

# NETWORK DATA

## (M2 ETE Econometrics II)

Koen Jochmans

Toulouse School of Economics

Last revised on April 2, 2024

©2023–2024 KOEN JOCHMANS

# Outline

1 Graphs

2 Spectral clustering

3 Stochastic block model

4 Fixed-effect estimation

5 Linear social-interaction models

## Basic terminology

Let  $V$  be a set of **vertices** (nodes) with  $|V| = n$ ; can normalize  $V = \{1, \dots, n\}$ .

Let  $E$  be a set of **edges** between nodes in  $V$ . This is a set of pairs of nodes  $(i, j) \in V \times V$ .

Then  $G = (V, E)$  is a **graph**.

$G$  is **(un)directed** if the pairs in  $E$  or (un)ordered.

$G$  is a **simple graph** if a pair  $(i, j)$  has at most one edge between them.

$G$  is a **multigraph** if the latter is not the case.

Edges can carry a weight,  $w_{i,j} \geq 0$ , in which case  $G$  is a **weighted** graph.

$G$  can be represented by its (weighted)  $n \times n$  **adjacency matrix**  $\mathbf{A}$ .

In the unweighted case,

$$(\mathbf{A})_{i,j} = \begin{cases} 0 & \text{if } (i,j) \notin E \\ 1 & \text{if } (i,j) \in E \end{cases}$$

In the weighted case,

$$(\mathbf{A})_{i,j} = \begin{cases} 0 & \text{if } (i,j) \notin E \\ w_{i,j} & \text{if } (i,j) \in E \end{cases}$$

We usually exclude self links, so the diagonal of  $\mathbf{A}$  contains only zeros.

In the undirected case,  $\mathbf{A}$  is symmetric.

The (weighted) **degree** of vertex  $i$  is

$$d_i = \sum_{j=1}^n w_{i,j} = \sum_{j:(i,j) \in E} w_{i,j}.$$

In the unweighted case  $d_i$  is the number of (if the graph is directed, outgoing) edges involving node  $i$ .

The **volume** of a set of nodes  $V_1 \subseteq V$  is

$$\text{vol}(V_1) = \sum_{i \in V_1} d_i = \sum_{i \in V_1} \sum_{j:(i,j) \in E} w_{i,j}.$$

A subset  $V_1$  of  $V$  is **connected** if any two vertices in  $V_1$  can be joined by a path of which all the intermediate points also lie in  $V_1$ .

The subset  $V_1$  is a **connected component** if  $V_1$  is connected and there are no connections between  $V_1$  and  $\bar{V}_1 = V \setminus V_1$ .

The collection of sets  $V_1, \dots, V_m$  (for some  $m$ ) form a **partition** of  $G$  if each of them is connected and

$$V_i \cap V_j = \emptyset \text{ for all } i \neq j \quad \text{and} \quad V_1 \cup \dots \cup V_m = V.$$

(For an undirected graph) the **Laplacian** of  $G$  is

$$\mathbf{L} = \mathbf{D} - \mathbf{A},$$

for

$$\mathbf{D} = \text{diag}(\mathbf{d}) = \text{diag}(d_1, \dots, d_n)$$

Note that, for any  $\mathbf{v} \in \mathbb{R}^n$ ,

$$\begin{aligned}\mathbf{v}' \mathbf{L} \mathbf{v} &= \sum_{i=1}^n v_i^2 d_i - \sum_{i=1}^n \sum_{j=1}^n v_i v_j w_{i,j} = \sum_{i=1}^n \sum_{j=1}^n v_i^2 w_{i,j} - \sum_{i=1}^n \sum_{j=1}^n v_i v_j w_{i,j} \\ &= \sum_{i=1}^n \sum_{j=1}^n v_i (v_i - v_j) w_{i,j} = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (v_i - v_j)^2 w_{i,j},\end{aligned}$$

where we have used the accounting identity  $(v_i - v_j)^2 = v_i(v_i - v_j) + v_j(v_j - v_i)$ .

So,  $\mathbf{L}$  is positive semi-definite. It is singular because

$$\mathbf{L} \mathbf{1}_n = \mathbf{D} \mathbf{1}_n - \mathbf{A} \mathbf{1}_n = \mathbf{d} - \mathbf{d} = 0 \mathbf{1}_n$$

and so its smallest eigenvalue is zero (with eigenvector  $\mathbf{1}_n$ ).

The **number of connected components** in graph  $G$  equals the multiplicity of the eigenvalue zero of  $\mathbf{L}$ .

Consider first a connected graph. Let  $\mathbf{v}$  be an eigenvector of  $\mathbf{L}$  associated with eigenvalue 0. Then

$$0 = \mathbf{v}' \mathbf{L} \mathbf{v} = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (v_i - v_j)^2 w_{i,j}.$$

As  $w_{i,j} \geq 0$  this sum can only be zero if  $(v_i - v_j)^2 w_{i,j} = 0$  for all  $(i, j) \in E$ . For such  $(i, j)$  we necessarily have  $v_i = v_j$ . Further, any third node  $k$  that can be reached via the path  $i \rightarrow j \rightarrow k$  then also has  $v_k = v_j = v_i$ . Extending the argument to longer paths and invoking that the graph is connected implies that  $\mathbf{v} = \mathbf{1}_n$ .

Next take  $m$  connected components. Then  $G$  is the collection of  $m$  connected subgraphs  $G_1, \dots, G_m$ , each with their own proper Laplacian,  $\mathbf{L}_1, \dots, \mathbf{L}_m$ . To each of these the above argument applies. Further, we may always re-arrange the nodes such that  $\mathbf{L} = \text{blockdiagonal}(\mathbf{L}_1, \dots, \mathbf{L}_m)$ . The eigenvectors of  $\mathbf{L}$  are the eigenvectors of the  $\mathbf{L}_1, \dots, \mathbf{L}_m$  (supplemented with zeros as additional entries), so that the eigenvalue 0 has  $m$  distinct eigenvectors associated with it; these are  $\mathbf{1}_{V_1}, \dots, \mathbf{1}_{V_m}$ , where  $\mathbf{1}_V$  has the  $n$  entries  $(\mathbf{1}_V)_i = \{i \in V\}$ . □

The normalized adjacency matrix is

$$\mathbf{P} = \mathbf{D}^{-1} \mathbf{A}$$

(the notation presumes that  $G$  is connected; otherwise use the Moore-Penrose pseudo inverse).

