aphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results 00000000000	Conclusion O	References

Statistics Seminars – University of Kent Bayesian estimation of the latent dimension and communities in stochastic blockmodels

Imperial College London

Francesco Sanna Passino, Nick Heard Department of Mathematics, Imperial College London francesco.sanna-passino16@imperial.ac.uk

February 6, 2020

000000	000000000	000000000000000000000000000000000000000	0000000000	0	
<u>^</u>					

Outline

1 Graphs, SBMs and RDPGs

- Statistical models for graphs
- SBMs and RDPGs
- Beyond RDPGs: the GRDPG
 - Introduction to the GRDPG
 - Network embeddings
 - Limit theorems for the GRDPG
 - Spectral estimation of the SBM
- **3** Bayesian modelling of embeddings
 - Joint estimation of d and K

- Model validation
- Curse of dimensionality
- Second-order clustering
- Inference
- Directed and bipartite graphs

4 Results

- Simulated data
- Santander bikes
- Enron e-mail network

5 Conclusion

Graphs, SBMs and RDPGs ●○○○○○	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion ○	References
UNDIRECTED G	RAPHS				

- Undirected graph $\mathbb{G} = (V, E)$ where:
 - V is the **node set**, n = |V|,
 - $E \subseteq V \times V$ is the **edge set**, containing dyads $(i, j), i, j \in V$.
- An edge is drawn if a node $i \in V$ connects to $j \in V$, written $(i, j) \in E$.
- From \mathbb{G} , an adjacency matrix $\mathbf{A} = \{A_{ij}\}$, of dimension $n \times n$, can be obtained:

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 & \cdots & 0 & 1 \\ 1 & 0 & 1 & \cdots & 1 & 0 \\ 0 & 1 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 0 & 1 & \cdots & 1 & 0 \end{pmatrix} \qquad \qquad A_{ij} = \begin{cases} 1 & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}$$

• Commonly, self-edges are not allowed, implying that A is a hollow matrix.

Bayesian estimation of the latent dimension and communities in stochastic blockmodels

Graphs, SBMs and RDPGs ○●○○○○	Beyond RDPGs: the GRDPG	Bayesian modellin	g of embed	dings	l	Result	s 0000	000	0	onclu	sion	References
A TOY EXAMPLE	E											
$\mathbb{G}=$		$\mathbf{A} =$	$\begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$	$egin{array}{c} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array}$	0 1 1 0 0 0 0 0 0 0 0 0	0 1 1 0 0 1 0 1 0 0 0	$ \begin{array}{c} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ $	0 0 0 0 0 0 0 1 0 1	0 0 0 1 0 1 0 1 1	0 0 0 0 0 0 0 1 0 0	$ \begin{array}{c} 0\\0\\0\\0\\0\\1\\1\\0\\0\end{array} \end{array} $	

Bayesian estimation of the latent dimension and communities in stochastic blockmodels

9 10

Graphs, SBMs and RDPGs ⊙●○○○○	Beyond RDPGs: the GRDPG	Bayesian modelling	g of embedo	lings	1	Result	s 0000	000	0	onclu	sion	References
A toy example												
$\mathbb{G}=$		$ ightarrow \mathbf{A} =$	$\begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 0 $	$egin{array}{c} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array}$	0 1 1 0 0 0 0 0 0 0 0 0	0 1 1 0 0 1 0 1 0 0	1 0 0 1 0 0 0 0 0 0	0 0 0 0 0 0 0 1 0 1	0 0 0 1 0 1 0 1 1	0 0 0 0 0 0 0 1 0 0	$\begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\$	

Bayesian estimation of the latent dimension and communities in stochastic blockmodels

10

9

Graphs, SBMs and RDPGs ○●○○○○	Beyond RDPGs: the GRDPG	Bayesian modellin	g of embedo	dings	l	Result	s 0000	000	0	onclu	sion	References
A TOY EXAMPLE	E											
$\mathbb{G}=% \left($		$\mathbf{A} =$	$\begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$	$egin{array}{c} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array}$	0 1 1 0 0 0 0 0 0 0 0 0 0	0 1 1 0 0 1 0 1 0 0	$ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ $	0 0 0 0 0 0 0 1 0 1	0 0 0 1 0 1 0 1 1	0 0 0 0 0 0 0 1 0 0	$ \begin{array}{c} 0\\0\\0\\0\\0\\0\\1\\1\\0\\0\end{array} \end{array} $	

Bayesian estimation of the latent dimension and communities in stochastic blockmodels

9 10



- Consider an undirected graph with symmetric adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$.
- Latent feature models (Hoff, Raftery, and Handcock, 2002): each node is assigned a latent position x_i in a *d*-dimensional latent space X.
- The edges are generated *independently* using a **kernel function** $\psi : \mathbb{X} \times \mathbb{X} \rightarrow [0, 1]$:

$$\mathbb{P}(A_{ij} = 1) = \psi(\boldsymbol{x}_i, \boldsymbol{x}_j), \ i < j, \ A_{ij} = A_{ji}.$$

- The latent positions are represented as a (n imes d)-dimensional matrix $\mathbf{X} = [m{x}_1, \dots, m{x}_n]^ op$.
- In random dot product graphs (RDPG) (Young and Scheinerman, 2007; Athreya et al., 2018), the kernel is the inner product of the latent positions, and \mathbb{X} is chosen such that $0 \le \mathbf{x}^{\top} \mathbf{y} \le 1 \forall \mathbf{x}, \mathbf{y} \in \mathbb{X}$:

$$\mathbb{P}(A_{ij} = 1) = \boldsymbol{x}_i^\top \boldsymbol{x}_j, \ i < j, \ A_{ij} = A_{ji}.$$

