

SC435

Introduction to Complex Networks

Ch. 6

6.1, 6.2, 6.3, 6.4, 6.4.1, 6.4.2, 6.6, 6.9
6.10, 6.10-1, 6.11, 6.11.1, 6.14.

Problems

6.1, 6.2, 6.3, 6.4, 6.8

Ch. 7

7.1, 7.2, 7.3, 7.4, 7.5, 7.6, 7.7

Problems

7.1, 7.2, 7.3, 7.4

Similarity

Similarity in networks

- Nodes that are connected to each other in (*social*) networks tend to be similar in their features.
 - friends may have similar features.
 - similar webpages may link to similar others
 - recommendation systems
 - circle of friends tell us about the person
- Assortativity: technical name of measuring similarity. We can make predictions about a person's qualities by inspecting their neighbors.
 - homophily: like attracts likes or due to social influence??
 - segregation and polarization of online communities on social media (Echo Chamber)
 - Degree assortativity: core-periphery structure.
- Disassortative: converse of assortative

questions

- In what way can the vertices in a network be similar?
- How can we quantify similarity? 

- Random Variable
- Probability Distribution
- $\langle x^n \rangle = \sum_x x^n p\{X=x\}$.
- Jointly distributed Random Variables.
Covariance, Independence, Correlation

Mean
Variance
Skewness
Kurtosis

$x_1 \dots x_n$

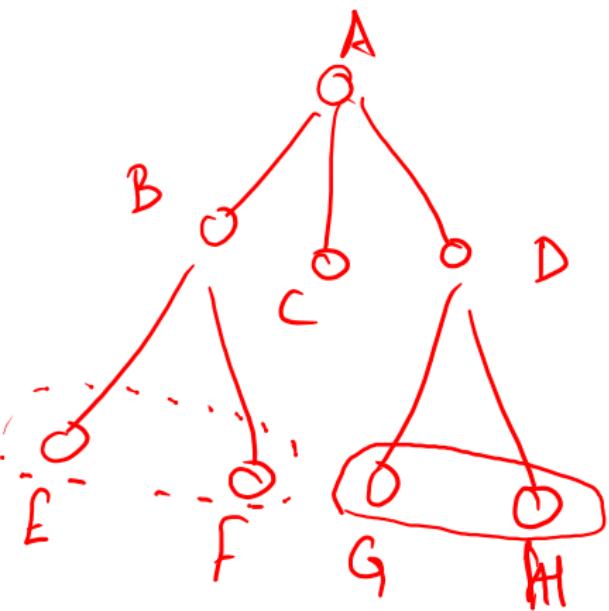
$$E[x] = \frac{1}{n} \sum_i x_i$$

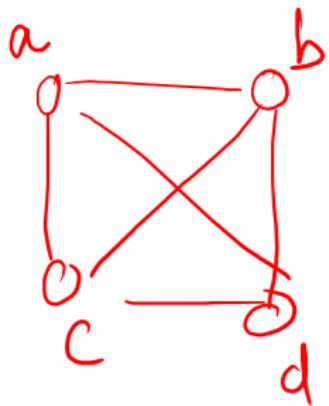
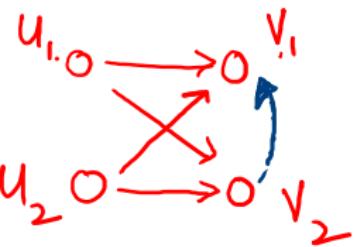
$$\text{Var}[x] = E[(x - E[x])^2].$$

$$\text{Cov}(x, y) = E[(x - E[x])(y - E[y])]$$

Constructing measures of similarity

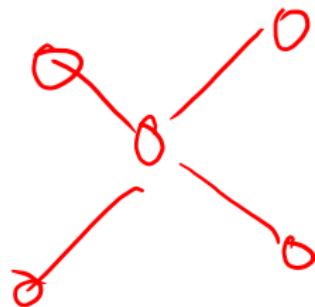
- Structural equivalence: Two vertices of a network are structurally equivalent if they share many of the same neighbors ⇒ **in social networks two nodes are similar if they share many of the same neighbors**
- Regular equivalence: Two vertices are regularly equivalent if they are equally related to equivalent others ⇒ **People with similar roles have same local neighborhood.**