Note that

$$(\mathbf{P})_{i,j} = \begin{cases} 0 & \text{if } (i,j) \notin E \\ w_{i,j}/d_i & \text{if } (i,j) \in E \end{cases}$$

and that the rows of  $\mathbf{P}$  all sum up to one.

Row  $i$  of  $\mathbf{P}$  gives the probability distribution of taking a random step through the graph  $\mathbf{G}$  when starting at node  $i$ .

$\mathbf{P}$  is the **transition matrix** of a Markov chain on  $G$ .

If  $G$  is connected (and not bipartite) then  $\mathbf{P}$  has a steady-state distribution given by

$$p_i = \frac{d_i}{\text{vol}(V)}$$

Can define (at least) two normalized version of the Laplacian matrix  $\mathbf{L}$ .

A symmetric normalization:

$$\mathbf{L}_{\text{sym}} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} = \mathbf{I}_n - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$$

and a ‘random-walk’ normalization:

$$\mathbf{L}_{\text{rw}} = \mathbf{D}^{-1} \mathbf{L} = \mathbf{I}_n - \mathbf{P}.$$

In both cases, can show similar properties as for  $\mathbf{L}$  before.

A connection between the two is as follows:

$\lambda \geq 0$  is an eigenvalue of  $\mathbf{L}_{\text{rw}}$  with eigenvector  $\mathbf{v}$  if and only if  $\lambda$  is an eigenvalue of  $\mathbf{L}_{\text{sym}}$  with eigenvector  $\mathbf{D}^{1/2}\mathbf{v}$ .

The number of zero eigenvalues continues to give the number of connected components in  $G$ .

For a connected graph  $G$  the **second smallest eigenvalue**  $\lambda_2 > 0$  is a measure of global connectivity of  $G$ .

This spectral gap is connected to the **Cheeger constant**

$$C = \min_{\{V_1 \subset V : 0 < \sum_{i \in V_1} d_i \leq \sum_{i \in \bar{V}_1} d_i\}} \frac{\sum_{i \in V_1} \sum_{j \in \bar{V}_1} w_{i,j}}{\sum_{i \in V_1} d_i}$$

through the inequalities

$$\frac{1}{2}C^2 \leq \lambda_2 \leq 2C.$$

$C \in [0, 1]$  measures how difficult it is to separate  $G$  into two disconnected components by removing edges from it.

The numerator in the definition of  $C$  is the total weight of the removed edges, the denominator is the total degree in the smallest of the two components.

Small positive  $C$  show that there is a bottleneck in  $G$ , and that the graph is only weakly connected.

# Outline

**1** Graphs

**2** Spectral clustering

**3** Stochastic block model

**4** Fixed-effect estimation

**5** Linear social-interaction models

# Motivation

Aim:

Partition set of nodes into  $m$  clusters (groups) so that they are similar within a cluster and heterogeneous across clusters.

Here, consider undirected (non-bipartite) graphs.

$m$  is taken as given.

Several types of clustering exist. Here, a flavor is given.

We will discuss spectral clustering based on the Laplacian and on the (random-walk) normalized Laplacian.

## Unnormalized clustering

Consider clustering  $V$  into clusters  $V_1, \dots, V_m$ .

One motivation: minimize the **ratio cut** of  $G$ :

$$\text{RatioCut}(V_1, \dots, V_m) = \sum_{k=1}^m \frac{\text{cut}(V_k, \bar{V}_k)}{|V_k|}$$

for

$$\text{cut}(V_k, \bar{V}_k) = \frac{1}{2} \sum_{i \in V_k} \sum_{j \in \bar{V}_k} w_{i,j} + \frac{1}{2} \sum_{i \in \bar{V}_k} \sum_{j \in V_k} w_{i,j} = \sum_{i \in V_k} \sum_{j \in \bar{V}_k} w_{i,j}$$

(using undirectedness in the last transition).

This minimization problem tries to minimize the weight assigned to edges that go across different clusters.

The weighting makes larger clusters more attractive.

For  $k = 1, \dots, m$ , introduce the  $n$ -vectors

$$(\mathbf{h}_k)_i = \begin{cases} 1/\sqrt{|V_k|} & \text{if } i \in V_k \\ 0 & \text{if } i \notin V_k \end{cases}.$$

For each  $k = 1, \dots, m$  we have (from calculations given above) that

$$\begin{aligned} \mathbf{h}'_k \mathbf{L} \mathbf{h}_k &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n ((\mathbf{h}_k)_i - (\mathbf{h}_k)_j)^2 w_{i,j} \\ &= \frac{1}{2} \sum_{i \in V_k} \sum_{j \notin V_k} \frac{1}{|V_k|} w_{i,j} + \frac{1}{2} \sum_{i \notin V_k} \sum_{j \in V_k} \frac{1}{|V_k|} w_{i,j} \\ &= \frac{\text{cut}(V_k, \bar{V}_k)}{|V_k|}. \end{aligned}$$

Therefore,

$$\text{RatioCut}(V_1, \dots, V_m) = \sum_{k=1}^m \frac{\text{cut}(V_k, \bar{V}_k)}{|V_k|} = \text{trace } \mathbf{H}' \mathbf{L} \mathbf{H}$$

where  $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_m)$ .

Note that, for  $k, \ell \in \{1, \dots, m\}^2$

$$\mathbf{h}'_k \mathbf{h}_\ell = \sum_{i=1}^n (\mathbf{h}_k)_i (\mathbf{h}_\ell)_i = \begin{cases} 1 & \text{if } k = \ell \\ 0 & \text{if } k \neq \ell \end{cases},$$

so that  $\mathbf{H}' \mathbf{H} = \mathbf{I}_m$  meaning that  $\mathbf{H}$  is an  $n \times m$  orthonormal matrix.

Minimizing the RatioCut can then be written as the discrete optimization problem

$$\min_{V_1, \dots, V_m} \text{trace } \mathbf{H}' \mathbf{L} \mathbf{H}$$

subject to  $\mathbf{H}$  being orthonormal and being of the form given on the previous slide.

This problem is NP hard.

A **relaxation** of the discrete problem is to solve

$$\min_{\mathbf{H}} \text{trace } \mathbf{H}' \mathbf{L} \mathbf{H}$$

subject to  $\mathbf{H}$  being orthonormal (but letting its columns take on arbitrary values in  $\mathbb{R}^n$ )

This is a well-known problem.

The solution is that  $\mathbf{H}$  are the **eigenvectors** associated with the  $m$  smallest eigenvalues of  $\mathbf{L}$ .

