### Expansion

In a graph on  $n$  nodes where expansion, for all pairs of nodes, there is a path of length  $O((\log n)/\alpha)$

Random graph  $G_{n,p}$ :

for  $\log n > np > c$ ,  $\text{diam}(G_{n,p}) = O(\log n / (\log(np)))$ .

random graphs have good expansions  
so it takes a logarithmic no. of steps for BFS to visit all nodes.

↳ if things / properties are different than to the Gnp model, we should be surprised.

The small-world model:

high clustering + ~~large~~ small diameter?



the 'controversy'

consequence of expansion

short paths:  $O(\log n)$  ← the smallest diameter we can see if we keep the degree constant. But the clustering is low.

BUT networks (like social networks) have 'local' structure:

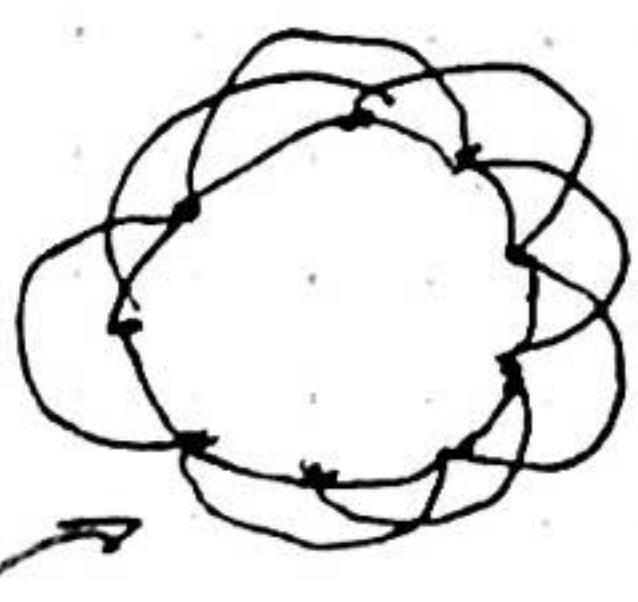
triadic closure: friend of a friend is my friend.

High clustering but diameter is also large.

Can we have both?

Small world model

(1) Start with a lower-dimensional regular lattice which has high clustering coefficients.



(2) Rewire/introduce randomness ('shortcuts') add/remove edges to create shortcuts to join remote parts of the lattice. For each edge - with prob.  $p$  move the other end wire to a random node.

Kronecker Graph model: generating large realistic graphs.

Idea: recursive graph generation.

Self-similarity: object is similar to a part of itself: the whole has the same shape as one or more of the parts.

Kronecker product of matrices A and B is given by:

$$C = A \otimes B = \begin{pmatrix} a_{1,1}B & a_{1,2}B & \dots & a_{1,m}B \\ a_{2,1}B & a_{2,2}B & \dots & a_{2,m}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1}B & a_{n,2}B & \dots & a_{n,m}B \end{pmatrix}_{N \times M \times K \times L}$$

Each  $B$  is expanded in place.

Kronecker product of two graphs  $\rightarrow$  Kronecker product of their adjacency matrices.

Kronecker graph is obtained by growing sequence of graphs by iterating the Kronecker product over the initiator matrix  $k_1$ .

$$k_i^{[m]} = k_m = k_1 \otimes k_1 \otimes \dots \otimes k_1 = k_{m-1} \otimes k_1$$

(can only generate graphs with  $n^2$  nodes) is not much a shortcoming when the graph is large enough.

Stochastic Kronecker Graph : (bring stochastic aspect to the generative model)

Create  $W \times N$ , probability matrix  $\Theta_1$ . compute the  $k$ th Kronecker power  $\Theta_k$ . for each entry  $P_{uv}$  of  $\Theta_k$  include an edge  $(u-v)$  in  $G_k$  with probability  $P_{uv}$

$\mathcal{P}$

this is an expensive operation  
when there are 1 million nodes.

$O(n^2)$

$\approx$  no. of nodes

the fix : exploit the sparsity of adj. matrix  
 $\hookrightarrow$  exploit the recursive structure of Kronecker graphs.  
drop edges onto the graph one by

Fast Kronecker ~~graph~~ <sup>one</sup> generator algorithm

for generating directed graphs.

Insert 1 edge on graph  $G$  on  $n=2^m$  nodes:

Create normalised matrix  $L_{uv} = \Theta_{uv} / (\sum_{op} \Theta_{op})$

for  $i=1 \dots m$ .

- start with  $x=0, y=0$

- pick a row/column  $(u-v)$  with prob.  $L_{uv}$

- descend into quadrant  $(u,v)$  at level  $i$  of  $G$ .

$$x+ = u \cdot 2^{m-i}, y+ = v \cdot 2^{m-i}$$

Add an edge  $(x,y)$  to  $G$

## SNAPsnap.py

### Motifs and structured roles in network

Subgraphs - break graphs down into non-isomorphic subgraphs and "count" how often they occur in the graph.