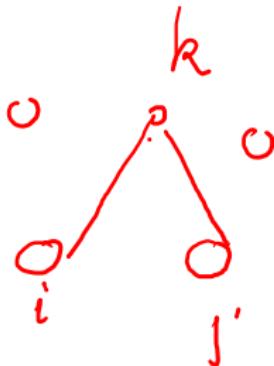


$\{c, d, b\}$
 $\{a, c, d\}$.

$$A = \begin{pmatrix} u_1 & u_2 & v_1 & v_2 \\ u_1 & 0 & 0 & 0 \\ u_2 & 0 & 0 & 0 \\ v_1 & 1 & 1 & 0 & 0 \\ v_2 & 1 & 1 & 0 & 0 \end{pmatrix} \quad |$$



$$\sigma_{ij} = \left[\quad \right]$$



$$\begin{aligned}\sigma_{ij} &= \text{number of common neighbors of vertex } i \text{ & } j \\ &= \sum_k A_{ik} A_{kj}\end{aligned}$$

Jaccard Similarity.

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

N : neighbor

$$\frac{n_{ij}}{\sum_k A_{ik} + \sum_k A_{jk} - n_{ij}} = \frac{n_{ij}}{k_i + k_j - n_{ij}} = \sigma_{ij}$$

Cosine Similarity

$$\bar{x} \cdot \bar{y} = |\bar{x}| |\bar{y}| \cos\theta.$$

$$\begin{matrix} A_{ik} & [&] \\ A_{jk} & [&] \end{matrix}$$

$$\cos\theta = \frac{n_{ij}}{\sqrt{\sum_k A_{ik}^2} \sqrt{\sum_k A_{jk}^2}} = \frac{n_{ij}}{\sqrt{r_i r_j}}$$

Expected overlap if the connection is purely at random

Similarity

Structural
Equivalence

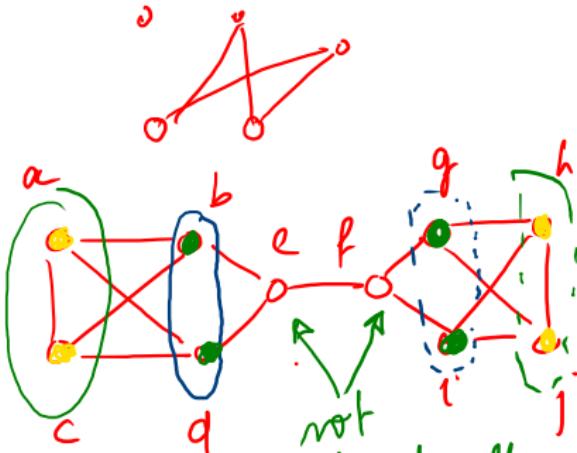
overlap
amongst neighbors

partial overlap
Hamming dist.

Regular
Equivalence

↑
Social roles

Jaccard
similarity
Cosine
Pearson



$$f \rightarrow \{e, g, i\}$$

$$e \rightarrow \{b, d, f\}.$$

Jaccard

$$\underline{n_{ij}} = \frac{\sum_k A_{ik} A_{kj}}{k_i + k_j - n_{ij}}$$

$$\text{Cosine} = \frac{n_{ij}}{\sqrt{\sum_k A_{ik}^2} \sqrt{\sum_k A_{jk}^2}}$$

$$A = \begin{bmatrix} & \xrightarrow{i} \\ & \downarrow \\ & \end{bmatrix}$$

$$\left[\begin{array}{c} \curvearrowright \\ E \end{array} \right] \left[\begin{array}{c} \curvearrowright \\ F \end{array} \right]$$

Statistical methods (comparison)