Let  $\mathbf{q}_i$  be the  $i$ th row of  $\mathbf{H}$ . This is not a classifier.

To obtain one we solve

$$\min_{V_1, \dots, V_m} \sum_{k=1}^m \frac{\sum_{i \in V_k, j \in V_k} \|\mathbf{q}_i - \mathbf{q}_j\|^2}{|V_k|} = \min_{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m} \sum_{i=1}^n \min_g \|\mathbf{q}_i - \boldsymbol{\mu}_g\|^2$$

which minimizes the within cluster sum of squares. This is the  **$k$ -means** algorithm.

This step is ad hoc without further justification (!)

## Normalized clustering

An alternative to RatioCut is

$$\text{Ncut}(V_1, \dots, V_m) = \sum_{k=1}^m \frac{\text{cut}(V_k, \bar{V}_k)}{\text{vol}(V_k)},$$

which features a different normalization in the summands.

Recall that we want to look for a partitioning so that (i) edges within a cluster have a high weight and (ii) edges across different clusters have a low weight.

(ii) amounts to making  $\text{cut}(V_k, \bar{V}_k)$  small.

(i) amounts to making

$$\sum_{(i,j) \in V_k} w_{i,j} = \sum_{i \in V_k} \sum_{j=1}^n w_{i,j} - \sum_{i \in V_k} \sum_{j \in \bar{V}_k} w_{i,j} = \text{vol}(V_k) - \text{cut}(V_k, \bar{V}_k)$$

large. This is done by both minimizing the cut and maximizing the volume.

Ncut achieves both (i) and (ii) while Ratiocut only looks at (ii).

Like RatioCut, the Ncut problem has a connection to the Laplacian matrix.

Redefine the columns of  $\mathbf{H}$  as

$$(\mathbf{h}_1)_i = \begin{cases} 1/\sqrt{\text{vol}(V_k)} & \text{if } i \in V_k \\ 0 & \text{if } i \notin V_k \end{cases} .$$

Then

$$\text{Ncut}(V_1, \dots, V_m) = \text{trace } \mathbf{H}' \mathbf{L} \mathbf{H}$$

subject to  $\mathbf{H}' \mathbf{D} \mathbf{H} = \mathbf{I}_m$  and  $\mathbf{H}$  being of the form given above.

Can substitute in  $\dot{\mathbf{H}} = \mathbf{D}^{1/2} \mathbf{H}$  to obtain the alternative representation

$$\text{Ncut}(V_1, \dots, V_m) = \text{trace } \dot{\mathbf{H}}' \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} \dot{\mathbf{H}} = \text{trace } \dot{\mathbf{H}}' \mathbf{L}_{\text{sym}} \dot{\mathbf{H}}$$

subject to  $\dot{\mathbf{H}}' \dot{\mathbf{H}} = \mathbf{I}_m$  and  $\dot{\mathbf{H}}$  being of the form given above.

Hence, relaxing discreteness yields the solution that  $\dot{\mathbf{H}}$  are the eigenvectors of  $\mathbf{L}_{\text{sym}}$  associated with the  $m$  smallest eigenvalues.

But then  $\mathbf{H}$  are the corresponding eigenvectors of  $\mathbf{L}_{\text{rw}}$ .

An alternative view on Ncut comes from the following observation.

Let  $\{X_t\}$  be a stationary Markov process on  $G$ . Then

$$\text{Ncut}(V_k, \bar{V}_k) = \Pr(X_t \in V_k | X_{t-1} \in \bar{V}_k) + \Pr(X_t \in \bar{V}_k | X_{t-1} \in V_k).$$

This is the probability of switching to/from a cluster.

Normalized spectral clustering aims to make the probability of such an event small.

A proof of the above comes from the observation that, for any set  $A \subset V$ ,

$$\Pr(X_t \in A, X_{t-1} \in \bar{A}) = \sum_{i \in \bar{A}} \sum_{j \in A} \Pr(X_t = j, X_{t-1} = i) = \sum_{i \in \bar{A}} \sum_{j \in A} \pi_i \mathbf{P}_{i,j}$$

but

$$\pi_i = \frac{d_i}{\text{vol}(V)}, \quad \mathbf{P}_{i,j} = \frac{w_{i,j}}{d_i},$$

and so

$$\Pr(X_t \in A, X_{t-1} \in \bar{A}) = \frac{1}{\text{vol}(V)} \sum_{i \in \bar{A}} \sum_{j \in A} w_{i,j}.$$

Therefore,

$$\Pr(X_t \in A | X_{t-1} \in \bar{A}) = \frac{\Pr(X_t \in A, X_{t-1} \in \bar{A})}{\Pr(X_{t-1} \in \bar{A})} = \frac{1}{\text{vol}(\bar{A})} \sum_{i \in \bar{A}} \sum_{j \in A} w_{i,j},$$

which equals

$$\frac{\text{cut}(\bar{A}, A)}{\text{vol}(\bar{A})}.$$

Reversing the roles of  $A$  and  $\bar{A}$  in the above and collecting terms yields the result.

# Outline

- 1 Graphs**
- 2 Spectral clustering**
- 3 Stochastic block model**
- 4 Fixed-effect estimation**
- 5 Linear social-interaction models**

## Classic SBM

Consider an unweighted undirected graph  $G$ . (If weighted, can always reduce to binary case.)

Suppose that nodes are of one of  $m$  latent types.

Can think about this as (independent) latent variables  $\alpha_i$  that take values from the set  $\{1, \dots, m\}$ .

Conditional on the types of all nodes, edges are placed independently.

The probability of an edge between  $(i, j)$  depends on  $\alpha_i, \alpha_j$ .

Thus

$$\Pr(\mathbf{A}) = \sum_{(z_1, \dots, z_n) \in \{1, \dots, m\}^n} \Pr(\mathbf{A}, \alpha_1 = z_1, \dots, \alpha_n = z_n)$$

where the summand factors as

$$\prod_{i=1}^n \Pr(\alpha_i = z_i) \prod_{i < j} \Pr((\mathbf{A})_{i,j} | \alpha_i = z_i, \alpha_j = z_j).$$

This is a simple but flexible model that allows for (dis-)assortative matching between nodes.

Parametrization of heterogeneity is low dimensional:

$m \times m$  symmetric matrix  $\mathbf{B}$  of bernoulli success probabilities for pairs of types:

$$(\mathbf{B})_{z_1, z_2} = (\mathbf{B})_{z_2, z_1} = \Pr(w_{i,j} = 1 | \alpha_i = z_1, \alpha_j = z_2)$$

and an  $m$ -vector of type probabilities.