• In RDPGs: $d = \operatorname{rank}\{\mathbb{E}(\mathbf{A})\} = \operatorname{rank}(\mathbf{X}\mathbf{X}^{\top}).$

A TOY EXAMPLE: COMMUNITY DETECTION



Graphs, SBMs and RDPGs ○○○○●○	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
CLUSTERING N	ODES IN UNDIREC	TED GRAPHS			

- The **stochastic blockmodel** (SBM) (Holland, Laskey, and Leinhardt, 1983) is the classical model for community detection in graphs.
- Assume K communities, and a matrix $\mathbf{B} \in [0, 1]^{K \times K}$ of within-community probabilities.
- Each node is assigned a community $z_i \in \{1, \ldots, K\}$ with probability $\boldsymbol{\theta} = (\theta_1, \ldots, \theta_K)$, from the K 1 probability simplex.
- The probability of a link depends on the **community allocations** z_i and z_j of the nodes:

$$\mathbb{P}(A_{ij} = 1) = B_{z_i z_j}, \ i < j, \ A_{ij} = A_{ji}.$$

• The likelihood for an observed symmetric adjacency matrix A is:

$$L(\mathbf{A}|\boldsymbol{z}, \mathbf{B}) = \prod_{1 \le i < j \le n} B_{z_i z_j}^{A_{ij}} (1 - B_{z_i z_j})^{1 - A_{ij}}.$$

THE SBM AS A SPECIAL CASE OF RDPG

- The stochastic blockmodel can be interpreted as a special case of a RDPG.
- For simplicity, initially assume that **B** is *positive semi-definite*.
- Assume that $B_{kh} = \boldsymbol{\mu}_k^\top \boldsymbol{\mu}_h$ for some $\boldsymbol{\mu}_k, \boldsymbol{\mu}_h \in \mathbb{X}$.
- If all the nodes in community k are assigned the latent position μ_k , then:

$$\mathbb{P}(A_{ij} = 1) = B_{z_i z_j} = \boldsymbol{\mu}_{z_i}^\top \boldsymbol{\mu}_{z_j}, \ i < j, \ A_{ij} = A_{ji}.$$

- In this framework: $d = \operatorname{rank}\{\mathbb{E}(\mathbf{A})\} = \operatorname{rank}(\mathbf{X}\mathbf{X}^{\top}) = \operatorname{rank}(\mathbf{B}) \leq K$.
- Inference on SBMs as RDPGs:
 - Latent dimension *d*,
 - Number of communities *K*,
 - Community allocations $\boldsymbol{z} = (z_1, \dots, z_n)$,
 - Latent positions μ_1, \ldots, μ_K .

Conclusion

BEYOND RDPGs: THE GENERALISED RANDOM DOT PRODUCT GRAPH

Definition (Generalised random dot product graph, GRDPG, Rubin-Delanchy et al., 2017)

Let d_+, d_- be non-negative integers such that $d = d_+ + d_-$. Let $\mathbb{X} \subseteq \mathbb{R}^d$ such that $\forall \boldsymbol{x}, \boldsymbol{x}' \in \mathbb{X}, 0 \leq \boldsymbol{x}^\top \mathbf{I}(d_+, d_-) \boldsymbol{x}' \leq 1$, where

$$\mathbf{I}(p,q) = \operatorname{diag}(\underbrace{1,\ldots,1}_{p},\underbrace{-1,\ldots,-1}_{q}).$$

Let \mathcal{F} be a probability measure on \mathbb{X} , $\mathbf{A} \in \{0,1\}^{n \times n}$ be a symmetric matrix and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{X}^n$. Then $(\mathbf{A}, \mathbf{X}) \sim \text{GRDPG}_{d_+, d_-}(\mathcal{F})$ if $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{iid}{\sim} \mathcal{F}$ and for i < j, independently

$$\mathbb{P}(A_{ij}=1) = \boldsymbol{x}_i^\top \mathbf{I}(d_+, d_-) \boldsymbol{x}_j.$$

To represent the K-community SBM as a GRDPG, 𝔅 can be chosen to have mass concentrated at μ₁,..., μ_K ∈ ℝ^d such that μ_i[⊤]I(d₊, d₋)μ_j = B_{ij} ∀ i, j ∈ {1,...,K}.

IDENTIFIABILITY OF THE GRDPG

- The GRDPG has two sources of non-identifiability (Cape, Tang, and Priebe, 2018).
- Uniqueness up to indefinite orthogonal transformations For any matrix $\mathbf{Q} \in \mathbb{O}(d_+, d_-)$, the indefinite orthogonal group with signature (d_+, d_-) ,

$$(\mathbf{Q}\boldsymbol{\mu}_{z_i})^{\top}\mathbf{I}(d_+, d_-)(\mathbf{Q}\boldsymbol{\mu}_{z_j}) = \boldsymbol{\mu}_{z_i}^{\top}\mathbf{I}(d_+, d_-)\boldsymbol{\mu}_{z_j},$$

which implies that the likelihood is invariant to any such transformation.

Uniqueness up to artificial dimension blow-up For $(\mathbf{A}, \mathbf{X}) \sim \text{GRDPG}_{d_{\perp}, d_{-}}(\mathcal{F})$, there exists \mathcal{F}^{\star} on $\mathbb{R}^{d^{\star}}$, with $d^{\star} > d$, such that

$$(\mathbf{A}, \mathbf{X}) \stackrel{d}{=} (\mathbf{A}^{\star}, \mathbf{X}^{\star}) \text{ with } (\mathbf{A}^{\star}, \mathbf{X}^{\star}) \sim \mathrm{GRDPG}_{d_{+}^{\star}, d_{-}^{\star}}(\mathcal{F}^{\star}).$$

In the SBM setting, this essentially means that **any** matrix $\mathbf{B} \in [0, 1]^{K \times K}$ with rank d can be obtained as an inner product between latent positions on **arbitrarily large** dimensions.