→ network significance profile : a feature vector with values for all subgraph types.

+ = over represented

- = under represented

insights: networks from same domain have similar significance profiles.

Network motifs : recurring, significant patterns of

• helps understand how networks work



more frequent than expected

found many times compared to randomly generated

and reaction of the network

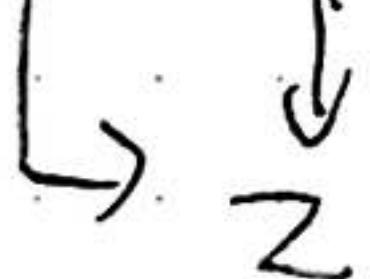
small connected subgraph

in a given situation.

compared to

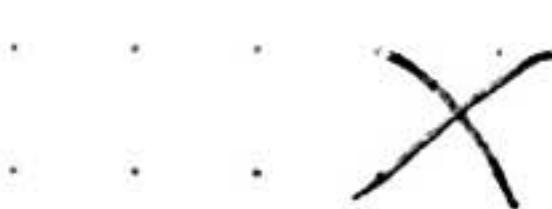
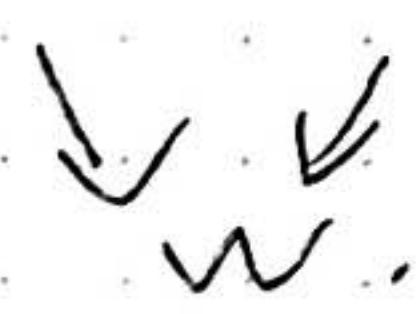
→ feed forward loops found in networks

of neurons where they neutralize biological noise



→ parallel loops : found in feed nets.

X  
↙ ↘ → Single-input modules : found in gene control networks.  
Y    Z

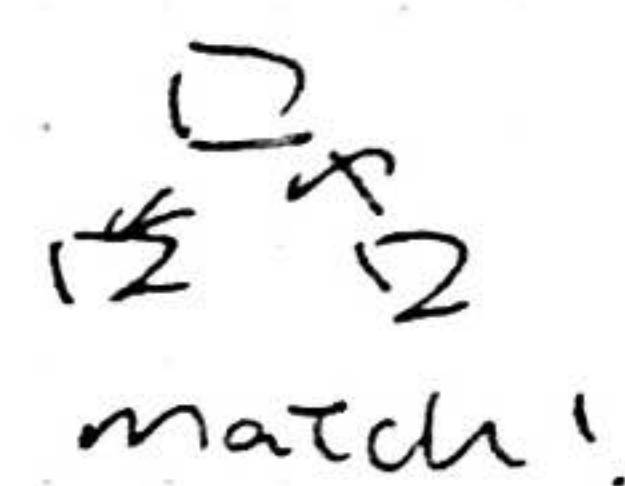
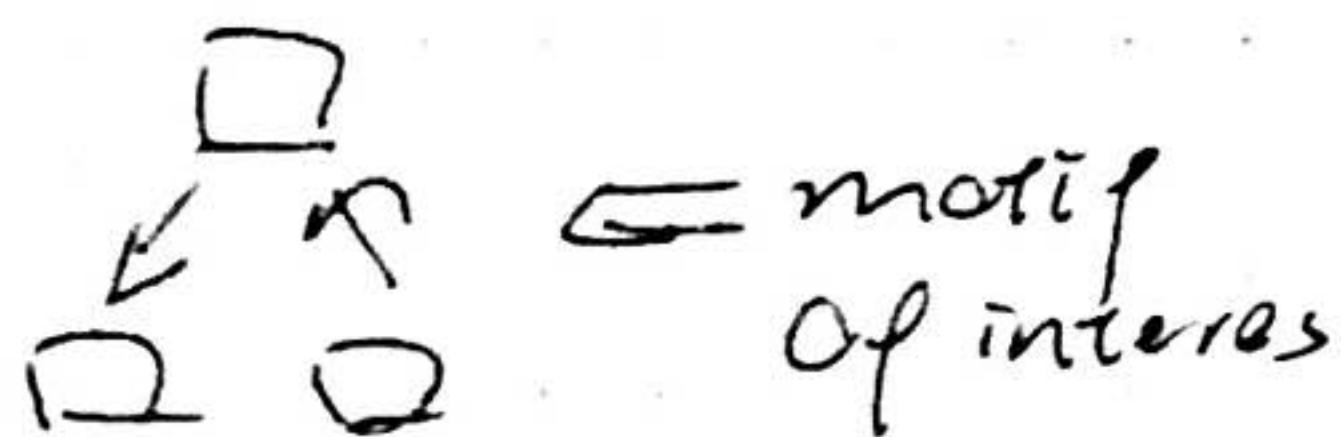


a b c d

Single input module.

