# Consistency of co-clustering for exchangeable array data

#### David Choi (joint work with Patrick Wolfe, UCL)

Carnegie Mellon University

Oct 6, 2013

#### What is a network?

Network data records interactions between units, such as

- Communication/friendships/etc. between people
- Chemical reactions between proteins

Common depictions for network data:



## Identifying communities in networks

Clustering + post-hoc checks to find community structure



Political blog network, Adamic 2005

- Network equivalent of clustering
- Post-hoc checks are usually required to interpret clusters
- Checks may be effort-intensive should be reserved for statistically significant findings
- But statistical significance is not well understood for community detection

## Example from Ad-Health dataset

Survey of high school friendships







students grouped by school year (black lines) students grouped by race (blue lines) students grouped by race (blue lines) and within race by likelihood-based clustering (red lines)

Could this be due to noise?

### What is known about community detection

Assuming a particular parametric generative model (stochastic blockmodel):

- ► Polynomial time algorithm (spectral clustering) finds correct clusters, with fraction of errors → 0
- ▶ Exhaustive search (profile likelihood) finds correct clusters, with total number of errors  $\rightarrow 0$ 
  - Also yields asymptotically normal parameter estiamtes

**Issue:** Assuming a stochastic blockmodel seems unrealistic – what if "true communities" don't exist in reality?

#### Main result

- Assume network data is generated a general nonparametric model, and then fit by a community detection model known as "stochastic co-blockmodel"
- Then the resulting estimate is a near-optimal "piecewise constant" approximation to the generative model

Implication: detected communities are consistent, in a sense analogous to **histograms** 





### Idea of nonparametric model



Adjacency matrix



Model  $\omega_{\alpha}: [0,1]^2 \mapsto [0,1]$ 

Figures from Lovasz, L. "Very Large Graphs" (2013)

#### Description of nonparametric model

- $\omega_{\alpha}$  maps  $[0,1]^2 \mapsto [0,1]$ , parameterized by  $\alpha$
- ► The directed adjacency matrix A is generated as follows:

$$\begin{split} \alpha &\sim \mathsf{Unif}[0,1],\\ \xi_i \stackrel{\mathrm{iid}}{\sim} \mathsf{Unif}[0,1]\\ \zeta_i \stackrel{\mathrm{iid}}{\sim} \mathsf{Unif}[0,1]\\ A_{ij} |\xi_i, \zeta_j \stackrel{\mathrm{iid}}{\sim} \mathsf{Bernoulli}(\omega_\alpha(\xi_i, \zeta_j)) \end{split}$$

ω<sub>α</sub> cannot be distinguished from (x, y) → ω<sub>α</sub>(π<sub>1</sub>(x), π<sub>2</sub>(y)) for any measure preserving transformations π<sub>1</sub>, π<sub>2</sub>.

Interpretation:  $\xi_i, \zeta_i$  are latent factors controlling *i*'s propensities in sending and receiving links

## Picture of generative process



- $\xi$  and  $\zeta$  are latent
  - Otherwise estimation of  $\omega_{\alpha}$  would be straightforward
- ► Eq(ω<sub>α</sub>): equivalence class of ω<sub>α</sub>, induced by measure-preserving maps on [0, 1]
- Mapping  $\alpha \mapsto \omega_{\alpha}$  is not identifiable from a single network

Note: results will not require  $\omega_{\alpha}$  to be smooth, only measurable

#### A natural model for exchangeable data

Model encompasses all infinitely exchangeable arrays

Definition (Exchangeable arrays) Binary array  $\{A_{ij}\}_{i,j=1}^{\infty}$  is separately exchangeable if

 $\mathbb{P}(A_{ij} = X_{ij}, i, j \in [n]) = \mathbb{P}(A_{ij} = X_{\eta_1(i)\eta_2(j)}, i, j \in [n])$ 

for all *n*, all permutations  $\eta_1, \eta_2$  of [*n*], and all *X*.

#### Theorem (Aldous-Hoover)(di Fenetti)

Let  $\{A_{ij}\}_{i,j=1}^{\infty}$  be a separately exchangeable binary array. There exists  $\omega$  which generates A.

Given a static snapshot of a network, natural to require the model to be invariant to a permutation of the adjacency matrix

We'll discuss later: are sparse graphs exchangeable?

### Approximating $\omega_{\alpha}$ by a piecewise constant function

Question: Can we fit a piecewise constant approximation to the generative model?



