



Community Detection via Random and Adaptive Sampling

Alexandre Proutiere (KTH)

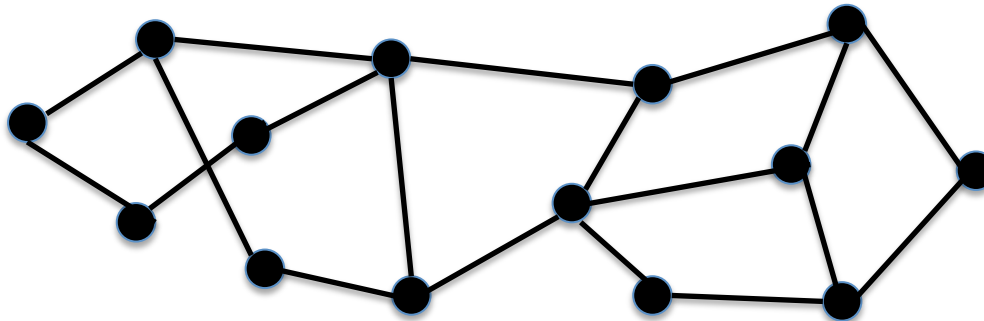
Lunteren Conference 2015

Community detection in networks

Objective: Extract K communities in a network of n nodes from random *observations*. $K \ll n$. Here finite K , very large n .