## Induced subgraphs

a graph formed from a subset  $X$  of the vertices of graph  $G$  and all of the edges connecting pairs of vertices in subset  $X$ .



## Recurrence

Allow overlapping of motifs.

## Significance

subgraphs that occur in a real network more often than in a random network have functional significance.

### Z-score

$$Z_i = (N_i^{\text{real}} - \bar{N}_i^{\text{rand}}) / \text{std}(N_i^{\text{rand}})$$

### Network significance profile (SP)

$$SP_i = Z_i / \sqrt{\sum_j Z_j^2} \quad \leftarrow \text{a vector of normalized } Z\text{-scores}$$

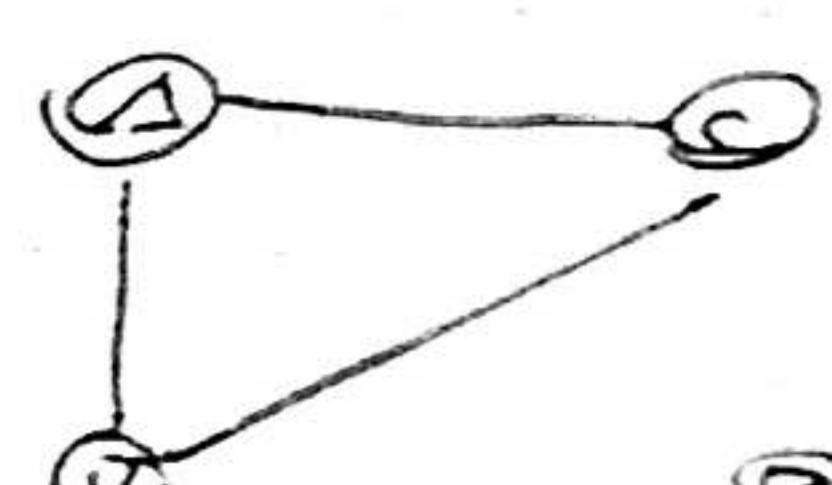
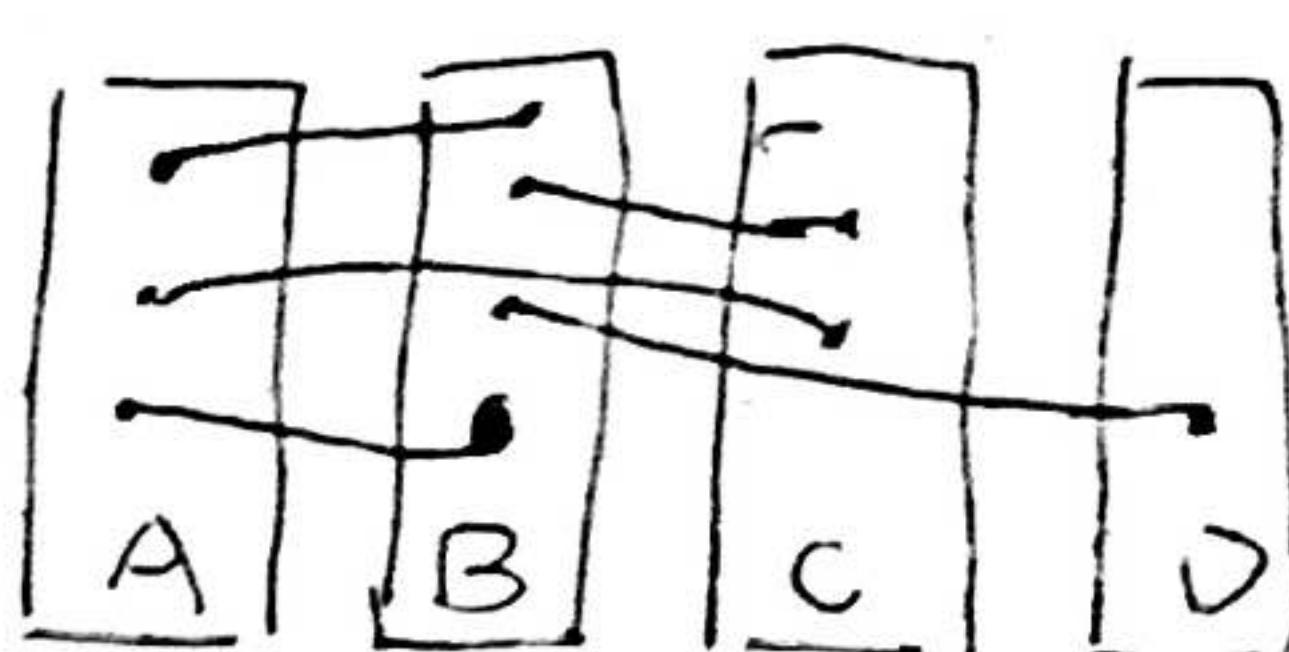
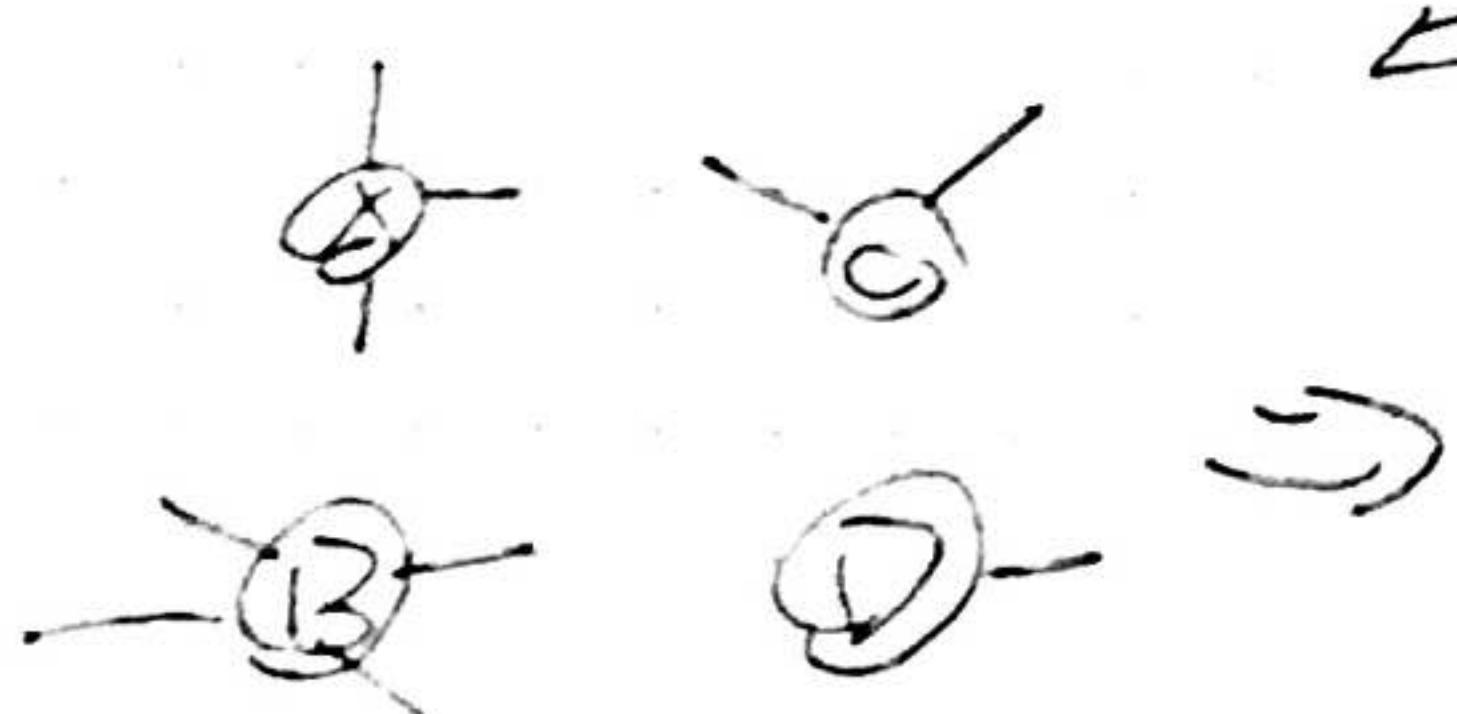
so how to generate the random graph for comparison?