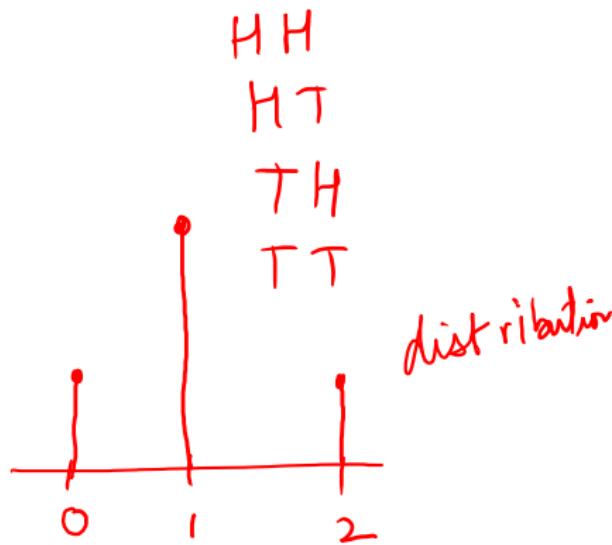
Covariance

X : number of heads.

$$P\{X=0\} = \frac{1}{4}$$

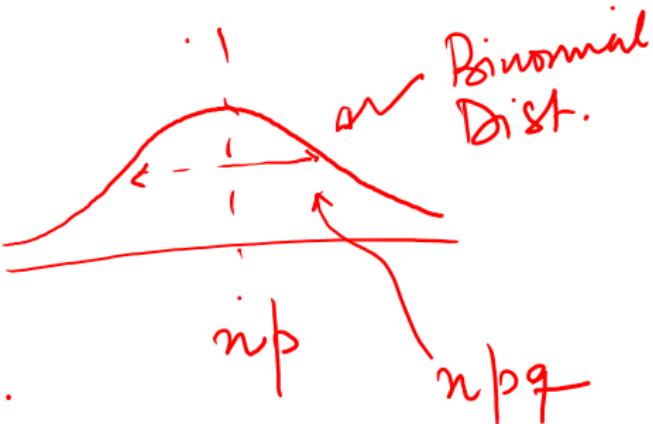
$$P\{X=1\} = \frac{2}{4}$$

$$P\{X=2\} = \frac{1}{4}$$



n flips X heads

$$\binom{n}{x} p^x (1-p)^{n-x}$$



$$E[X] = \sum_{x=0}^n x p \{ X=x \}.$$

$$E[X^2] = \sum_{x=0}^n x^2 p \{ X=x \}$$

$$\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

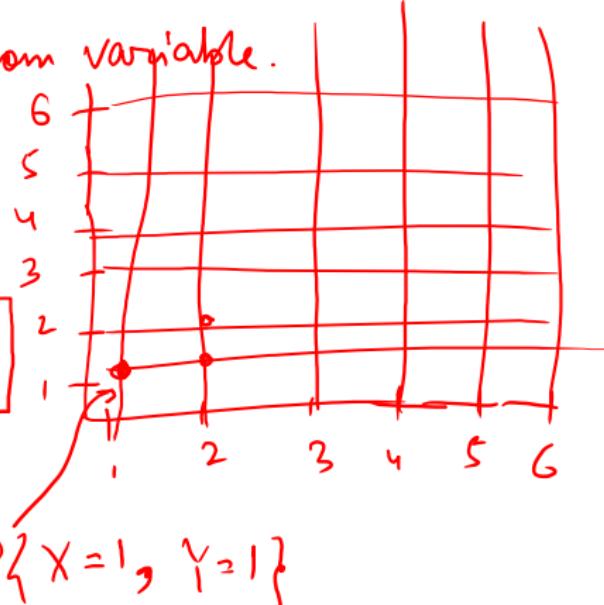
Jointly distributed random variable.

$$P\{X=x, Y=y\}.$$

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

$$\left| \begin{array}{c|cc} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \hline \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{array} \right| \quad (+)$$
$$\left| \begin{array}{c|cc} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \hline \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{array} \right| \quad (-)$$