Observations

1. A graph of interactions or similarities



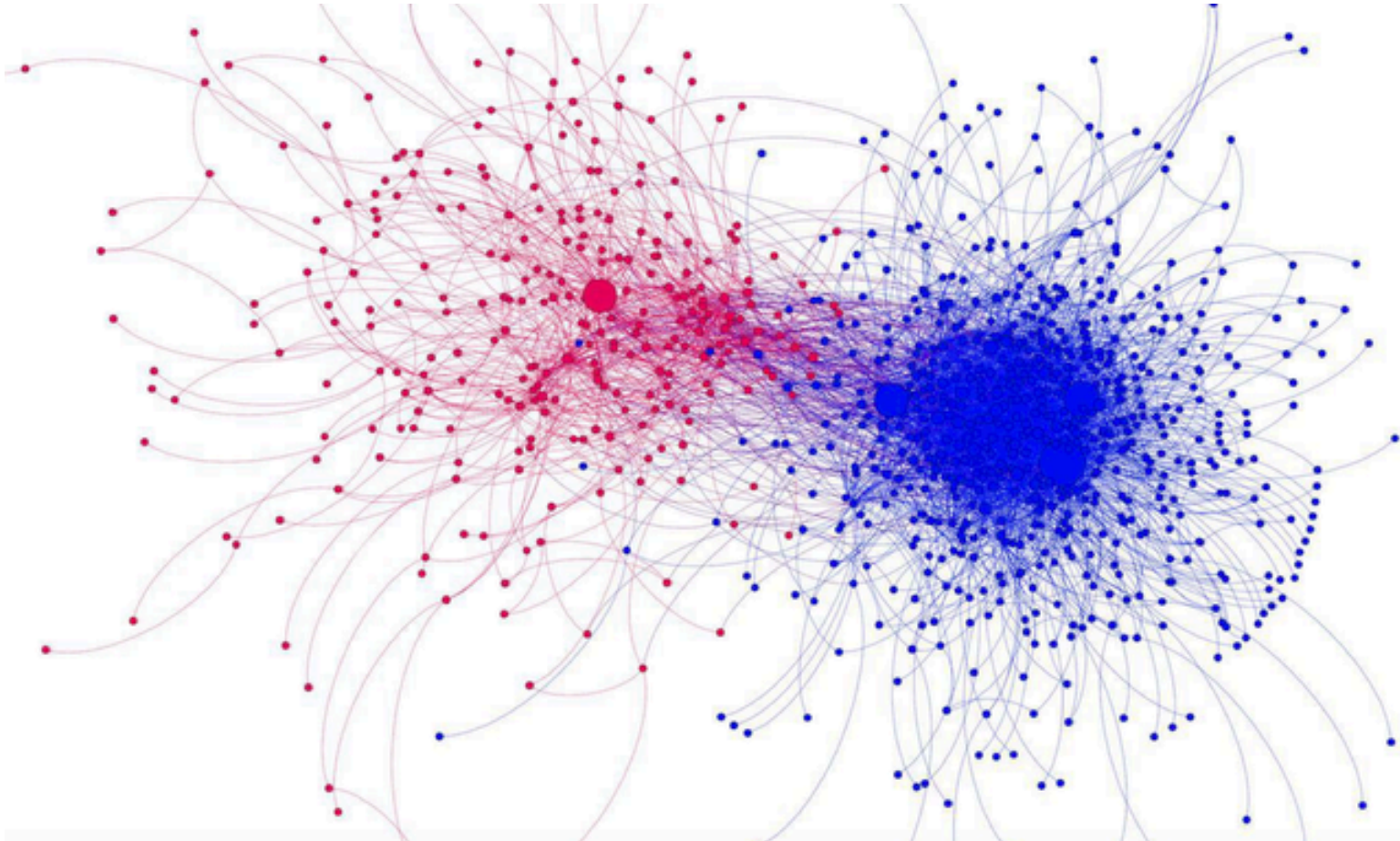
2. General sampling framework

Applications

- Social networks: recommendation systems, targeted advertisement
- Biology: the role of proteins
- Distributed computing: balanced partitions
- Communication networks: caching, pro-active resource allocation (user mobility)
- ...

Applications

- Politics: e.g. tweets related to Ferguson (Emma Pierson)



Applications

- Related work: classify researchers who worked on community detection in three clusters (physicists, mathematicians, computer scientists)

Arora, Rao, Newman, Coja-Oghlan, Jerrum, Chen, Frieze, McSherry, Dyer, Sorkin, Kannan, Vempala, Vetta, Fortunato, Decelle, Krzakala, Karp, Condon, Reichart, Sanghavi, Nadakuditi, Girvan, Mosel, Sly, Rohe, Chatterjee, Yu, Massoulie, Lelarge, Vazirani, Karger, Feld, Fischer, Kleinberg, Gibson, Raghavan, Hopcroft, Khan, Kulis, Santo, Wellman, Hogan, Berg, White, Boorman, Kelley, Xie, Kumar, Mathieu, Schudy, Alon, Krivelevich, Sudakov, Xu, Achlioptas, Kahale, Feige, Zdeborova, Carson, Giesen, Mitshe, Proutiere, Shamir, Tsur, Hassibi, Oymak, Ames, Parrilo, Holland, Laskley, Pothen, Simon, Liou, Girvan, Chauhan, Leone, Ball, Karrer, Abbe, ...

This talk

Under which conditions communities can be accurately detected? And how?

How can we deal with extremely large networks ($\gg 10^8$ nodes)?
(hard to even store the adjacency matrix in the RAM)

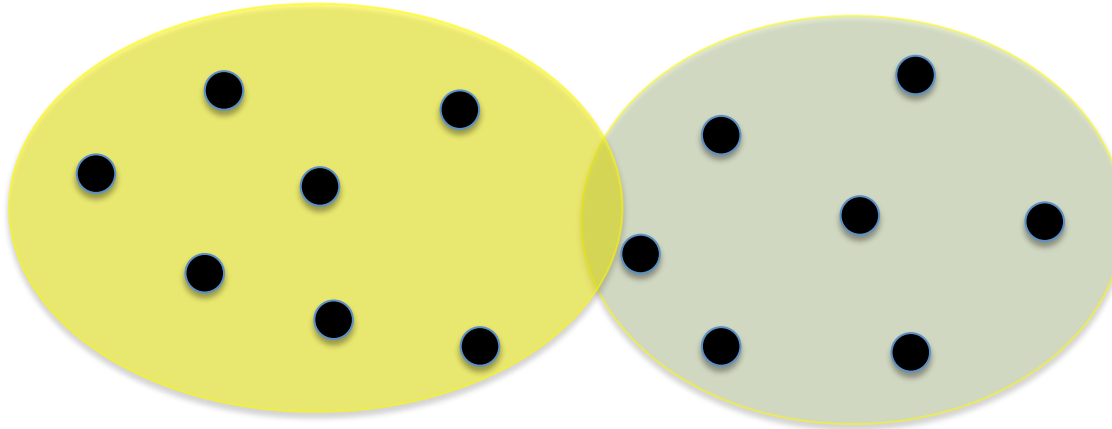
Outline

1. Community detection in the Stochastic Block Model
2. General sampling framework
3. Streaming, memory-limited algorithms