generate a random graph given degree sequence,

$k_1, k_2, \dots, k_N$ .

### (Configuration model)

(ass of local structures in this way there are given or large graphs.)



nodes with spokes

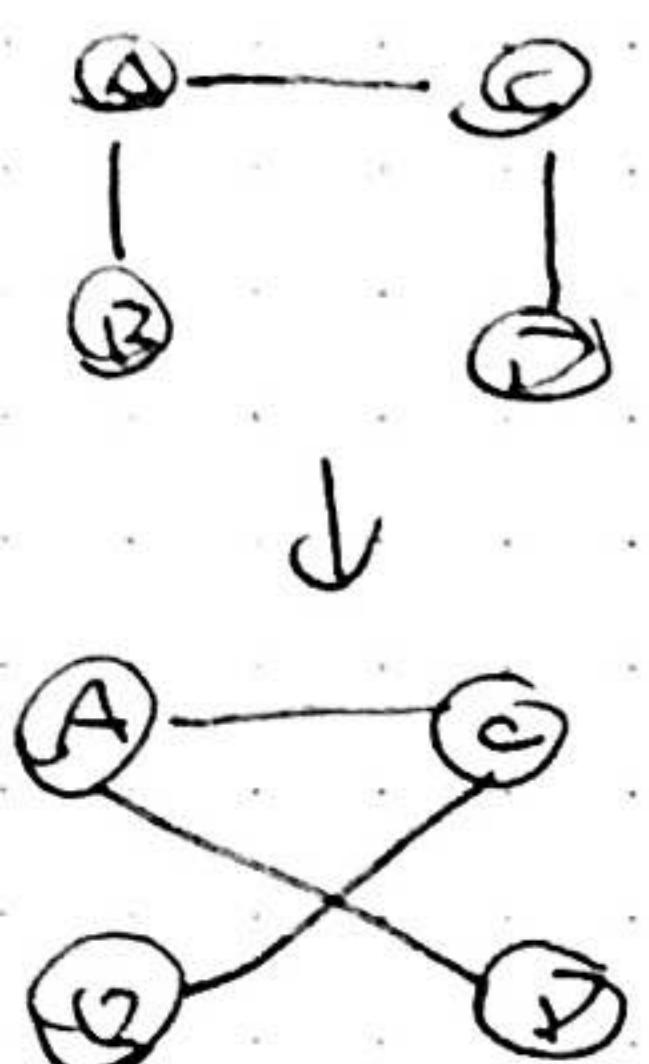
randomly pair up  
'mini'-nodes

Resulting  
graph

Alternative for spokes : switching

start from a given graph  $G$

Repeating the switching step  $Q$ : (El times)



- select a pair of edges  $A \rightarrow B$ ,  $C \rightarrow D$  are random.
- exchange the end points to give  $A \rightarrow D$ ,  $C \rightarrow B$ .
- Exchange edges, only if no multiple edges or self-edges are generated

Result: a randomly rewired graph

\*  $Q$  needs to be large enough ( $Q = 100$ ) for the process to converge.

Graphlets: node feature vectors

connected non-isomorphic  
subgraphs

induced subgraphs of any frequency

→ used to obtain a node-level subgraph metric

Automorphism Orbits

$\alpha \hookrightarrow$  takes into account the symmetries of a subgraph.  
Graphlet Degree Vector (GDV): a vector with  
the frequency of the node in each orbit position  
(counts # (graphlets) that a node touches are  
in a particular orbit).

→ provides a measure of a node's local  
network topology.

Finding size- $k$  motifs/graphlets requires  
solving two challenges:

- 1) enumerating all size- $k$  connected subgraphs
- 2) counting # (occurrences of each subgraph type).

NP-complete!

ESU algorithm:

Two sets:  $V_{\text{subgraph}}$ : currently constructed subgraph (motif)

$V_{\text{extension}}$ : set of candidate nodes to extend the motif.

Idea: starting with a node  $v$ , add those nodes  $u$  to  $V_{\text{extension}}$  such that have two properties.

→  $u$ 's node-id must be larger than that of  $v$ .

→  $u$  may only be neighbored to some newly added node  $w$  but not of any node already in  $V_{\text{subgraph}}$ .

ESU is implemented as recursive function:  
it can be displayed as a tree-structure  
of depth  $k$  called the ESU-tree.

use ESU-tree to count subgraphs.

↳ classify subgraphs placed in the ESU-Tree  
Leaves into non-isomorphic size- $k$  classes.

## Graph Isomorphism

Graphs  $G$  and  $H$  are isomorphic if there exists a bijection  $f: V(G) \rightarrow V(H)$  such that:

any two nodes  $u$  and  $v$  of  $G$  are adjacent in  $G$ , iff  $f(u)$  and  $f(v)$  are adjacent in  $H$ .

## Structure roles in network

Roles are 'functions' of nodes in a network and they are measured by structural behaviors.