Francesco Sanna Passino

Definition (Adjacency spectral embedding, ASE)

For $d \in \{1, \ldots, n\}$, consider the spectral decomposition

 $\mathbf{A} = \hat{\boldsymbol{\Gamma}} \hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Gamma}}^\top + \hat{\boldsymbol{\Gamma}}_\perp \hat{\boldsymbol{\Lambda}}_\perp \hat{\boldsymbol{\Gamma}}_\perp^\top,$

where $\hat{\mathbf{A}}$ is a $d \times d$ diagonal matrix containing the top d eigenvalues in magnitude, in decreasing order, $\hat{\mathbf{\Gamma}}$ is a $n \times d$ matrix containing the corresponding orthonormal eigenvectors, and the matrices $\hat{\mathbf{A}}_{\perp}$ and $\hat{\mathbf{\Gamma}}_{\perp}$ contain the remaining n - d eigenvalues and eigenvectors. The adjacency spectral embedding $\hat{\mathbf{X}} = [\hat{x}_1, \dots, \hat{x}_n]^{\top}$ of \mathbf{A} in \mathbb{R}^d is

$$\hat{\mathbf{X}} = \hat{\mathbf{\Gamma}} |\hat{\mathbf{\Lambda}}|^{1/2} \in \mathbb{R}^{n \times d},$$

where the operator $|\cdot|$ applied to a matrix returns the absolute value of its entries.

• $\hat{\mathbf{X}}\mathbf{I}(d_+, d_-)\hat{\mathbf{X}}^{\top}$ represents an estimate of $\mathbb{E}(\mathbf{A}) = \mathbf{X}\mathbf{I}(d_+, d_-)\mathbf{X}^{\top} \rightarrow \mathbf{link}$ prediction.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
Network embi	EDDINGS				

Definition (Laplacian spectral embedding, LSE)

For $d \in \{1, \dots, n\}$, consider the (modified) normalised Laplacian matrix

$$\mathbf{L} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}, \ \mathbf{D} = \operatorname{diag}\left(\sum_{j=1}^{n} A_{ij}\right),$$

and its spectral decomposition

$$\mathbf{L} = \tilde{\boldsymbol{\Gamma}} \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Gamma}}^\top + \tilde{\boldsymbol{\Gamma}}_\perp \tilde{\boldsymbol{\Lambda}}_\perp \tilde{\boldsymbol{\Gamma}}_\perp^\top.$$

The Laplacian spectral embedding $ilde{\mathbf{X}} = [ilde{m{x}}_1, \dots, ilde{m{x}}_n]^ op$ of \mathbf{A} in \mathbb{R}^d is

$$ilde{\mathbf{X}} = ilde{\mathbf{\Gamma}} | ilde{\mathbf{\Lambda}}|^{1/2}.$$

• The modified Laplacian $\mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ (Rohe, Chatterjee, and Yu, 2011) is preferred to the version $\mathbf{I}_n - \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ since its eigenvalues lie in (-1, 1), providing a convenient interpretation for disassortative networks (Rubin-Delanchy, Adams, and Heard, 2016).



- Let $\boldsymbol{\xi}$ be a random vector such that $\boldsymbol{\xi} \sim F$, where F is supported on \mathbb{X} and $\boldsymbol{\xi}$ has full rank second order moment matrix $\boldsymbol{\Delta} = \mathbb{E}(\boldsymbol{\xi}\boldsymbol{\xi}^{\top}) \in \mathbb{R}^{d \times d}$, for d fixed, constant and known.
- Introduce a *sparsity factor* ρ_n , requiring $\rho_n = 1$ or $\rho_n \to 0$.
- The latent positions $\boldsymbol{x}_1^{(n)} = \rho_n^{1/2} \boldsymbol{\xi}_1^{(n)}, \dots, \boldsymbol{x}_n^{(n)} = \rho_n^{1/2} \boldsymbol{\xi}_n^{(n)}$ at each step are assumed to be independent replicates of the random vector $\rho_n^{1/2} \boldsymbol{\xi}$.
- Consequently, \mathcal{F} is assumed to factorise into a product F_{ρ}^{n} of n identical marginal distributions that are equal to F up to scaling.

Theorem (ASE two-to-infinity norm bound)

Consider $(\mathbf{A}^{(n)}, \mathbf{X}^{(n)}) \sim \operatorname{GRDPG}_{d_+, d_-}(F_{\rho}^n)$. There exists a universal constant $\varepsilon > 0$ such that, provided that $n\rho_n = \omega\{(\log n)^{4\varepsilon}\}$, there exists $\mathbf{Q}_n \in \mathbb{O}(d_+, d_-)$ such that

$$\left\|\mathbf{Q}_n \hat{\boldsymbol{x}}_i^{(n)} - \boldsymbol{x}_i^{(n)}\right\|_{2 \to \infty} = \max_i \left\|\mathbf{Q}_n \hat{\boldsymbol{x}}_i^{(n)} - \boldsymbol{x}_i^{(n)}\right\| = O_{\mathbb{P}}\left\{\frac{(\log n)^{\varepsilon}}{n^{1/2}}\right\}.$$

 $X = O_{\mathbb{P}}\{f(n)\} \text{ if, for any constant } \varepsilon > 0, \exists n_{\varepsilon} \in \mathbb{N} \text{ and } C_{\varepsilon} > 0, \text{s.t. } \mathbb{P}\{|X| \leq C_{\varepsilon}f(n)\} \geq 1 - n^{-\varepsilon} \forall n \geq n_{\varepsilon}.$ **13/45**Francesco Sanna Passino
Imperial College London