From joint work with Se-Young Yun and Marc Lelarge

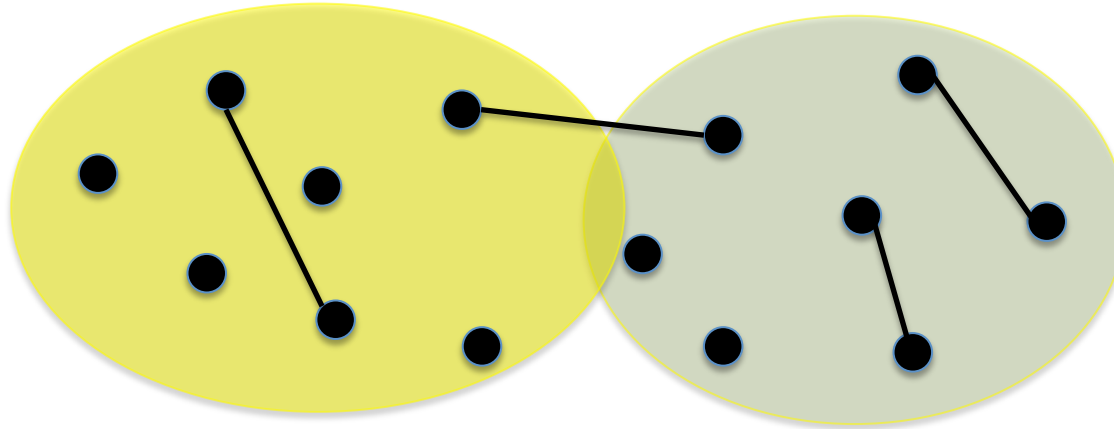
1. Community Detection in the Stochastic Block Model

Stochastic Block Model (SBM)



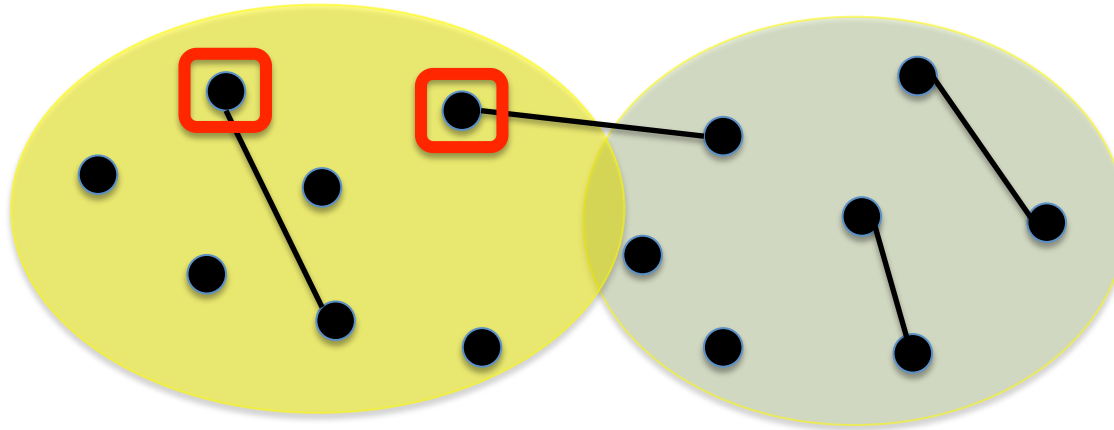
- The graph is built by considering each pair of nodes once
 - If in the same community: put an edge with probability p
 - Else: put an edge with probability $q < p$

Stochastic Block Model (SBM)