### Roles vs. Groups in Networks

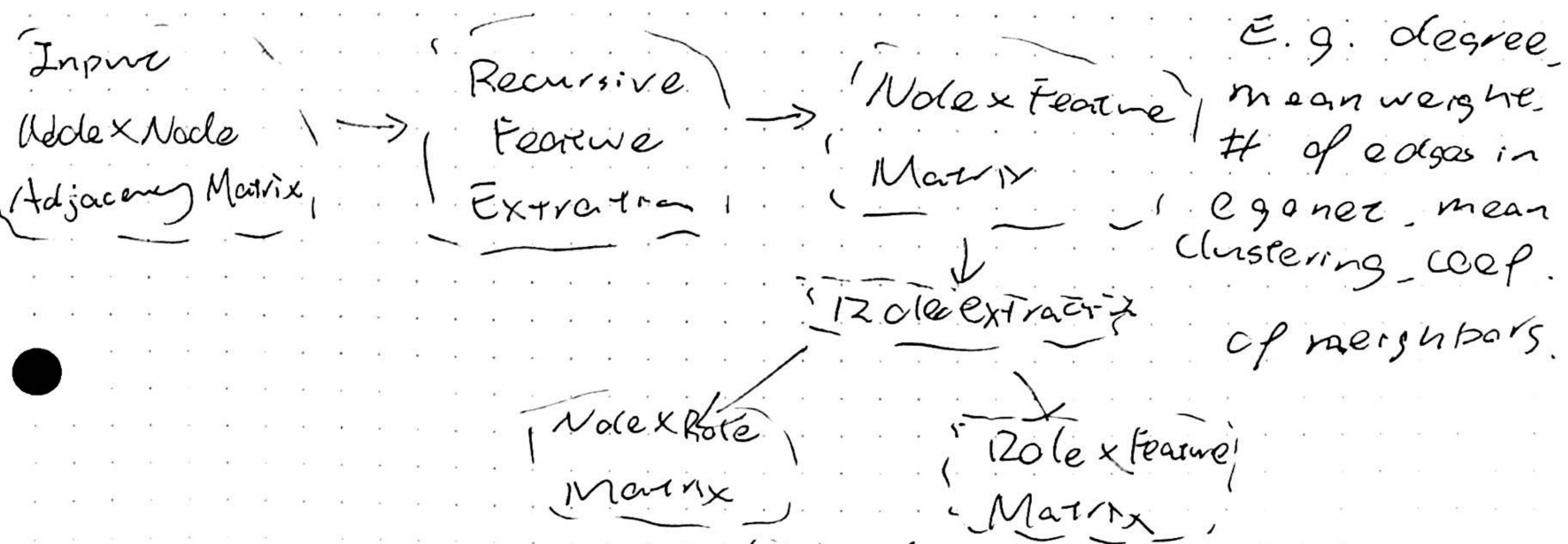
Nodes with the same role need not be in direct, or even indirect interaction with each other. (unlike groups / communities).

roles: a group of nodes with similar structural complementary properties.

communities/groups: a group of nodes that are well connected to each other.

RoleX: automatic discovery of nodes' structural roles in networks.

→ unsupervised learning.



Recursive Feature Extraction: turns network connectivity into structural features.

<sup>idea</sup>  
idea: aggregate features of a node and use them to generate new recursive features.

Base set of a nodes neighborhood features.

→ local features

→ ego net features

↓  
includes the node, its neighbors and any edges in the induced subgraph on these nodes.

Community Structure in Network

Granovetter's theory suggests that networks are composed of tightly connected sets of nodes.

Network communities: clusters - groups - modules, sets of nodes with lots of internal connections and few external ones (to the rest of the network).

sets of tightly connected nodes

Modularity Q

a measure of how well a network is partitioned into communities. Given a partitioning of the network into groups disjoint  $S \in S$ .

$$Q \propto \sum_{S \in S} [\underbrace{(\text{#edges within groups})}_{\uparrow} - (\text{expected # edges within groups})]$$

need a null model

Null model: Given real  $G$  on  $n$  nodes and  $m$  edges. Construct revised network  $G'$

- same degree distribution but uniformly random connections

Differences before  $\rightarrow$  consider  $G'$  as a multigraph

The expected no. of edges between nodes  $i$  and  $j$  of degrees  $k_i$  and  $k_j$  equals:

$$k_i \cdot \frac{k_j}{2m} = \frac{k_i \cdot k_j}{2m} \quad (\sum_{i \in G} k_i = 2m)$$

$$Q(G, S) = \frac{1}{2m} \sum_{S \subseteq G} \sum_{i \in S} \sum_{j \in S} (A_{ij} - \frac{k_i k_j}{2m})$$

↑  
normalizing

const.:  $-1 \leq Q \leq 1$

Modularity values take range

if positive  $\rightarrow$  community structure (0.3 - 0.7)

Idea: We can identify communities by maximizing modularity

(how do we find them?)

Louvain Algorithm (greedy algorithm.)

High level overview:

Phase 1: Partition: modularity is optimized by allowing only local changes to node-communities memberships.

Phase 2: Aggregation: the identified communities are aggregated into supernodes to build a new network.

Go to Phase 1.

(The passes are repeated iteratively until no increase in modularity is possible)

Louvain; 1<sup>st</sup> phase (Partitioning).

→ put each node in a graph into a distinct community (one node per community).

For each node  $i$ , the algorithm performs two calculations

- ① compute the modularity delta ( $\Delta Q$ ) when putting node  $i$  into the community of some neighbor  $j$ .
  - ② Move  $i$  to a community of node  $j$  that yields the largest gain in  $\Delta Q$ .