Conclusion

LIMIT THEOREMS FOR ASE (RUBIN-DELANCHY ET AL., 2017)

Theorem (ASE central limit theorem)

Consider the sequence of graphs $(\mathbf{A}^{(n)}, \mathbf{X}^{(n)}) \sim \operatorname{GRDPG}_{d_+,d_-}(F_{\rho}^n)$, such that $n\rho_n = \omega\{(\log n)^{4\varepsilon}\}$ for the universal constant $\varepsilon > 0$. For any integer m > 0, choose points $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_m \in \mathbb{X}$ in the support of F, and points $\boldsymbol{q}_1, \ldots, \boldsymbol{q}_m \in \mathbb{R}^d$. Then there exists a sequence of random matrices $\mathbf{Q}_n \in \mathbb{O}(d_+, d_-)$ such that

$$\mathbb{P}\left\{\left.\bigcap_{i=1}^{m}n^{1/2}\left(\mathbf{Q}_{n}\hat{\boldsymbol{x}}_{i}^{(n)}-\boldsymbol{x}_{i}^{(n)}\right)\leq\boldsymbol{q}_{i}\right|\boldsymbol{\xi}_{1}^{(n)}=\boldsymbol{x}_{1},\ldots,\boldsymbol{\xi}_{m}^{(n)}=\boldsymbol{x}_{m}\right\}\longrightarrow\prod_{i=1}^{m}\Phi\left\{\boldsymbol{q}_{i},\boldsymbol{\Sigma}(\boldsymbol{x}_{i})\right\},$$

where $\Phi{q, \Sigma}$ is the cumulative distribution function of a multivariate normal distribution with mean 0 and covariance Σ , evaluated at q, and

$$\boldsymbol{\Sigma}(\boldsymbol{x}) = \begin{cases} \mathbf{I}(d_+, d_-) \boldsymbol{\Delta}^{-1} \mathbb{E}[\{\boldsymbol{x}^\top \mathbf{I}(d_+, d_-)\boldsymbol{\xi}\}\{1 - \boldsymbol{x}^\top \mathbf{I}(d_+, d_-)\boldsymbol{\xi}\}\boldsymbol{\xi}\boldsymbol{\xi}^\top] \boldsymbol{\Delta}^{-1} \mathbf{I}(d_+, d_-) & \text{if } \rho_n = 1\\ \mathbf{I}(d_+, d_-) \boldsymbol{\Delta}^{-1} \mathbb{E}[\{\boldsymbol{x}^\top \mathbf{I}(d_+, d_-)\boldsymbol{\xi}\}\boldsymbol{\xi}\boldsymbol{\xi}^\top] \boldsymbol{\Delta}^{-1} \mathbf{I}(d_+, d_-) & \text{if } \rho_n \to 0 \end{cases}$$

• The theorem has *crucial* relevance in practice.

PRACTICAL UTILITY OF THE LIMIT THEOREMS

- If d is known, conditioning on K, the ASE CLT implies that Gaussian mixture modelling gives a consistent estimate of the locations μ_1, \ldots, μ_K in SBMs.
- Intuitively, the algorithm approximately holds because, taking a graph with n nodes, and restricting the attention to the first m nodes, with m < n:

$$\mathbf{Q}_n \hat{\mathbf{x}}_i \longrightarrow \mathbb{N}\{\boldsymbol{\mu}_{z_i}, n^{-1/2} \boldsymbol{\Sigma}(\boldsymbol{\mu}_{z_i})\}, n \to \infty, i = 1, \dots, m.$$

- Importantly, *K*-means, with Euclidean distance, which has been traditionally extensively used in spectral clustering, is **suboptimal** and **unsound** for identifiability reasons.
- Similar asymptotic results are also available for the Laplacian spectral embedding.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG ○○○○○○●○○	Bayesian modelling of embeddings	Results	Conclusion ○	References
ASE AND SRM					

 Simulate a 2-block stochastic blockmodel using the within-community probability matrix

$$\mathbf{B} = \begin{bmatrix} 0.02 & 0.03\\ 0.03 & 0.01 \end{bmatrix}.$$

- Eigenvalues: $\lambda_1 \approx 0.045$ and $\lambda_2 \approx -0.015 \Rightarrow$ GRDPG (**B** is indefinite).
- Simulate the community allocations under two settings:
 - $\boldsymbol{\theta} = (0.5, 0.5)$ (balanced communities),
 - $\boldsymbol{\theta} = (0.9, 0.1)$ (unbalanced communities).
- Simulate two adjacency matrices A_1 and A_2 under both settings, for n = 4,000.
- Take ASE of A_1 and A_2 in \mathbb{R}^2 , say $\hat{\mathbf{X}}_1$ and $\hat{\mathbf{X}}_2$.





Francesco Sanna Passino

ASE and SBMs: an example of the role of \mathbf{Q}_n

- In the simulation, μ_1 and μ_2 are known.
- The **purple** point cloud $\hat{\mathbf{X}}_2$ is reconfigured, and aligned to the **orange** point cloud $\hat{\mathbf{X}}_1$, using two (indefinite) orthogonal transformations estimated from the two ASEs.
- The two representations of the **purple** point cloud are **equivalent**.
- In the CLT, **Q**_n is *unidentifiable*, but it *materially affects (Euclidean) distances* between points.
- The picture confirms that GMMs are preferable over *K*-means.