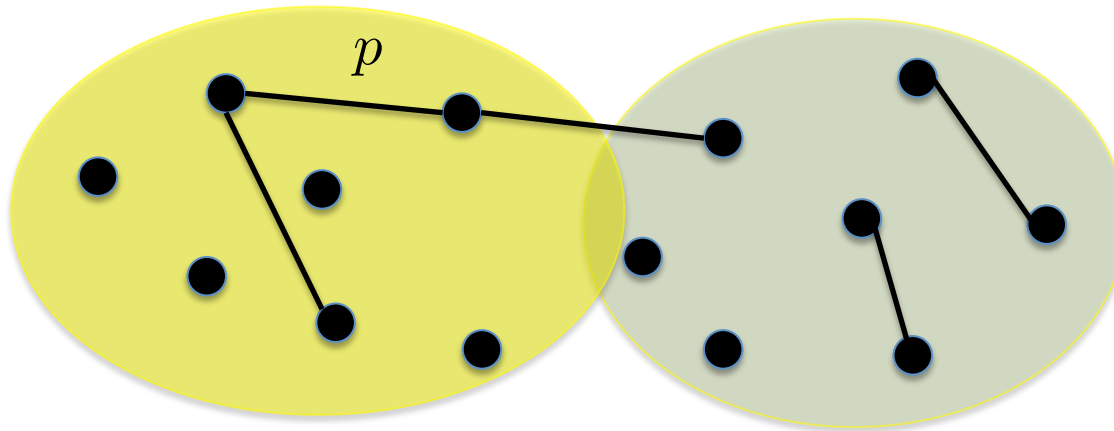
- The graph is built by considering each pair of nodes once
 - If in the same community: put an edge with probability p
 - Else: put an edge with probability $q < p$

Stochastic Block Model (SBM)



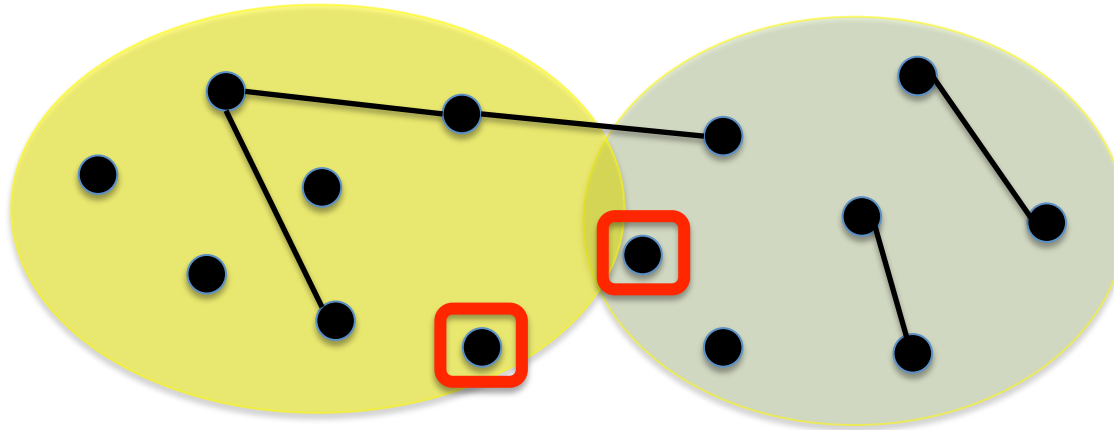
- The graph is built by considering each pair of nodes once
 - If in the same community: put an edge with probability p
 - Else: put an edge with probability $q < p$

Stochastic Block Model (SBM)