Phase I runs until no movement yields a gain  
(a local maxima of the modulus is attained)

\* note that the output of the algorithm depends on the order in which the nodes are considered

$$\Delta Q = \Delta Q(i \rightarrow c) + \Delta Q(D \rightarrow i)$$

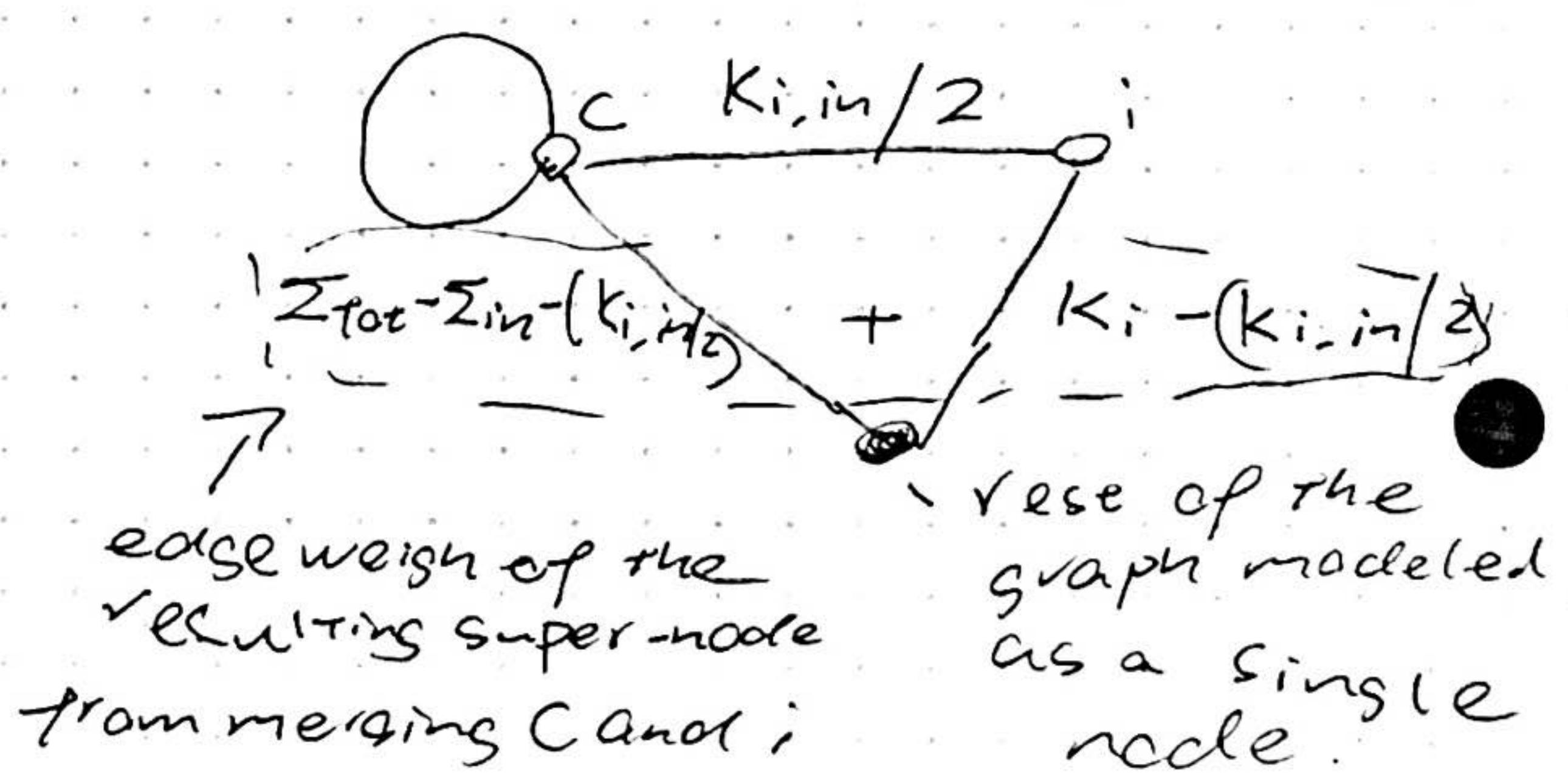
→ move point i  
into community C → move Point i  
out of community D.

$$\Delta Q(i \rightarrow C) = \left[ \frac{\sum_{in} k_{i,in}}{2m} - \left( \frac{\sum_{tot} k_i}{2m} \right)^2 \right]$$

$$= \left[ \frac{\sum i^2}{2m} - \left( \frac{\sum x_i}{2m} \right)^2 - \left( \frac{k_i}{2m} \right)^2 \right]$$

$\overrightarrow{P_i}$   
modular  
contribution  
before merging  
node i.

modulating  
Contribution after  
merging node;



## Louvain: 2<sup>nd</sup> phase (Restructuring)

The communities obtained in the first phase are contracted into super-nodes, and the network is created accordingly.

→ super-nodes are connected if there is at least one edge between the nodes of the corresponding communities.