In this model there are  $m$  latent communities. Spectral clustering correctly classifies nodes to communities as  $n \rightarrow \infty$ .

## Spectral clustering in SBM (fixed effects)

The intuition behind the consistency of spectral clustering has two steps:

- (i) The eigenvectors of the (normalized) Laplacian converge in probability to those of the expected Laplacian.
- (ii) The expected Laplacian has only  $m$  eigenvectors, which identify group membership.

Point (i) is a technical issue. We focus on Point (ii).

We do this conditional on the variables  $\alpha_1, \dots, \alpha_n$ ; so we treat these as fixed.

Note that the  $n \times n$  matrix  $\mathbf{A}$  has (conditional) expectation

$$\mathbf{A}_0 = \mathbb{E}(\mathbf{A} | \alpha_1, \dots, \alpha_n) = \mathbf{Z}\mathbf{B}\mathbf{Z}'$$

for  $n \times m$  matrix  $\mathbf{Z}$  whose  $i$ th column is composed of all zeros except at location  $\alpha_i$ .

Will assume that all columns of  $\mathbf{Z}$  contain at least one 1 (so all communities have at least one member, and so the matrix has full column rank) and that  $\mathbf{B}$  has full rank.

The normalized Laplacian associated with  $\mathbf{A}_0$  is

$$\mathbf{L}_0 = \mathbf{I}_n - \mathbf{D}_0^{-1/2} \mathbf{A}_0 \mathbf{D}_0^{-1/2}$$

for  $\mathbf{D}_0$  the diagonal degree matrix.

$\mathbf{L}_0$  and  $\bar{\mathbf{L}}_0 = \mathbf{I}_n - \mathbf{L}_0$  have the same eigenvectors so focus on the latter here.

Note that  $\mathbf{D}_0 = \text{diag}(\mathbf{ZBZ}'\mathbf{1}_n)$  and, therefore, that

$$(\mathbf{D}_0)_{i,i} = \sum_{z=1}^m \mathbf{B}_{Z_i,z} n_z$$

for  $n_z = \sum_{i=1}^n \{\alpha_i = z\}$ . Hence, the entries of the  $n \times n$  diagonal matrix  $\mathbf{D}_0$  can only take on  $m$  different values. Collect these values in the  $m \times m$  diagonal matrix

$$(\bar{\mathbf{D}}_0)_{z,z} = \sum_{z'=1}^m \mathbf{B}_{z,z'} n_{z'}$$

Then

$$\bar{\mathbf{L}}_0 = \mathbf{Z}\bar{\mathbf{D}}_0^{-1/2} \mathbf{B}\bar{\mathbf{D}}_0^{-1/2} \mathbf{Z}' = \mathbf{Z}\bar{\mathbf{B}}\mathbf{Z}'$$

for  $\bar{\mathbf{B}} = \bar{\mathbf{D}}_0^{-1/2} \mathbf{B}\bar{\mathbf{D}}_0^{-1/2}$ .

Next, consider the (small)  $m \times m$  matrix

$$\mathbf{C}\bar{\mathbf{B}}\mathbf{C}$$

for  $\mathbf{C} = (\mathbf{Z}'\mathbf{Z})^{1/2}$ . Note that  $\mathbf{C}\bar{\mathbf{B}}\mathbf{C}$  is a full rank and symmetric. So we have the eigendecomposition

$$\mathbf{C}\bar{\mathbf{B}}\mathbf{C} = \mathbf{V}\mathbf{S}\mathbf{V}'$$

where all  $m$  eigenvalues are non-zero.

With  $\mathbf{C}$  invertible, we have  $\bar{\mathbf{B}} = \mathbf{C}^{-1}\mathbf{V}\mathbf{S}\mathbf{V}'\mathbf{C}^{-1}$  and, therefore, also that

$$\bar{\mathbf{L}}_0 = \mathbf{Z}\bar{\mathbf{B}}\mathbf{Z}' = \mathbf{Z}(\mathbf{C}^{-1}\mathbf{V}\mathbf{S}\mathbf{V}'\mathbf{C}^{-1})\mathbf{Z}' = \mathbf{Z}\mathbf{Q}\mathbf{S}\mathbf{Q}'\mathbf{Z}'$$

for  $\mathbf{Q} = \mathbf{C}^{-1}\mathbf{V}$ .

Note that  $(\mathbf{Z}\mathbf{Q})'(\mathbf{Z}\mathbf{Q}) = \mathbf{Q}'\mathbf{Z}'\mathbf{Z}\mathbf{Q} = \mathbf{V}'(\mathbf{Z}'\mathbf{Z})^{-1/2}(\mathbf{Z}'\mathbf{Z})(\mathbf{Z}'\mathbf{Z})^{-1/2}\mathbf{V} = \mathbf{V}'\mathbf{V} = \mathbf{I}_m$ . So the columns of the  $n \times m$  orthonormal matrix  $\mathbf{Z}\mathbf{Q}$  are the eigenvectors of the matrix  $\bar{\mathbf{L}}_0$  (and so equally of  $\mathbf{L}_0$ ).

Finally, letting  $\mathbf{z}_i$  the  $i$ th row of  $\mathbf{Z}$ , we have that  $\mathbf{z}_i\mathbf{Q} = \mathbf{z}_j\mathbf{Q}$  if and only if

$$\mathbf{z}_i = \mathbf{z}_j \quad (\text{and so } \alpha_i = \alpha_j)$$

because  $\mathbf{Q}$  is invertible. Hence, clustering on the eigenvectors of  $\mathbf{L}_0$  yields correct assignment of nodes to communities.

Consistency of spectral clustering has been established in more general stochastic block models.

This includes the setting where  $m$  grows slowly with  $n$ , and also several forms of **sparsity**.

There is not a single definition of a sparse graph.

Roughly, sparse means that the number of edges does not grow proportional to  $n^2$  (which would give rise to a **dense** network). This is often formulated in the (expected) degree  $d_i$  growing only slowly with  $n$ .

In the stochastic block model, usually done by including a (common) sparsity parameter that shrinks the probability of edge formation to zero as  $n$  grows.

A weighted graph can always be translated into an unweighted graph so the above equally applies to the weighted case.

In principle, with a consistent classification in hand, we can treat cluster membership as known and then estimate the model parameters by standard methods.