17/45

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion ○	References
Spectral esti	MATION OF THE ST	OCHASTIC BLOCKMO	DEL		

• Based the asymptotic properties derived in Rubin-Delanchy et al., 2017, the following algorithm should be used for consistent estimation of the latent positions in stochastic blockmodels, when *d* and *K* are known:

Algorithm: Spectral estimation of the stochastic blockmodel (spectral clustering) Input: adjacency matrix \mathbf{A} (or the Laplacian matrix \mathbf{L}), dimension d, and number of

communities $K \geq d$.

1 compute spectral embedding $\hat{\mathbf{X}} = [\hat{x}_1, \dots, \hat{x}_n]^\top$ or $\tilde{\mathbf{X}} = [\tilde{x}_1, \dots, \tilde{x}_n]^\top$ into \mathbb{R}^d ,

2 fit a Gaussian mixture model with K components, **Result:** return cluster centres $\mu_1, \ldots, \mu_K \in \mathbb{R}^d$ and node memberships z_1, \ldots, z_n .

- What about *d* and *K*? In practice the two parameters are estimated sequentially.
 - The latent dimension *d* is chosen according to the scree-plot criterion (Jolliffe, 2002), or the universal singular value thresholding method (Zhu and Ghodsi, 2006).
 - The number of communities K is *usually* chosen using information criteria, conditional on d.

• This talk discusses a novel framework for joint estimation of d and K.

Estimation of *d*: "overshooting"

- Main issues for estimation of *d* and *K*:
 - Sequential approach is sub-optimal: the estimate of K depends on choice of d.
 - Theoretical results only hold for d fixed and known.
 - Distributional assumptions when *d* is misspecified are **not available**.
 - What is the distribution of the last m d columns of the embedding, for m > d?
- How to deal with uncertainty in the estimate of *d*? "Overshooting".
 - Obtain embeddings $\mathbf{X} = [\boldsymbol{x}_1, \dots, \boldsymbol{x}_n]^\top \in \mathbb{R}^{n \times m}, \ \boldsymbol{x}_i \in \mathbb{R}^m$ (ASE or LSE) for some m.
 - Here X represents an estimate of the latent positions (ASE or LSE), dropping "hats" and "tildes".
 - Ideally, m must be $d \leq m \leq n,$ so it can be given an arbitrarily large value.
 - ${\ensuremath{\, \bullet \,}}$ The parameter m is always assumed to be fixed and obtained from a preprocessing step.
 - Choosing an appropriate value of m is arguably **much easier** than choosing the correct d.
 - Under the estimation framework that will be proposed, the correct d can be recovered for any choice of m, as long as $d \le m$.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings ⊙●○○○○○○○○○○	Results	Conclusion O	References

A BAYESIAN MODEL FOR NETWORK EMBEDDINGS

- Choose integer $m \leq n$ and obtain embedding $\mathbf{X} \in \mathbb{R}^{n \times m} \to m$ arbitrarily large.
- Bayesian model for simultaneous estimation of d and $K \rightarrow \text{allow for } d = \text{rank}(\mathbf{B}) \leq K$.

$$\begin{aligned} \boldsymbol{x}_{i}|d, \boldsymbol{z}_{i}, \boldsymbol{\mu}_{\boldsymbol{z}_{i}}, \boldsymbol{\Sigma}_{\boldsymbol{z}_{i}}, \boldsymbol{\sigma}_{\boldsymbol{z}_{i}}^{2} \sim \mathbb{N}_{m} \left(\begin{bmatrix} \boldsymbol{\mu}_{\boldsymbol{z}_{i}} \\ \boldsymbol{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{z}_{i}} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\sigma}_{\boldsymbol{z}_{i}}^{2} \mathbf{I}_{m-d} \end{bmatrix} \right), \ i = 1, \dots, n, \\ (\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})|d \overset{iid}{\sim} \operatorname{NIW}_{d}(\boldsymbol{0}, \kappa_{0}, \nu_{0} + d - 1, \boldsymbol{\Delta}_{d}), \ k = 1, \dots, K, \\ \boldsymbol{\sigma}_{kj}^{2} \overset{iid}{\sim} \operatorname{Inv-} \chi^{2}(\lambda_{0}, \sigma_{0}^{2}), \ j = d + 1, \dots, m, \\ d|\boldsymbol{z} \sim \operatorname{Uniform} \{1, \dots, K_{\varnothing}\}, \\ \boldsymbol{z}_{i}|\boldsymbol{\theta} \overset{iid}{\sim} \operatorname{Discrete}(\boldsymbol{\theta}), \ i = 1, \dots, n, \ \boldsymbol{\theta} \in \mathcal{S}_{K-1}, \\ \boldsymbol{\theta}|K \sim \operatorname{Dirichlet} \left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right), \\ K \sim \operatorname{Geometric}(\omega). \end{aligned}$$

where K_{\varnothing} is the number of non-empty communities.

• Alternative: $d \sim \text{Geometric}(\delta)$.

• Yang et al., 2019, independently and simultaneously proposed a similar frequentist model.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References



Figure 3. Scatterplot of the columns X_1 and X_2 of the ASE.

Figure 4. Scatterplot of the columns X_3 and X_4 of the ASE.

- Simulated GRDPG-SBM with n = 2,500, d = 2, K = 5.
- Nodes allocated to communities with probability $\theta_k = \mathbb{P}(z_i = k) = 1/K$.







Figure 5. Within-cluster and overall means of $X_{:15}$.

Figure 6. Within-cluster variance of $\mathbf{X}_{:25}$.

- Means are approximately 0 for columns with index > *d*.
- Different cluster-specific variances even for columns with index > d.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References



Figure 7. Within-cluster correlation coefficients of $\mathbf{X}_{:30}$.