→ The weight of the edge between two super-nodes is the sum of the weights from all edges between their corresponding communities.

Phase I is then run on the super-node network.

## Overlapping communities

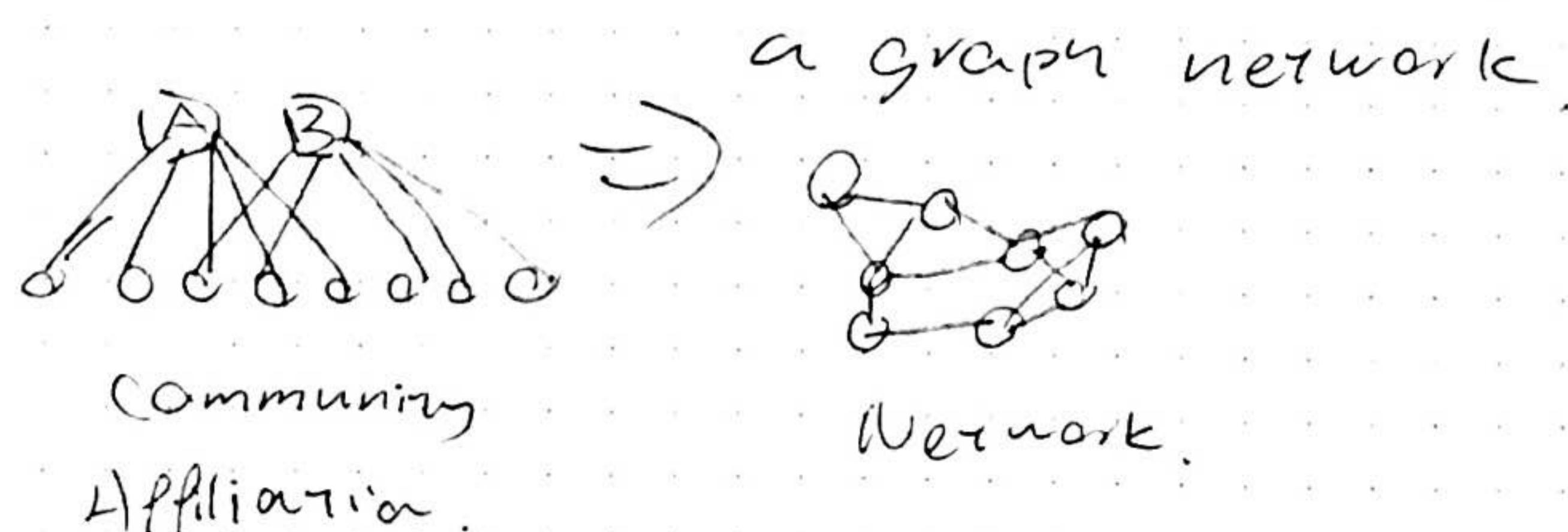
1) Define a generative model for graphs that is based on node community affiliations

- Community Affiliation Graph Model (AGM)

2) Given graph  $G$ , make the assumption that  $G$  was generated by AGM.

- Find the best AGM that could have generated  $G$ . And this way we discover overlapping communities.

AGM. Generative process: from bipartite graph to



Given parameters ( $V, C, M, \{p_{ij}\}$ ).

- Nodes in community  $c$  connect to each other by flipping a coin with probability  $p_{cc}$ .

- Nodes that belong to multiple communities have multiple coin flips.
- If they 'miss' the first time, they get another chance through the next community

$$P(u, v) = \prod_{c \in \text{Mu} \cap \text{Mv}} (1 - p_c)$$

Detecting communities:

given a graph, find the model  $\tilde{F}$

1) Affiliation graph  $M$

2) No. of communities  $C$

3) Parameters  $p_c$ .

Graph fitting  $\rightarrow$  how to estimate  $F$  given  $G$ .

$$\text{MLE} : \arg \max_F P(G|F)$$

$$P(G|F) = \prod_{(u,v) \in G} P(u, v) \prod_{(u,v) \notin G} (1 - P(u, v))$$

Likelihood of edges      Likelihood of edges  
in light of probabilities in the graph      not in the graph.  
we should relax hard membership assignments in  $M$

$\hookrightarrow$  BigRAM Model

prob. of nodes  $u, v$  linking is  $\propto$  to the strength of shared memberships.

$$P(u, v) = 1 - \exp(-\mathbf{f}_u \cdot \mathbf{f}_v^T)$$

$$L(F) = \sum_{(u,v) \in E} \log(1 - \exp(-\mathbf{f}_u \cdot \mathbf{f}_v^T)) - \sum_{(u,v) \notin E} \mathbf{f}_u \cdot \mathbf{f}_v^T$$

Optimization:

1. Start with random  $F$

2. update  $\mathbf{f}_u$  for node  $u$  while fixing the memberships of all other nodes in the cluster in the class