Statistical dependence 0-random.
between X & Y .



$$\rho = \frac{E[(X - E[X])(Y - E[Y])]}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

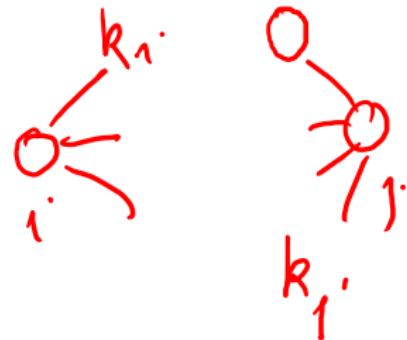
Pearson Correlation coefficient

Expected number of common neighbors if the connections are at random

A

$$\frac{k_j}{n-1} \approx \frac{k_j}{n}$$

$$\frac{k_i k_j}{n}$$



$$n_{ij} - \frac{k_i k_j}{n}$$

$$= \sum_k A_{ik} A_{kj} - \frac{k_i k_j}{n} \leftarrow \sum_k (A_{ik} - \langle A_i \rangle) (A_{jk} - \langle A_j \rangle)$$

$$A_i \rightarrow [0]$$

$$\begin{bmatrix} \dots & \langle A_i \rangle \\ \dots & \langle A_j \rangle \end{bmatrix} \boxed{\langle A_i \rangle = \frac{1}{n} \sum_k A_{ik}}$$

$$A_{ij} \rightarrow [$$

Similarity

Measures of Structural equivalence (extent)

- Count of the number of common neighbors

$$n_{ij} = \sum_k A_{ik} A_{kj} = [A^2]_{ij}$$

- Cosine Similarity:

$$\begin{aligned}\sigma_{ij} &= \frac{\sum_k A_{ik} A_{kj}}{\sqrt{\sum_k A_{ik}^2} \sqrt{\sum_k A_{jk}^2}} \\ \sigma_{ij} &= \frac{n_{ij}}{\sqrt{k_i k_j}} \quad (\text{simple, unweighted})\end{aligned}$$

- Jaccard Similarity:

$$J_{ij} = \frac{n_{ij}}{k_i + k_j - n_{ij}}$$

Pearson similarity

$$\boxed{\sigma_{ij} = n_{ij} - \frac{k_i k_j}{n}}$$

Expected common
neighbors if the
connections are drawn
at random.

$$= \sum_k A_{ik} A_{kj} - \frac{k_i k_j}{n}$$

$$= \sum_k A_{ik} A_{kj} - n \underbrace{\langle A_i \rangle \langle A_j \rangle}_{}$$

$$= \sum_k [A_{ik} A_{kj} - \langle A_i \rangle \langle A_j \rangle]$$

$$A_{ik} = A_{ki}$$

$$A^2_{ik} = A_{ik}$$

$$k_i = \sum_k A_{ik}$$

$$\langle A_i \rangle = \frac{1}{n} \sum_k A_{ik}$$

$$= \sum_k [A_{ik} A_{kj} - \overbrace{\langle A_i \rangle \langle A_j \rangle}^{\uparrow n \langle A_i \rangle \langle A_j \rangle = n \cdot \frac{1}{n} \sum_k A_{ik} \langle A_j \rangle} + \langle A_i \rangle \langle A_j \rangle - \overbrace{\langle A_i \rangle \langle A_j \rangle}^{\sum_k A_{ik} \langle A_j \rangle}]$$

$$= \sum_k [A_{ik} A_{kj} - A_{ik} \langle A_j \rangle + \langle A_i \rangle \langle A_j \rangle - A_{jk} \langle A_i \rangle]$$

$$= \sum_k (A_{ik} - \langle A_i \rangle) (A_{jk} - \langle A_j \rangle) \sim \text{Covariance}$$