Parameters  $\phi \equiv (\mu, \nu, \theta)$  describe  $\omega_{\phi}$ :

- Vectors μ, ν : boundaries of the piecewise constant regions (i.e., nonuniform grid)
- Matrix  $\theta$ : heights of the piecewise constant regions

#### Fitting criteria for approximation

Fit by finding optimal clustering

We can fit  $\phi = (\mu, \nu, \theta)$  to an observed adjacency matrix A by various criteria:

Likelihood: 
$$L_A(\mu, \nu, \theta) = \max_{S, T} \sum_{i, j} \log \mathbb{P}(A_{ij} \mid \theta_{S(i)T(j)})$$
  
 $\ell_2 \text{ error: } R_A(\mu, \nu, \theta) = \min_{S, T} \sum_{i, j} (A_{ij} - \theta_{S(i)T(j)})^2$ 

Mappings S and T assign nodes to K clusters, and are constrained to have assignment proportions matching  $\mu$  and  $\nu$ .

### Metrics for approximation quality

In what sense might  $\omega_{\phi}$  be close to  $\omega_{\alpha}$ ?

Fitting criteria correspond to risk functions that measure distance to the unknown  $\omega:$ 

$$L_{\omega}(\phi) = \inf_{\omega \in \text{Eq}(\omega_{\alpha})} \mathbb{E} \operatorname{KL}(\omega(\xi, \zeta), \omega_{\phi}(\xi, \zeta))$$
$$R_{\omega}(\phi) = \inf_{\omega \in \text{Eq}(\omega_{\alpha})} \mathbb{E}(\omega(\xi, \zeta) - \omega_{\phi}(\xi, \zeta))^{2}$$

 $Eq(\omega_{\alpha})$ : equivalence class of  $\omega_{\alpha}$  (all measure preserving maps of [0,1])

KL: Kullback-Leibler divergence

#### First main result

Excess risk goes to zero

#### Theorem

For  $\overline{\phi}$  maximizing<sup>1</sup> the true  $L_{\omega}$ , and  $\hat{\phi}$  maximizing its proxy  $L_A$ , such that  $\overline{\theta}$  and  $\hat{\theta}$  are bounded away from 1 or 0,

$$L_\omega(ar\phi) - L_\omega(\hat\phi) = O_P\left(rac{1}{n^{1/4}}
ight).$$

Also, for  $\overline{\phi}$  minimizing  $R_{\omega}$ , and  $\hat{\phi}$  minimizing  $R_{A}$ ,

$$R_\omega(\hat{\phi}) - R_\omega(ar{\phi}) = O_P\left(rac{1}{n^{1/4}}
ight).$$

The estimate will be asymptotically optimal in terms of minimizing risk.

<sup>1</sup>over class frequencies  $\mu, \nu$  which are multiples of 1/n

#### Significance of detected community structure Corollary of Theorem 1

The estimate  $\hat{\phi}$  approximates  $\phi^*$ , for some some blockmodel  $\phi^*$  induced by partitioning  $\omega_{\alpha}$ :



By approximate we mean:

- $\blacktriangleright$  Identical region boundaries  $\mu$  and  $\nu$
- Heights of piecewise constant regions are close:

$$\|\theta^* - \hat{\theta}\|_2^2 = O_P\left(\frac{1}{n^{1/4}}\right).$$

Implication: if you see an extreme partition in the data, then one also exists in the model  $% \left( {{{\left[ {{{\left[ {{{\left[ {{{c_{1}}} \right]}}} \right]}_{\rm{c}}}}}} \right)$ 

#### Discussion (1/2)Why new approach was needed

In standard learning problems, we want to choose φ to minimize expected loss compared to ω<sub>α</sub>:

$$R_{\omega}(\phi) = \mathbb{E}(\omega_{lpha}(\xi,\zeta) - \omega_{\phi}(\xi,\zeta))^2$$

• We can't evaluate  $R_{\omega}$ , so we use noisy samples instead:

$$R_{\mathcal{A}}(\phi) = \frac{1}{n^2} \sum_{i,j} (A_{ij} - \omega_{\phi}(\xi_i, \zeta_j))^2$$

- We'd try to show that  $R_A(\phi) \approx R_\omega(\phi)$  uniformly over  $\phi$ 
  - Note: the n<sup>2</sup> samples (ξ<sub>i</sub>, ζ<sub>j</sub>) are not independent. Better to think of as n independent samples.