## Identification in small networks (random effects)

We can see the stochastic block model as a (nonlinear) model with two-way heterogeneity:

$$(\mathbf{A})_{i,j} = w_{i,j} = \varphi(\alpha_i, \alpha_j, \varepsilon_{i,j}),$$

where the  $\alpha_i$  are i.i.d. across  $i$  and the  $\varepsilon_{i,j}$  are i.i.d. across dyads, independent of the  $\alpha_i$ .

Similar to a finite-mixture model in panel data.

Edge weights are dependent through their joint dependence on the  $\alpha_i$  only.

By using such conditional-independence we can nonparametrically identify

The distribution of  $w_{i,j} | \alpha_i, \alpha_j$  and,

the distribution of  $\alpha_i$ ,

from the joint distribution of edge weights in small graphs.

Easiest to see under the assumption that the functions

$$G_z(w) = \mathbb{P}(w_{i,j} \leq w | \alpha_i = z) = \sum_{z'=1}^m \mathbb{P}(w_{i,j} \leq w | \alpha_i = z, \alpha_j = z') p_{z'}$$

are linearly independent.

In that case,  $n \geq 4$  suffices for identification.

Consider a grid of  $m$  points  $w_1, \dots, w_m$  and construct the  $m \times m$  matrix

$$\mathbf{G} = \begin{pmatrix} G_1(w_1) & \dots & G_m(w_1) \\ \vdots & \dots & \vdots \\ G_1(w_m) & \dots & G_m(w_m) \end{pmatrix}.$$

By linear independence we can always find points such that this matrix has full rank.

Next, for distinct integers  $i, j, k, \ell \in \{1, \dots, n\}^4$ , let

$$(\mathbf{M}_-)_{m_1, m_2} = \mathbb{P}(w_{\textcolor{red}{i},j} \leq w_{m_1}, w_{\textcolor{red}{i},k} \leq w_{m_2}),$$

$$(\mathbf{M}_w)_{m_1, m_2} = \mathbb{P}(w_{\textcolor{red}{i},j} \leq w_{m_1}, w_{\textcolor{red}{i},k} \leq w_{m_2}, w_{\textcolor{red}{i},\ell} \leq w),$$

for  $(m_1, m_2) \in \{1, \dots, m\}^2$  and any  $w$ .

Then conditional independence implies that

$$\mathbf{M} = \mathbf{G} \mathbf{D} \mathbf{G}'$$

and

$$\mathbf{M}_w = \mathbf{G} \mathbf{D}^{1/2} \mathbf{D}_w \mathbf{D}^{1/2} \mathbf{G}'$$

for  $\mathbf{D} := \text{diag}(p_1, \dots, p_m)$  and  $\mathbf{D}_w = \text{diag}(G_1(w), \dots, G_m(w))$  for any  $w$ .

Use an eigendecomposition of  $\mathbf{M}$  to construct a matrix  $\mathbf{V}$  so that

$$\mathbf{V}\mathbf{M}\mathbf{V}^\top = \mathbf{I}_m.$$

Note that

$$\mathbf{V}\mathbf{M}\mathbf{V}^\top = \mathbf{V}(\mathbf{G}\mathbf{D}\mathbf{G}')\mathbf{V}^\top = (\mathbf{V}\mathbf{G}\mathbf{D}^{1/2})(\mathbf{D}^{1/2}\mathbf{G}'\mathbf{V}^\top) = \mathbf{Q}\mathbf{Q}' = \mathbf{I}_m$$

for  $\mathbf{Q} := \mathbf{V}\mathbf{G}\mathbf{D}^{1/2}$  an  $m \times m$  orthonormal matrix of full rank.

Next,

$$\mathbf{V}\mathbf{M}_w\mathbf{V}^\top = (\mathbf{V}\mathbf{G}\mathbf{D}^{1/2})\mathbf{D}_w(\mathbf{D}^{1/2}\mathbf{G}'\mathbf{V}^\top) = \mathbf{Q}\mathbf{D}_w\mathbf{Q}'$$

for any  $w$ .

By linear independence of the columns of  $\mathbf{G}$  we can recover the matrix  $\mathbf{Q}$  up to sign and permutation of its columns as the joint diagonalizer of the collection of matrices  $\mathbf{V}\mathbf{M}_{w_1}\mathbf{V}^\top, \dots, \mathbf{V}\mathbf{M}_{w_m}\mathbf{V}^\top$ .

Let this matrix be denoted as  $\bar{\mathbf{Q}} = \mathbf{Q}\Delta\Sigma$ .

Here,  $\Sigma$  is a permutation matrix and  $\Delta$  is a diagonal matrix with only 1 or  $-1$  as diagonal entries.

Next, construct the matrix

$$(\mathbf{N}_w)_{m_1, m_2} = \mathbb{P}(w_{\textcolor{red}{i}, k} \leq w_{m_1}, w_{\textcolor{red}{j}, \ell} \leq w_{m_2}, w_{\textcolor{red}{i}, \textcolor{red}{j}} \leq w).$$

Then

$$\mathbf{N}_w = \mathbf{G} \mathbf{D} \mathbf{H}_w \mathbf{D} \mathbf{G}'$$

for

$$(\mathbf{H}_w)_{z, z'} = \mathbb{P}(w_{i, j} \leq w | \alpha_i = z, \alpha_j = z').$$

But

$$\mathbf{V} \mathbf{N}_w \mathbf{V}^\top = (\mathbf{V} \mathbf{G} \mathbf{D}^{1/2}) \mathbf{D}^{1/2} \mathbf{H}_w \mathbf{D}^{1/2} (\mathbf{D}^{1/2} \mathbf{G}' \mathbf{V}^\top) = \mathbf{Q} \mathbf{D}^{1/2} \mathbf{H}_w \mathbf{D}^{1/2} \mathbf{Q}'$$

and so the entries of

$$\bar{\mathbf{Q}}' \mathbf{V} \mathbf{N}_w \mathbf{V}^\top \bar{\mathbf{Q}} = \boldsymbol{\Sigma}' \boldsymbol{\Delta} (\mathbf{D}^{1/2} \mathbf{H}_w \mathbf{D}^{1/2}) \boldsymbol{\Delta} \boldsymbol{\Sigma}$$

give the

$$\mathbb{P}(w_{i, j} \leq w | \alpha_i = z, \alpha_j = z') \sqrt{p_z p_{z'}}$$

up to sign and labelling. Note that for  $w = +\infty$  this gives  $\sqrt{p_z p_{z'}}$  up to the same sign and labelling, so that the ratio yields  $\mathbb{P}(w_{i, j} \leq w | \alpha_i = z, \alpha_j = z')$  while the diagonals yield the  $p_z$ .