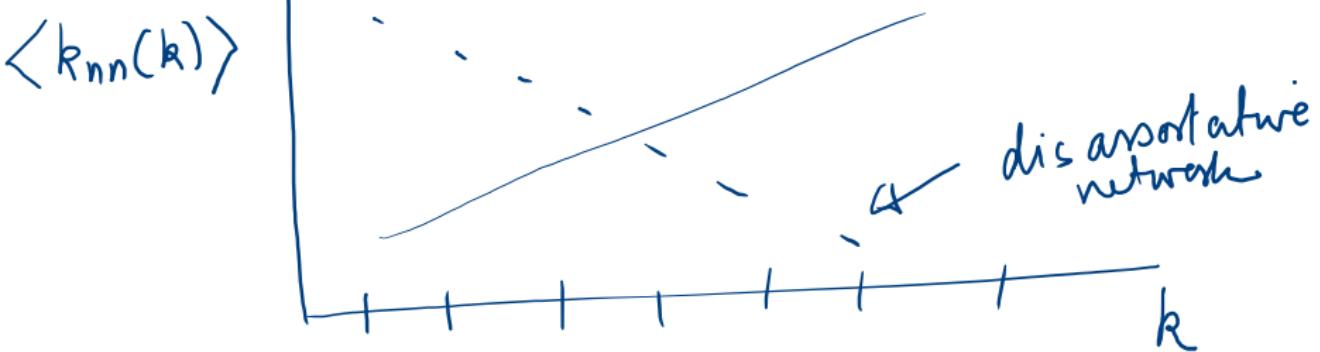
$$A = \begin{bmatrix} & & & & \end{bmatrix} \quad \begin{array}{c} \uparrow \\ A_{1k} \\ \uparrow \\ A_{2k} \\ \uparrow \\ \sigma_{12} \end{array}$$
$$\Gamma = \begin{bmatrix} & & & & \end{bmatrix} \quad \text{heatmaps.}$$

$$\frac{n_{ij}}{k_i k_j} \quad \left| \begin{array}{l} \text{hamming distance} \\ d_{ij} = \sum_k (A_{ik} - A_{jk})^2 \\ \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases} \end{array} \right.$$

$$k_{nn}(i) = \frac{1}{k_i} \sum_j A_{ij} k_j$$

↑ nearest neighbor

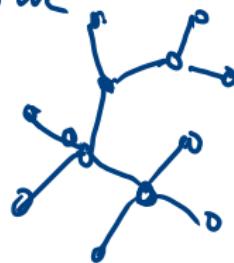
$$\langle k_{nn}(k) \rangle = \frac{1}{N_R} \sum_{i=1}^n k_{nn}(i) \delta_{k_i, k}$$

Assortative network — Core-periphery structure



Disassortative network

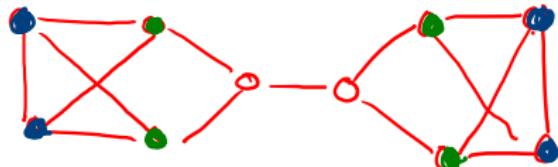


Regular equivalence

Two vertices are regularly equivalent if they are connected to vertices that are regularly equivalent.

$$\sigma_{ij} = \alpha \sum_{k,l} A_{ik} A_{jl} \sigma_{kl} + \delta_{ij}$$

$$[\sigma = \alpha A \sigma A]$$



$$\begin{aligned} \sigma_{ij} &= \alpha \sum_{k,l} A_{ik} A_{jl} \sigma_{kl} \\ &= \alpha \sum_{k,l} \underbrace{\sigma_k}_{i} \underbrace{\sigma_l}_{j} \end{aligned}$$

$$\sigma = \alpha A \sigma A + I$$

$$\sigma(t=0) = \bar{0}$$

$$\sigma(t=1) = I$$

$$\sigma(t=2) = \alpha A^2 + I$$

$$\begin{aligned}\sigma(t=3) &= \alpha^2 A^4 + \alpha A^2 + I \\ &\vdots\end{aligned}$$