Figure 8. Marginal likelihood as a function of *d*.

- Reasonable to assume correlation $\rho_{ij}^{(k)} = 0$ for i, j > d.
- Marginal likelihood has maximum at the true value of *d*.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results 00000000000	Conclusion O	References
2					

CURSE OF DIMENSIONALITY



Figure 9. Within-block variance and total variance for the adjacency embedding obtained from a simulated SBM with d = 2, K = 5, n = 500, and well separated means $\mu_1 = [0.7, 0.4]$, $\mu_2 = [0.1, 0.1]$, $\mu_3 = [0.4, 0.8]$, $\mu_4 = [-0.1, 0.5]$ and $\mu_5 = [0.3, 0.5]$, and $\theta = (0.2, 0.2, 0.2, 0.2, 0.2)$.

• For some k and $k': \sigma_{kj}^2 \approx \sigma_{k'j}^2$ for $j \gg d$ and $k \neq k'$.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
Second order	CLUSTERING				

- Bayesian model parsimony: K underestimated for $d \ll m$.
- Possible solution: second order clustering $v = (v_1, \ldots, v_K)$ with $v_k \in \{1, \ldots, H\}$.

• If
$$v_k = v_{k'}$$
, then $\sigma_{kj}^2 = \sigma_{k'j}^2$ for $j > d$:

$$\begin{aligned} \boldsymbol{x}_{i}|d, \boldsymbol{z}_{i}, \boldsymbol{v}_{\boldsymbol{z}_{i}}, \boldsymbol{\mu}_{\boldsymbol{z}_{i}}, \boldsymbol{\Sigma}_{\boldsymbol{z}_{i}}, \boldsymbol{\sigma}_{\boldsymbol{v}_{\boldsymbol{z}_{i}}}^{2} \sim \mathbb{N}_{m} \left(\begin{bmatrix} \boldsymbol{\mu}_{\boldsymbol{z}_{i}} \\ \boldsymbol{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{z}_{i}} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\sigma}_{\boldsymbol{v}_{\boldsymbol{z}_{i}}}^{2} \mathbf{I}_{m-d} \end{bmatrix} \right), \ i = 1, \dots, n, \\ v_{k}|K, H \sim \text{Discrete}(\boldsymbol{\phi}), \ k = 1, \dots, K, \\ \boldsymbol{\phi}|H \sim \text{Dirichlet} \left(\frac{\beta}{H}, \dots, \frac{\beta}{H} \right), \\ H|K \sim \text{Uniform}\{1, \dots, K\}. \end{aligned}$$

- The parameter vk defines clusters of clusters.
- Empirical results show that the model is able to handle $d \ll m$.
- If H = 1, the model is a special case of Raftery and Dean, 2006 \rightarrow ordinal variable selection in clustering.



Francesco Sanna Passino

26/45 Imperial College London

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
Inference					

- Integrate out nuisance parameters μ_k , Σ_k , σ_{jk}^2 and $\theta \rightarrow$ inference on d, K, H and z.
- Inference via MCMC: collapsed Metropolis-within-Gibbs sampler \rightarrow 7 moves.
 - Propose a change in the community allocations z,
 - Propose to split (or merge) two communities,
 - Propose to create (or remove) an empty community,
 - Propose a change in the latent dimension d,
 - ${\ensuremath{\,\circ\,}}$ Propose a change in the second order community allocations v,
 - Propose to split (or merge) two second-order communities,
 - Propose to create (or remove) an empty second-order community.
- Initialisation: *K*-means clustering, choose *K* from scree-plot + uninformative priors (with zero means and variances comparable in scale with the observed data).
- Posterior for *d* is usually similar to a **point mass** → might be worth exploring constrained and unconstrained model.
- The latent dimension *d* could also be treated as a nuisance parameter and **marginalised out** (often not computationally feasible).

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
Extension to 1	DIRECTED AND BIF	PARTITE GRAPHS			

- Consider a **directed graph** with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$.
- The d-dimensional "directed" adjacency embedding (DASE) of ${f A}$ in ${\Bbb R}^{2d}$, is defined as:

$$\hat{\mathbf{X}} = \hat{\mathbf{U}}\hat{\mathbf{D}}^{1/2} \oplus \hat{\mathbf{V}}\hat{\mathbf{D}}^{1/2} = \begin{bmatrix} \hat{\mathbf{U}}\hat{\mathbf{D}}^{1/2} & \hat{\mathbf{V}}\hat{\mathbf{D}}^{1/2} \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{X}}_s & \hat{\mathbf{X}}_r \end{bmatrix},$$

where $\mathbf{A} = \hat{\mathbf{U}}\hat{\mathbf{D}}\hat{\mathbf{V}}^{\top} + \hat{\mathbf{U}}_{\perp}\hat{\mathbf{D}}_{\perp}\hat{\mathbf{V}}_{\perp}^{\top}$ is the SVD decomposition of \mathbf{A} , where $\hat{\mathbf{D}} \in \mathbb{R}^{d \times d}_{+}$ is a diagonal matrix containing the top d singular values in decreasing order, and $\hat{\mathbf{U}} \in \mathbb{R}^{n \times d}$ and $\hat{\mathbf{V}} \in \mathbb{R}^{n \times d}$ contain the corresponding left and right singular vectors.