## Discussion (2/2)

Why new approach was needed

What happens now that  $(\xi_i, \zeta_j)$  are latent?

•  $\omega_{\alpha}$  is indistinguishable under measure-preserving maps:

$$R_{\omega}(\phi) = \inf_{\omega \in \mathsf{Eq}(\omega_{lpha})} \mathbb{E}(\omega(\xi,\zeta) - \omega_{\phi}(\xi,\zeta))^2$$

• We have to estimate  $\phi$  and the latent  $\xi, \zeta$  jointly:

$$R_{\mathcal{A}}(\phi) = \min_{\xi_i,\zeta_j} \frac{1}{n^2} \sum_{i,j} (A_{ij} - \omega_{\phi}(\xi_i,\zeta_j))^2,$$

Only *n* samples – won't suffice for union bound

Our approach: optimization over  $\xi_i, \zeta_j$  has additional structure

### Second main result

Convex hulls of partition spaces  $\mathcal{F}_{A}$  and  $\mathcal{F}_{\omega}$ 



## Second main result

Convex hulls of partition spaces  $\mathcal{F}_A$  and  $\mathcal{F}_\omega$ 



**Theorem**:  $\max_{\mu,\nu} d_{\text{Hausdorff}} \left( \text{conv}(\mathcal{F}_{\omega}), \text{conv}(\mathcal{F}_{A}) \right) = O_{P} \left( \frac{1}{n^{1/4}} \right).$ 

## Relation to excess risk bound (1/2)

Prelude: understanding convex hulls



- Notion of supporting hyperplane in each direction  $\theta$
- The set of supporting hyperplanes in all directions induces the convex hull

#### Relation to excess risk bound (2/2)

Simple algebra (expanding the square):

$$R_{A}(\mu,\nu,\theta) = \min_{S,T} \sum_{i,j} \left( \theta_{S(i)T(j)} - A_{ij} \right)^{2},$$
  
=  $\sum_{i,j} \theta_{S(i)T(j)}^{2} - 2 \max_{S,T} \langle \theta, A/ST \rangle + \sum_{i,j} A_{ij}^{2}$   
supporting hyperplane

- ► R<sub>A</sub>(θ) is determined by boundary of convex hull of F<sub>A</sub> in direction θ
- Similar result relates  $R_{\omega}$  to conv( $\mathcal{F}_{\omega}$ ).
- Convergence of convex hulls  $\rightarrow$  convergence of risk functions

## Intuition of proof

The supporting hyperplane

$$\max_{S,T} \langle \theta, A/ST \rangle$$

involves an optimization over assignments S, T of all N nodes to clusters

- ► Show if √N nodes are assigned optimally, and remainder are assigned "greedily", clustering is near-optimal<sup>2</sup>
- This implies the intrinsic dimensionality is  $\sqrt{N}$
- ► Hence, *N* samples will suffice to estimate supporting hyperplane

<sup>&</sup>lt;sup>2</sup>Alon et. al. "Random sampling and approximation of MAX-CSPs." 2003.

#### Relation to graph limits literature

- Approach is taken from Borgs et. al. (2008) work on graph limits
  - Defines convergence of dense graphs to be convergence of all subgraph frequencies
  - Closure (i.e. the limit objects) shown to be ω<sub>α</sub>, termed "graphons"
  - Proved law of large numbers under this sense of convergence
- Relation to this work (roughly):
  - They show  $\mathcal{F}_A \to \mathcal{F}_\omega$ , instead of  $\operatorname{conv}(\mathcal{F}_A) \to \operatorname{conv}(\mathcal{F}_\omega)$ .
  - Their proof requires Szemeredi lemma and has exponentially slower rate
- Multiple notions of convergence exist for sparse graphs, with closure not yet known

#### Dense vs. sparse graphs

Is exchangeablity assumption always appropriate?