# Outline

- 1 Graphs**
- 2 Spectral clustering**
- 3 Stochastic block model**
- 4 Fixed-effect estimation**
- 5 Linear social-interaction models**

## Models on networks

If, in an (undirected) network of  $n$  agents we observe interactions between all

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

possible pairs of nodes we can consider estimation of nonlinear fixed-effect models.

Say

$$y_{i,j} | x_{i,j}, \alpha_i, \alpha_j$$

has a parametric distribution. Can do bias-corrected maximum likelihood estimation.

This is like a panel data problem where  $N, T \rightarrow \infty$  at the same rate; for each of  $n$  nodes we observe  $(n - 1)$  outcomes, so we get more measurements on each individual as  $n \rightarrow \infty$ .

(Fixed effect estimation from many small networks will not be possible.)

The case where we observe more **limited interaction** between nodes is much more complicated.

## A linear model

Consider an undirected (multi-) graph where each edge  $(i, j)$  has associated with it an outcome  $y_{i,j}$ , with

$$\mathbb{E}(y_{i,j} | \alpha_i, \alpha_j, w_{i,j} > 0) = \alpha_i - \alpha_j$$

Say we observe  $m$  edges.

Stack all outcomes in  $m$ -vector  $\mathbf{y}$ .

Let  $\mathbf{B}$  be the  $m \times n$  matrix where the row for edge  $e = (i, j)$  has  $(\mathbf{B})_{e,i} = 1$  and  $(\mathbf{B})_{e,j} = -1$ . This is the (oriented) **incidence matrix** of the graph; we note that

$$\mathbf{L} = \mathbf{B}' \mathbf{B}$$

connects it to the Laplacian matrix.

We can write

$$\mathbf{y} = \mathbf{B}\boldsymbol{\alpha} + \boldsymbol{\varepsilon}, \quad \mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{B}) = \mathbf{0}.$$

The conditional-mean restriction is one of **network exogeneity**.

## A bipartite special case

Suppose that the node set  $V$  can be partitioned into two sets  $V_1, V_2$  and that edges are only formed between the two subsets.

Can then reparametrize

$$\alpha_i = \begin{cases} \eta_i & \text{if } i \in V_1 \\ -\gamma_i & \text{if } i \in V_2 \end{cases}$$

An example is worker  $i$  works in firm  $j$  at time  $t$ ; so that log-wage  $y_{i,j,t}$  decomposes as

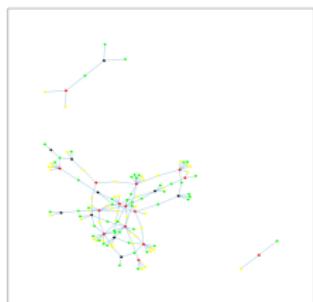
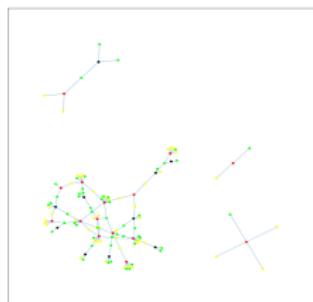
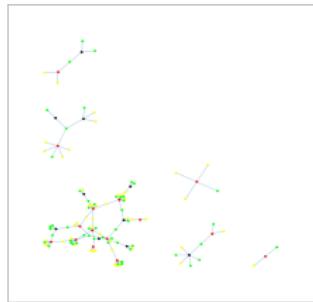
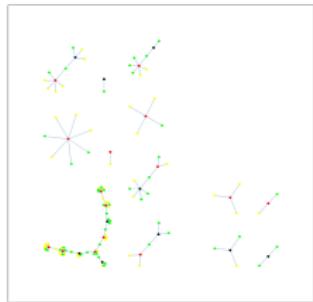
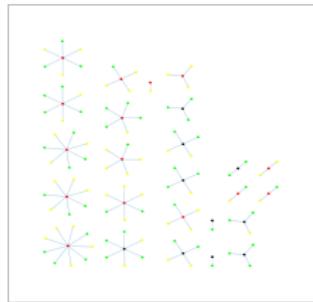
$$y_{i,j,t} = \eta_i + \gamma_j(i,t) + \varepsilon_{i,j,t}.$$

Here, the number of edges between any  $i$  and  $j$  is the number of periods that  $i$  is employed by  $j$ .

The connectivity structure of the induced graph drives the statistical properties of the fixed-effect estimator.

For example, if no worker ever switches employer we cannot learn anything.

## (Dis-)connected graphs



The above plot is for a stationary model where workers and firms are of binary types (high and low quality), and mobility of workers is introduced through exogenous lay-offs.

A one-period network has as many components as firms.

Cannot disentangle worker effects from firm effects.

Adding time periods makes workers switch to different firms.

Makes the graph more connected.

With enough mobility, the graph becomes connected.

Connectivity allows for the least-squares estimator to be well defined (subject to one normalization on the fixed effects) but is not enough for it to have good properties.

## Least-squares estimator

Each row of  $\mathbf{B}$  sums to zero so we normalize the fixed effects so that

$$\sum_{i=1}^n \sum_{j=1}^n (\mathbf{A})_{i,j} (\alpha_i + \alpha_j) = 0.$$

If the graph is connected then the (constrained) least-squares estimator exists and is equal to

$$\hat{\boldsymbol{\alpha}} = \mathbf{D}^{-1/2} (\mathbf{L}_{\text{sym}})^* \mathbf{D}^{-1/2} \mathbf{B}' \mathbf{y}.$$

By strict exogeneity this estimator is unbiased.

Further note that

$$\text{var}(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha} | \mathbf{B}) = \mathbf{D}^{-1/2} (\mathbf{L}_{\text{sym}})^* \mathbf{D}^{-1/2} \mathbf{B}' \mathbb{E}(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' | \mathbf{B}) \mathbf{B} \mathbf{D}^{-1/2} (\mathbf{L}_{\text{sym}})^* \mathbf{D}^{-1/2}.$$

Notably, with  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I}_m)$ ,

$$\text{var}(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha} | \mathbf{B}) = \sigma^2 \mathbf{D}^{-1/2} (\mathbf{L}_{\text{sym}})^* \mathbf{D}^{-1/2}.$$

So the precision of the estimator is driven by the Laplacian of the graph.

We have

$$\text{var}(\hat{\alpha}_i | \mathbf{B}) = \sigma^2 \frac{(\mathbf{L}_{\text{sym}}^*)_{i,i}}{d_i}$$

The behavior of this variance depends in a very complicated way on the connectivity structure of the network.