• Extended model:

$$m{x}_i | d, K, z_i \sim \mathbb{N}_{2m} \left(egin{bmatrix} m{\mu}_{z_i} \ m{0} \ m{\mu}_{z_i}' \ m{0} \end{bmatrix}, egin{bmatrix} m{\Sigma}_{z_i} & m{0} & m{0} & m{0} \ m{0} & \sigma^2_{z_i} \mathbf{I}_{m-d} & m{0} & m{0} \ m{0} & m{0} & m{\Sigma}'_{z_i} & m{0} \ m{0} & m{0} & m{\Sigma}'_{z_i}' & m{0} \ m{0} & m{0} & m{\sigma}^{2\prime}_{z_i} \mathbf{I}_{m-d} \end{bmatrix}
ight).$$

• **Co-clustering**: different clusters for sources and receivers \rightarrow bipartite graphs.



• Simulate bipartite 250×300 graph with K = 5 and K' = 3 obtained from $\mathbf{B} \in [0, 1]^{K \times K'}$ with $B_{k\ell} \sim \text{Beta}(1.2, 1.2), \boldsymbol{\theta} = (1/K, \dots, 1/K), \boldsymbol{\theta}' = (1/K', \dots, 1/K')$, and d = 2.



Francesco Sanna Passino

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
-					



• Means are approximately 0 for columns with index > d, even for a relatively small graph.

30/45

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References



Figure 14. Within-cluster variances of $\hat{\mathbf{X}}_s$.

Figure 15. Within-cluster variances of $\hat{\mathbf{X}}_r$.

- Different cluster-specific variances even for columns with index > d.
- Some evidence of second-order clustering.

Francesco Sanna Passino

Graphs, SBMs and RDPGsBeyond RDPGs: the GRDPGBay0000000000000000000

Bayesian modelling of embeddings

Results

Conclusion

References

SIMULATED DATA: PARAMETER ESTIMATION

(d K)	Model	1	m = 25	5
(<i>a</i> , n)	Model	\bar{d}	\bar{K}_{\varnothing}	\bar{H}_{\varnothing}
	constrained, ASE	2.00	2.00	1.99
(2, 2)	unconstrained, ASE	2.00	2.00	1.99
(2, 2)	constrained, LSE	2.01	2.03	1.99
	unconstrained, LSE	2.02	2.02	1.99
	constrained, ASE	2.00	5.05	1.77
(2, 5)	unconstrained, ASE	2.00	5.07	1.80
(2, 3)	constrained, LSE	2.05	5.10	3.11
	unconstrained, LSE	2.07	5.11	3.10
	constrained, ASE	6.00	7.04	2.10
(6, 7)	unconstrained, ASE	6.00	7.05	2.20
(0, 1)	constrained, LSE	6.00	7.10	2.47
	unconstrained, LSE	6.00	7.07	2.39

Table 1. Results of the inferential procedure for undirected SBMs simulated using different (d, K) pairs, n = 1,000.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion	References
000000	000000000	00000000000	0000000000	0	

SIMULATED DATA: PARAMETER ESTIMATION

(d K)	Model	1	m = 25	
(a, \mathbf{n})	Model	\bar{d}	\bar{K}_{\varnothing}	\bar{H}_{\varnothing}
	constrained, ASE	8.97	9.01	2.08
(0,0)	unconstrained, ASE	9.00	9.01	1.98
(9, 9)	constrained, LSE	9.00	9.02	2.12
	unconstrained, LSE	9.00	9.04	2.11
	constrained, ASE	9.00	12.02	1.96
(0, 12)	unconstrained, ASE	9.00	12.01	1.90
(9, 12)	constrained, LSE	9.00	12.03	2.60
	unconstrained, LSE	9.00	12.02	2.53
	constrained, ASE	10.00	14.78	1.25
(10, 15)	unconstrained, ASE	10.00	14.11	1.27
(10, 10)	constrained, LSE	10.00	14.81	1.81
	unconstrained, LSE	10.00	15.01	1.87

Table 2. Results of the inferential procedure for undirected SBMs simulated using different (d, K) pairs, n = 1,000.

Simulated data: effect of second-order clustering

(d K)	m	H random				H =	K	
(<i>a</i> , n)		\hat{d}	\hat{K}_{\varnothing}	\bar{H}_{\varnothing}	ARI	\hat{d}	\hat{K}_{\varnothing}	ARI
	15	3	5	1.669	1.000	3	5	1.000
(2.5)	50	3	5	1.577	1.000	3	4	0.768
(3 , 5)	150	3	5	1.467	1.000	3	4	0.768
	500	3	5	1.006	1.000	3	4	0.768
	15	9	12	1.979	1.000	9	12	1.000
(9,12)	50	9	12	1.912	1.000	9	12	1.000
	150	9	12	1.875	1.000	9	11	0.942
	500	9	12	1.388	1.000	9	5	0.517

Table 3. Results for the MCMC sampler on simulated undirected SBMs for different values of m, with and without second order clustering, n = 1,000, assuming the unconstrained model.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion ○	References
Santander cy	CLES DATA				



- Aka Boris bikes.
- Santander cycles \rightarrow bike sharing system in central London.
- £2 for access for 24 hours, first 30 minutes of each ride are free. Limited speed.
- Data freely available at https://cycling. data.tfl.gov.uk/, powered by TfL.
- One week of data: 5 11 September, 2018.
- |V| = 783 nodes/stations, |E| = 69,153 (excluding self-loops).

• Undirected graph:

 $A_{ij} = \left\{ \begin{array}{ll} 1 & \text{if at least one journey between stations } i \text{ and } j \text{ is completed}, \\ 0 & \text{otherwise.} \end{array} \right.$

Image: CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=71800653.

Francesco Sanna Passino

35/45 Imperial College London



SANTANDER CYCLES DATA: NUMBER OF CLUSTERS



Figure 16. Adjacency embedding – Posterior histogram of K_{\emptyset} and H_{\emptyset} , unconstrained model, MAP for *d* in red.