- $\blacktriangleright$  Exchangeability is bad for modeling asymptotics when graph density  $\rightarrow$  0.
  - Assumes that  $\omega_{\alpha}$  is fixed with *n*
  - ▶ One approach: assume model is  $\rho_n \omega_\alpha$ , where  $\rho_n \to 0$  and normalize risk by  $\rho_n$ 
    - Our results carry over; risk bound of O<sub>P</sub>(1/√d) if degree d = ω(log n)
    - Issue: model becomes very smooth for large n
- For finite data, when is the sparse asymptotic regime appropriate?
  - Opinion: issue is one of enforcing sparsity and connectedness at once
    - ▶ If average degree is *d*, then expected number of isolates under  $\omega_{\alpha}$  is at least  $e^{-d}$

### Conclusions

#### Summary

- Blockmodels may be useful for exploratory data analysis of network data
- Even if generative model is not even approximately a blockmodel
- Akin to taking histograms

Future work

- Symmetric graphs
- Sparse graphs with  $\omega(\log n)$  degree.
- Nonparametric estimation: letting model complexity increase with n so that total risk → 0.
- Conjecture: proof techniques may be useful for other latent variable models.

#### Backup slides

#### Description of co-blockmodel

Let  $\mu, \nu$  be points on *K*-simplex. Let  $\theta \in [0, 1]^{K \times K}$ .

#### Stochastic co-blockmodel

Given  $\mu, \nu, \theta$ , the graph  $G = (V_1, V_2, E)$  and its adjacency matrix  $A \in \{0, 1\}^{m \times n}$  are generated as follows:

- 1. Generate latent variables  $S = (S_1, ..., S_m) \sim \mu$ , and  $T = (T_1, ..., T_m) \sim \nu$
- 2. Let  $A_{ij}$  be Bernoulli with parameter  $\theta_{S(i)T(j)}$ . Connect (i, j) if  $A_{ij} = 1$ .

Intuition: Discrete number of classes, with class probabilities  $\mu$  and  $\nu$ , and connections probabilities  $\theta.$ 

Note: Only identifiable up to label-switching.

Similar to Borgs et al. (2008), Alon et al. (2003)

$$\begin{split} \max_{S,T} \langle \theta, A/ST \rangle &= \max_{ST} \sum_{i=1}^{m} \sum_{j=1}^{n} \theta_{S(i)T(j)} A_{ij} \\ &\to \langle \theta, A/\hat{S}\hat{T} \rangle, \, \text{where } \hat{S} \text{ and } \hat{T} \text{ given by:} \end{split}$$

Similar to Borgs et al. (2008), Alon et al. (2003)

$$\begin{split} \max_{S,T} \langle \theta, A/ST \rangle &= \max_{ST} \sum_{i=1}^{m} \sum_{j=1}^{n} \theta_{S(i)T(j)} A_{ij} \\ &\to \langle \theta, A/\hat{S}\hat{T} \rangle, \, \text{where } \hat{S} \text{ and } \hat{T} \text{ given by:} \end{split}$$

 $\blacktriangleright$  Randomly sample rows  ${\cal I}$  and columns  ${\cal J}$ 



Similar to Borgs et al. (2008), Alon et al. (2003)

$$\begin{split} \max_{S,T} &\langle \theta, A/ST \rangle = \max_{ST} \sum_{i=1}^{m} \sum_{j=1}^{n} \theta_{S(i)T(j)} A_{ij} \\ &\to \langle \theta, A/\hat{S}\hat{T} \rangle, \, \text{where } \hat{S} \text{ and } \hat{T} \text{ given by:} \end{split}$$



Similar to Borgs et al. (2008), Alon et al. (2003)

$$\begin{split} \max_{S,T} &\langle \theta, A/ST \rangle = \max_{ST} \sum_{i=1}^{m} \sum_{j=1}^{n} \theta_{S(i)T(j)} A_{ij} \\ &\to \langle \theta, A/\hat{S}\hat{T} \rangle, \, \text{where } \hat{S} \text{ and } \hat{T} \text{ given by:} \end{split}$$



Similar to Borgs et al. (2008), Alon et al. (2003)

$$\begin{split} \max_{S,T} \langle \theta, A/ST \rangle &= \max_{ST} \sum_{i=1}^{m} \sum_{j=1}^{n} \theta_{S(i)T(j)} A_{ij} \\ &\to \langle \theta, A/\hat{S}\hat{T} \rangle, \, \text{where } \hat{S} \text{ and } \hat{T} \text{ given by:} \end{split}$$



Show that  $\hat{S}$  and  $\hat{T}$  involve  $K^{|\mathcal{I}|+|\mathcal{J}|} \ll K^n$  possible labelings

## Simulation results

#### Data generated by non-blockmodel



Top: excess risk.

Bottom: KL divergence of blockmodel approximation (grey line is optimal).  $\beta$ : True model more resembles blockmodel for large  $\beta$ .