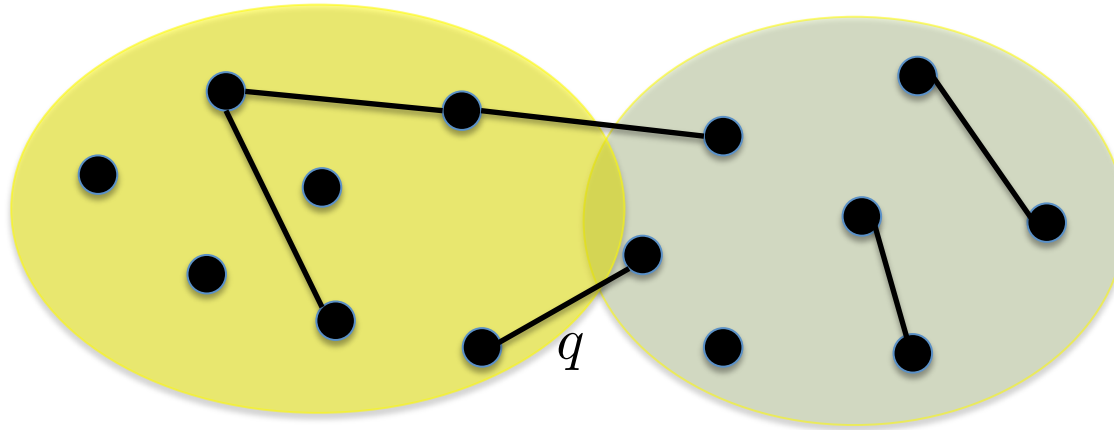
- The graph is built by considering each pair of nodes once
 - If in the same community: put an edge with probability p
 - Else: put an edge with probability $q < p$

Stochastic Block Model (SBM)



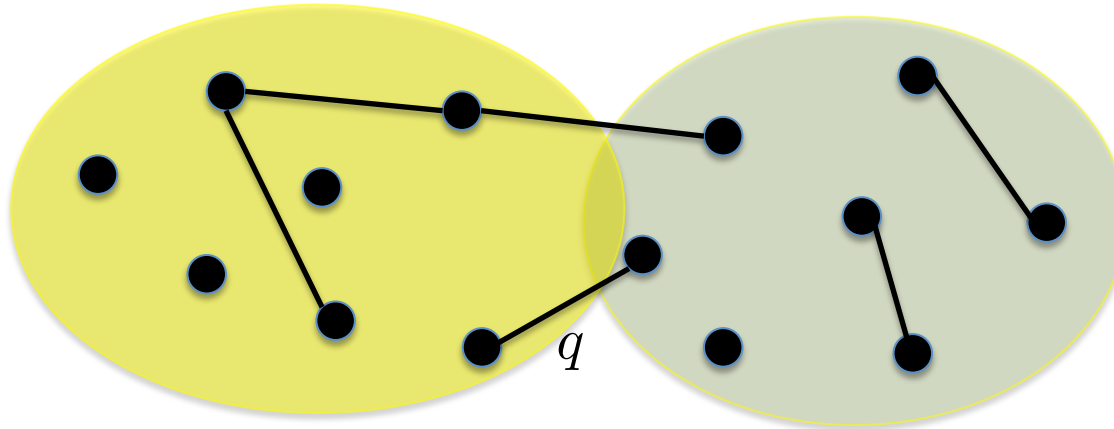
- The graph is built by considering each pair of nodes once
 - If in the same community: put an edge with probability p
 - Else: put an edge with probability $q < p$

Stochastic Block Model (SBM)



- The graph is built by considering each pair of nodes once
 - If in the same community: put an edge with probability p
 - Else: put an edge with probability $q < p$

Stochastic Block Model (SBM)



- Network size: n nodes, n tends to ∞
- Sparse interaction: $p, q = o(1)$
 - Very sparse $p, q \sim 1/n$
 - Sparse $p, q \sim f(n)/n$, $f(n) = \omega(1)$
- Dense interaction: $p, q = O(1)$

Performance metrics

Proportion of misclassified nodes under π : $\varepsilon^\pi(n)$

1. **Asymptotic detection:** an algorithm *detects* the clusters if it does better than the algorithm that randomly assigns nodes to clusters
2. **Accurate asymptotic detection:** an algorithm π is asymptotically accurate if $\lim_{n \rightarrow \infty} \mathbb{E}[\varepsilon^\pi(n)] = 0$

Asymptotic detection in the SBM

- Two communities of equal sizes, sparse case $p = \frac{a}{n}$, $q = \frac{b}{n}$

Theorem 1 (Mossel-Neeman-Sly 2012)

If $a - b < \sqrt{2(a + b)}$, then asymptotic detection is impossible.