Figure 17. Laplacian embedding – Posterior histogram of K_{\emptyset} and H_{\emptyset} , **unconstrained** model, MAP for *d* in **red**.



SANTANDER CYCLES DATA: SCREE-PLOTS



Figure 18. Magnitude of eigenvalues of the adjacency matrix.

Figure 19. Magnitude of eigenvalues of the Laplacian matrix.

• Choice of *d* is consistent with the *elbow* of the scree-plot.

Graphs, SBMs and RDPGs

Beyond RDPGs: the GRDPG

Bayesian modelling of embeddings

Results



Figure 20. Adjacency embedding – Estimated communities for K = 11.

Francesco Sanna Passino

38/45 Imperial College London

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results ○○○○○●○○○	Conclusion ○	References

ENRON E-MAIL NETWORK



- Corpus of e-mails sent by the employees of Enron corporation.
- Data freely available at https://www.cs.cmu.edu/~enron/.
- Version of dataset: May 7, 2015.
- |V| = 184 nodes/employees, |E| = 3,010.
- Extensively analysed in Priebe et al., 2005.

Directed graph:

 $A_{ij} = \left\{ \begin{array}{ll} 1 & \text{if employee } i \text{ sends at least one e-mail to employee } j, \\ 0 & \text{otherwise.} \end{array} \right.$

Image: Paul Rand, https://commons.wikimedia.org/wiki/File:Logo_de_Enron.svg.

39/45



ENRON E-MAIL NETWORK: NUMBER OF CLUSTERS



Figure 21. ASE – Posterior histogram of K_{\emptyset} and H_{\emptyset} , **unconstrained** model, MAP for *d* in **red**.

Figure 22. ASE – Posterior histogram of K_{\emptyset} and H_{\emptyset} , constrained model, MAP for *d* in red.

Francesco Sanna Passino



ENRON E-MAIL NETWORK: NUMBER OF CLUSTERS



Figure 23. ASE – Posterior histogram of K_{\emptyset} , unconstrained model without second order clustering, MAP for *d* in red.

Figure 24. ASE – Posterior histogram of K_{\emptyset} , constrained model without second order clustering, MAP for *d* in red.

ENRON E-MAIL NETWORK: SCREE-PLOT



Figure 25. Singular values of the adjacency matrix.

• Choice of *d* is consistent with the *elbow* of the scree-plot.

CONCLUSION

- Community detection and stochastic blockmodels:
 - Bayesian model for simultaneous selection of K and d in generalised random dot product graphs,
 - Allow for initial misspecification of the arbitrarily large parameter *m*, then refine estimate *d*,
 - Gaussian mixture model (with constraints) based on spectral embedding,
 - Easy to extend to directed and bipartite graphs.
- More details:

Sanna Passino and Heard, 2019 – arXiv: 1904.05333.



Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
References I					

- Athreya, A. et al. (2018). "Statistical Inference on Random Dot Product Graphs: a Survey". In: *Journal of Machine Learning Research* 18.226, pp. 1–92.
- Cape, J., M. Tang, and C. E. Priebe (2018). "On spectral embedding performance and elucidating network structure in stochastic block model graphs". In: arXiv e-prints, arXiv:1808.04855, arXiv:1808.04855. arXiv: 1808.04855 [math.ST].
- Hoff, P. D, A. E. Raftery, and M. S. Handcock (2002). "Latent space approaches to social network analysis". In: *Journal of the American Statistical Association* 97.460, pp. 1090–1098.
- Holland, P. W., K. B. Laskey, and S. Leinhardt (1983). "Stochastic blockmodels: First steps". In: Social Networks 5.2, pp. 109 137.
- Jolliffe, I. T. (2002). Principal Component Analysis. Springer Series in Statistics. Springer.
- Priebe, C. E. et al. (2005). "Scan Statistics on Enron Graphs". In: *Computational & Mathematical Organization Theory* 11.3, pp. 229–247.
- Raftery, A. E. and N. Dean (2006). "Variable Selection for Model-Based Clustering". In: Journal of the American Statistical Association 101.473, pp. 168–178.
- Rohe, K., S. Chatterjee, and B. Yu (2011). "Spectral clustering and the high-dimensional stochastic blockmodel". In: *Annals of Statistics* 39.4, pp. 1878–1915.

Graphs, SBMs and RDPGs	Beyond RDPGs: the GRDPG	Bayesian modelling of embeddings	Results	Conclusion O	References
References II					

- Rubin-Delanchy, P., N.M. Adams, and N.A. Heard (2016). "Disassortativity of computer networks". In: 2016 IEEE Conference on Intelligence and Security Informatics (ISI), pp. 243–247.
- Rubin-Delanchy, P. et al. (2017). "A statistical interpretation of spectral embedding: the generalised random dot product graph". In: *ArXiv e-prints*. arXiv: 1709.05506.
- Sanna Passino, F. and N. A. Heard (2019). "Bayesian estimation of the latent dimension and communities in stochastic blockmodels". In: *arXiv e-prints*. arXiv: 1904.05333.
- Yang, C. et al. (2019). "Simultaneous dimensionality and complexity model selection for spectral graph clustering". In: *arXiv e-prints*. arXiv: 1904.02926.
- Young, S. J. and E. R. Scheinerman (2007). "Random Dot Product Graph Models for Social Networks". In: *Algorithms and Models for the Web-Graph*. Ed. by A. Bonato and F. R. K. Chung. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 138–149.
- Zhu, M. and A. Ghodsi (2006). "Automatic dimensionality selection from the scree plot via the use of profile likelihood". In: *Computational Statistics & Data Analysis* 51.2, pp. 918–930.

Bayesian estimation of the latent dimension and communities in stochastic blockmodels