No reason it would behave proportional to  $d_i^{-1}$ , which would correspond to the usual parametric rate.

Let  $\lambda_2$  be second smallest eigenvalue of  $\mathbf{L}_{\text{sym}}$  and

$$h_i = \left( \frac{1}{d_i} \sum_{j=1}^n \frac{w_{i,j}^2}{d_j} \right)^{-1}.$$

Can derive that

$$\frac{\sigma^2}{d_i} - \frac{2\sigma^2}{m} \leq \text{var}(\hat{\alpha}_i | \mathbf{B}) \leq \frac{\sigma^2}{d_i} \left( 1 + \frac{1}{\lambda_2 h_i} \right) - \frac{2\sigma^2}{m}.$$

So,

$$\text{var}(\hat{\alpha}_i) = \frac{\sigma^2}{d_i} + o(d_i^{-1})$$

provided that  $\lambda_2 h_i \rightarrow \infty$ .

In this case can get conventional asymptotic normality with  $d_i$  serving as effective sample size for estimating  $\alpha_i$ .

## Example: student/teacher data

Pupils in Grades 4 and 5 of elementary school over the period 2008–2012.

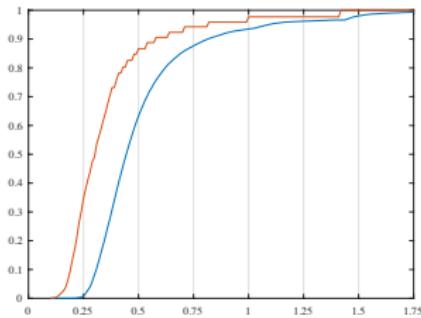
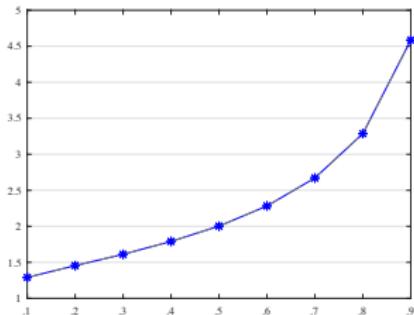
(One-mode projection has)  $m = 41,612$  edges and  $n = 11,945$  teachers.

$\lambda_2 = .0039$ .

	mean	stdev	10 <sup>th</sup> %	20 <sup>th</sup> %	30 <sup>th</sup> %	40 <sup>th</sup> %	50 <sup>th</sup> %	60 <sup>th</sup> %	70 <sup>th</sup> %	80 <sup>th</sup> %	90 <sup>th</sup> %
$d_i$	13.87	10.76	3.00	5.50	7.50	9.00	11.00	14.00	17.50	21.50	27.50
$h_i$	7.15	7.13	2.43	3.30	4.01	4.72	5.48	6.36	7.44	9.12	12.56

Conventional inference procedures not applicable.

Under homokedasticity we can compute the ratio of the exact variance versus its first-order approximation.



Left plot: deciles of the distribution of this ratio.

Right plot: distribution of width of 95% confidence intervals constructed using exact (blue) and approximate (red) standard error (using  $\sigma^2 = 1$ ).

## Variance components

Interest often lies in variances and covariances between fixed effects.

Of the form

$$\boldsymbol{\alpha}' \mathbf{W} \boldsymbol{\alpha}$$

for some weight matrix  $\mathbf{W}$ .

Plug-in estimator is

$$\hat{\boldsymbol{\alpha}}' \mathbf{W} \hat{\boldsymbol{\alpha}}$$

and

$$\mathbb{E}(\hat{\boldsymbol{\alpha}}' \mathbf{W} \hat{\boldsymbol{\alpha}} | \mathbf{B}, \mathbf{W}) = \boldsymbol{\alpha}' \mathbf{W} \boldsymbol{\alpha} + \sigma^2 \text{trace}(\mathbf{W} \text{var}(\hat{\boldsymbol{\alpha}} | \mathbf{B})).$$

The bias again depends on graph structure. Here, can be corrected for exactly by subtracting an unbiased estimator of the bias.

Nevertheless, only in the dense-network case is the bias-corrected estimator asymptotically normal!.

For example in a graph that contains a fixed number of clusters, the estimator follows a mixture of non-central chi-squared distributions. Inference here is highly **non-standard**. (student-teacher/worker-firm data does not fit this setup!)

# Outline

- 1 Graphs**
- 2 Spectral clustering**
- 3 Stochastic block model**
- 4 Fixed-effect estimation**
- 5 Linear social-interaction models**

## Peer effects

Observe data  $(y_i, x_i)$  for agents  $i = 1, \dots, n$ .

Agents engage in social interactions.

Action and/or outcome of agent  $i$  may be influenced by action/outcome or a characteristic of agent  $j$ .

Corresponds to a graph  $G$  with adjacency matrix  $\mathbf{A}$  and normalized  $\mathbf{P}$ .

Gives rise to **peer effects**:

**Correlated effects:**

Agents behave similarly because they are in the same environment.

**Exogenous effects, contextual effects, or spillover effects:**

The propensity of an agent to behave in a certain way varies with the distribution of characteristics in his/her reference group.

**Endogenous effects:**

The propensity of an agent to behave in a certain way varies with the prevalence of that behavior in the reference group.

## A linear model

Stack observations in  $n$ -vectors  $\mathbf{y} = (y_1, \dots, y_n)'$  and  $\mathbf{x} = (x_1, \dots, x_n)'$ .

A linear model for social interactions is

$$\mathbf{y} = \alpha \mathbf{1}_n + \rho \mathbf{P} \mathbf{y} + \beta \mathbf{x} + \delta \mathbf{P} \mathbf{x} + \boldsymbol{\varepsilon}$$

with

$$\mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{A}, \mathbf{x}) = \mathbf{0}.$$

Assume that  $|\rho| < 1$  so that the model has a reduced form:

$$\begin{aligned}\mathbf{y} &= (\mathbf{I}_n - \rho \mathbf{P})^{-1} (\alpha \mathbf{1}_n + \beta \mathbf{x} + \delta \mathbf{P} \mathbf{x} + \boldsymbol{\varepsilon}) \\ &= \left( \sum_{k=0}^{\infty} \rho^k \mathbf{P}^k \right) (\alpha \mathbf{1}_n + \beta \mathbf{x} + \delta \mathbf{P} \mathbf{x} + \boldsymbol{\varepsilon}) \\ &= \frac{\alpha}{1-\rho} \mathbf{1}_n + \beta \mathbf{x} + \lambda \sum_{k=1}^{\infty} \rho^{k-1} \mathbf{P}^k \mathbf{x} + \sum_{k=0}^{\infty} \rho^k \mathbf{P}^k \boldsymbol{\varepsilon}\end{aligned}$$

for  $\lambda = \rho\beta + \delta$ .