$$\sigma = (I - \alpha A)^{-1}$$

$$i \xrightarrow{\sigma} k$$
$$j \xrightarrow{\sigma} b_j k$$

$$\sigma_{ij} = \alpha \sum_{ik} A_{ik} \sigma_{kj} + \delta_{ij}$$

$$\boxed{\sigma = \alpha A \sigma + I}$$

$$\sigma(t=0) = \bar{0}$$

$$\sigma(t=1) = I$$

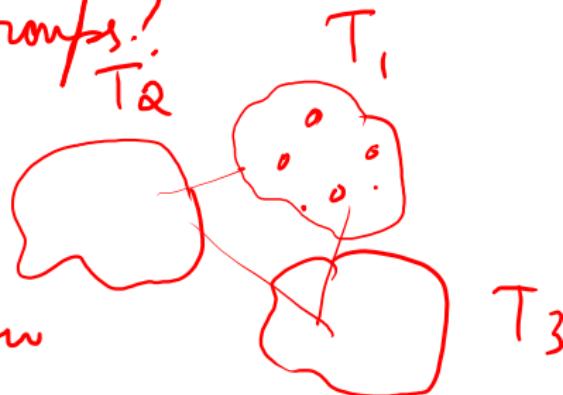
$$\sigma(t=2) = \alpha A + I$$

$$\sigma(t=3) = \alpha^2 A^2 + \alpha A + I$$

- In real networks are the connections based on similarity in terms of attributes?
 - Homophily / assortativity.
- How do we quantify it? ←
- Which ones in which groups?

Graph Partitioning →

Community Detection



$$Q = \frac{m_c - \langle m_c \rangle}{m}$$

↑

c : class
cluster
group
attribute

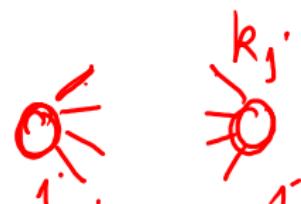
Modularity

Let c_i represent class $i = 1, \dots, n$

m_c = number of links between nodes with the same attribute

$$\frac{1}{2} \sum_{ij} A_{ij} \delta_{c_i, c_j}$$

$\langle m_c \rangle$: Expected number of links between nodes of the same type if connections are drawn at random.

$$\sum_{ij} \frac{k_i k_j}{2m} \delta_{c_i, c_j}$$


$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{c_i, c_j}$$

Measures of Structural equivalence (extent)

- Pearson correlation coefficient:

$$r_{ij} = \frac{\sum_k (A_{ik} - \langle A_i \rangle) \sum_k (A_{jk} - \langle A_j \rangle)}{\sqrt{\sum_k (A_{ik} - \langle A_i \rangle)^2} \sqrt{\sum_k (A_{jk} - \langle A_j \rangle)^2}} = \frac{\text{Cov}(\sigma_i, \sigma_j)}{\sigma_i \sigma_j}$$

- For unweighted, undirected graph

$$r_{ij} = \frac{n_{ij} - \frac{k_i k_j}{n}}{\sqrt{k_i - \frac{k_i^2}{n}} \sqrt{k_j - \frac{k_j^2}{n}}}$$

$$-1 \leq r_{ij} \leq 1$$

Measures of Structural equivalence (extent)

- normalize common neighbors by the expected number of common neighbors when they are picked at random

$$g_{ij} = \frac{n_{ij}}{\frac{k_i k_j}{n}} = \frac{n \sum_k A_{ik} A_{jk}}{\sum_k A_{ik} \sum_k A_{jk}}$$

- Euclidean (Hamming) distance:

$$d_{ij} = 1 - \frac{2n_{ij}}{k_i + k_j}$$

Similarity

Measures of Regular equivalence (extent)

- Two nodes i and j have high similarity score if they have neighbors k and l that themselves have high similarity. For an undirected network

$$\sigma_{ij} = \alpha \sum_{k,l} A_{ik} A_{jl} \sigma_{kl}$$

- Should have high self-similarity (σ_{ii})

$$\sigma_{ij} = \alpha \sum_{k,l} A_{ik} A_{jl} \sigma_{kl} + \delta_{ij}$$

- problem: assuming a null matrix as the initial condition after k iterations we have

$$\sigma^{(k)} = \sum_{r=0}^{k-1} \alpha^r A^{2r}$$

This measure of similarity is a weighted sum over the number of paths of even length between pairs of vertices.