Multiplying through by  $\mathbf{P}$  and taking conditional expectations yields

$$\mathbb{E}(\mathbf{P}\mathbf{y}|\mathbf{x}, \mathbf{A}) = \frac{\alpha}{1-\rho} \mathbf{1}_n + \beta \mathbf{P}\mathbf{x} + \lambda(\mathbf{P}^2\mathbf{x} + \rho\mathbf{P}^3\mathbf{x} + \rho^2\mathbf{P}^4\mathbf{x} + \dots).$$

This equation reveals that  $\mathbf{P}^2\mathbf{x}$  is a relevant (and valid) instrumental variable for  $\mathbf{P}\mathbf{y}$  provided that

$$\lambda = \rho\beta + \delta \neq 0 \text{ and}$$

$\mathbf{I}_n, \mathbf{P}, \mathbf{P}^2$  are not colinear.

When  $\rho \neq 0$  we have overidentifying restrictions coming from  $\mathbf{P}^3\mathbf{x}, \mathbf{P}^4\mathbf{x}, \dots$

## Linear-in-means model

Recalling that  $\mathbf{P}$  is a transition matrix on the network,  $\mathbf{P}^2\mathbf{x}$  is an average of characteristics of peers of peers.

No-multicollinearity then requires variation in peer-group composition.

In the complete network,

$$\mathbf{A} = \mathbf{1}_n \mathbf{1}'_n - \mathbf{I}_n, \quad \mathbf{D} = (n-1) \mathbf{I}_n$$

and so

$$\mathbf{P} = \mathbf{D}^{-1} \mathbf{A} = \frac{1}{n-1} (\mathbf{1}_n \mathbf{1}'_n - \mathbf{I}_n).$$

But also

$$\mathbf{P}^2 = \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1} \mathbf{A} = \frac{(n-2)}{(n-1)^2} \mathbf{P} + \frac{1}{(n-1)} \mathbf{I}_n,$$

yielding collinearity.

## Variation in peer-group size

Note that

$$\sum_{j=1}^n (\mathbf{P})_{i,j} x_j = \frac{1}{n-1} \sum_{j \neq i} x_j = \frac{1}{n-1} (n \bar{x} - x_i),$$

and

$$\sum_{j=1}^n (\mathbf{P})_{i,j} y_j = \frac{1}{n-1} \sum_{j \neq i} y_j = \frac{1}{n-1} (n \bar{y} - y_i),$$

So

$$y_i = \alpha + \rho \sum_{j=1}^n (\mathbf{P})_{i,j} y_j + \beta x_i + \delta \sum_{j=1}^n (\mathbf{P})_{i,j} x_j + \frac{1}{n-1} \sum_{j \neq i} x_j + \varepsilon_i$$

is equal to

$$y_i = \alpha + \rho \frac{1}{n-1} (n \bar{y} - y_i) + \beta x_i + \delta \frac{1}{n-1} (n \bar{x} - x_i) + \frac{1}{n-1} \sum_{j \neq i} x_j + \varepsilon_i.$$

This can be re-arranged to get

$$\left(1 + \rho \frac{1}{n-1}\right) y_i = \alpha + \rho \frac{n}{n-1} \bar{y} + \left(\beta + \delta \frac{1}{n-1}\right) x_i + \delta \frac{n}{n-1} \bar{x} + \varepsilon_i.$$

Averaging yields

$$\left(1 + \rho \frac{1}{n-1}\right) \bar{y} = \alpha + \rho \frac{n}{n-1} \bar{y} + \left(\beta + \delta \frac{1}{n-1}\right) \bar{x} + \delta \frac{n}{n-1} \bar{x} + \bar{\varepsilon}.$$

Differencing gives

$$\left(1 + \rho \frac{1}{n-1}\right) (y_i - \bar{y}) = \alpha + \left(\beta + \delta \frac{1}{n-1}\right) (x_i - \bar{x}) + (\varepsilon_i - \bar{\varepsilon}).$$

The coefficients in the (population) regression

$$(y_i - \bar{y}) = a_n + b_n (x_i - \bar{x}) + (\varepsilon_i - \bar{\varepsilon})$$

are identified and equal

$$a_n = \frac{\alpha}{1 + \rho \frac{1}{n-1}}, \quad b_n = \frac{\beta + \delta \frac{1}{n-1}}{1 + \rho \frac{1}{n-1}}.$$

Because  $b_n$  is known and depends on  $n$  and involves three different structural parameters, we can recover  $\beta, \delta, \rho$  from  $b_{n_1}, b_{n_2}, b_{n_3}$  if we observe three (or more) different network sizes  $n_1, n_2, n_3$ .

They are uniquely recoverable provided that  $\lambda \neq 0$ , where  $\lambda$  is the same parameter as before.

## Network effects

The above model does little to accommodate correlated effects.

When at the network level, such effects are network fixed effects.

Again take

$$\mathbf{y} = \alpha \mathbf{1}_n + \rho \mathbf{P}\mathbf{y} + \beta \mathbf{x} + \delta \mathbf{P}\mathbf{x} + \varepsilon$$

but see  $\alpha$  as a fixed effect for the network in question (se we are sampling many small networks here).

Can difference within the network:

$$(\mathbf{I} - \mathbf{P})\mathbf{y} = \rho(\mathbf{I} - \mathbf{P})\mathbf{P}\mathbf{y} + \beta(\mathbf{I} - \mathbf{P})\mathbf{x} + \delta(\mathbf{I} - \mathbf{P})\mathbf{P}\mathbf{x} + (\mathbf{I} - \mathbf{P})\varepsilon.$$

Reduced form then comes from inverting

$$(\mathbf{I} - \rho\mathbf{P})(\mathbf{I} - \mathbf{P})\mathbf{y} = \beta(\mathbf{I} - \mathbf{P})\mathbf{x} + \delta(\mathbf{I} - \mathbf{P})\mathbf{P}\mathbf{x} + (\mathbf{I} - \mathbf{P})\varepsilon.$$

Now need that  $\mathbf{I}_n, \mathbf{P}, \mathbf{P}^2, \mathbf{P}^3$  are linearly independent.