Measures of Regular equivalence (extent)

- **Modified definition of regular equivalence:** Nodes i and j are similar if i has a neighbor of k that is itself similar to j .

$$\sigma_{ij} = \alpha \sum_k A_{ik} \sigma_{kj} + \delta_{ij}$$

- Iterating again starting with $\sigma = 0$, we get:

$$\sigma = \sum_m (\alpha A)^m = (1 - \alpha A)^{-1}$$

Similarity of two nodes is the weighted sum of the number of paths of different length that connect them.

- Bias in favor of high degree node can be removed by dividing with node degree

$$\sigma_{ij} = \frac{1}{k_i} \left[\alpha \sum_k A_{ik} \sigma_{kj} + \delta_{ij} \right]$$

Homophily and Assortative Mixing

- People have a strong tendency to associate with others whom they perceive as being similar to themselves. This property is called *homophily* or *assortative mixing*.
- Disassortative mixing: tendency to people to associate with others who are unlike them.
- Political polarization, mixing on the basis of race, obesity etc. ([Assortative](#))
- Dating networks, food web (predator-prey), economic networks (producers/consumers) ([Disassortative](#))

Similarity

Mixing by Categorical attributes

- Characterize by some numbers the value of the mixing in the network. One of the ways to do this is called (**modularity**)
- Every vertex has a label (c_i)
- How much more often do attributes match across edges than what is expected at random? (Pearson correlation coeff.)
- Modularity

$$Q = \frac{m_c - \langle m_c \rangle}{m} = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

- The total number of edges that run between nodes of the same class

$$\sum_{ij} A_{ij} \delta(c_i, c_j)$$

- expected number of edges between all pairs of vertices of the same type if edges are placed at random

$$\frac{k_i k_j}{2m}$$

Mixing by Categorical attributes

- Modularity

$$Q = \frac{m_c - \langle m_c \rangle}{m} = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

$Q = 0$ (**single class or completely random**)

$Q > 0$ (**Assortative mixing**)

$Q < 0$ (**Disassortative**)

Mixing by Categorical attributes

- Modularity matrix B

$$[B]_{ij} = A_{ij} - \frac{k_i k_j}{2m}$$

- Q is never equal to 1

$$Q_{max} = \frac{1}{2m} \left(2m - \sum_{ij} \frac{k_i k_j}{2m} \right)$$

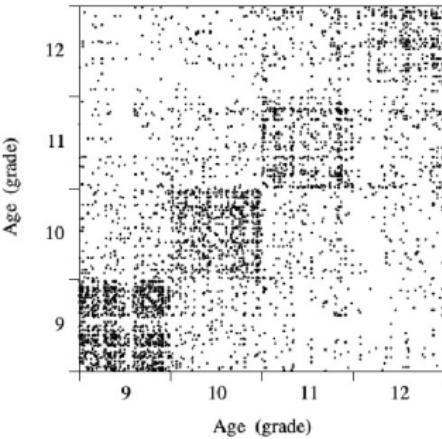
- Value of Q can be significantly smaller than 1 even for perfectly mixed networks.
So we normalize it by the maximum value ([Assortativity coefficient \$Q/Q_{max}\$](#))

Similarity

Mixing by ordered characteristics

- We can also have mixing if characteristics are approximately the same.

$$r = \frac{\sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) x_i x_j}{\sum_{ij} \left(k_i \delta_{ij} - \frac{k_i k_j}{2m} \right) x_i x_j}$$



Similarity

Assortative mixing by degree

- Assortative: core-periphery
- Disassortative: star-like structure

$$r = \frac{\sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) k_i k_j}{\sum_{ij} \left(k_i \delta_{ij} - \frac{k_i k_j}{2m} \right) k_i k_j}$$