Theorem 2 (Massoulié 2013)

If $a - b > \sqrt{2(a + b)}$, then there exists an algorithm leading to clusters that are positively correlated with the true clusters.

Conjectured by Decelle-Krzakala-Moore-Zdeborova 2012

Non-rigorous spectral analysis

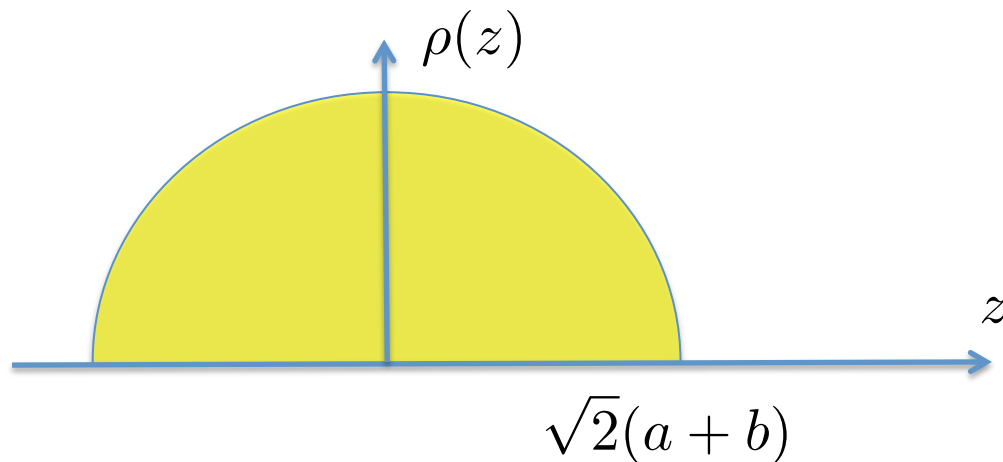
(M. Newman)

- Average adjacency matrix

$$\mathbb{E}[A] = \frac{1}{2}(a+b)11^T + \frac{1}{2}(a-b)uu^T$$

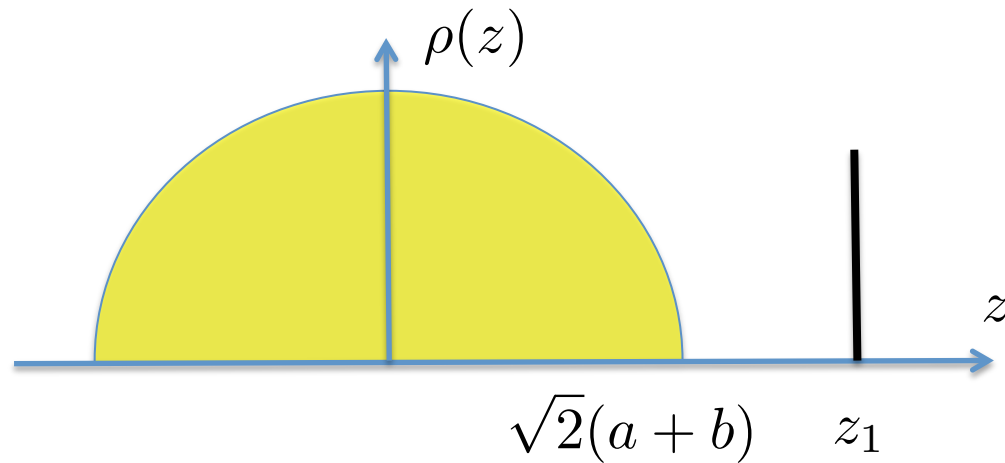
$$1 = \frac{1}{\sqrt{n}}(1, \dots, 1)^T, \quad u = \frac{1}{\sqrt{n}}(1, \dots, 1, -1, \dots, -1)^T$$

- Noisy observation: $A = \mathbb{E}[A] + X$
- Spectral density of noise matrix X (Wigner semicircle law)



Non-rigorous spectral analysis

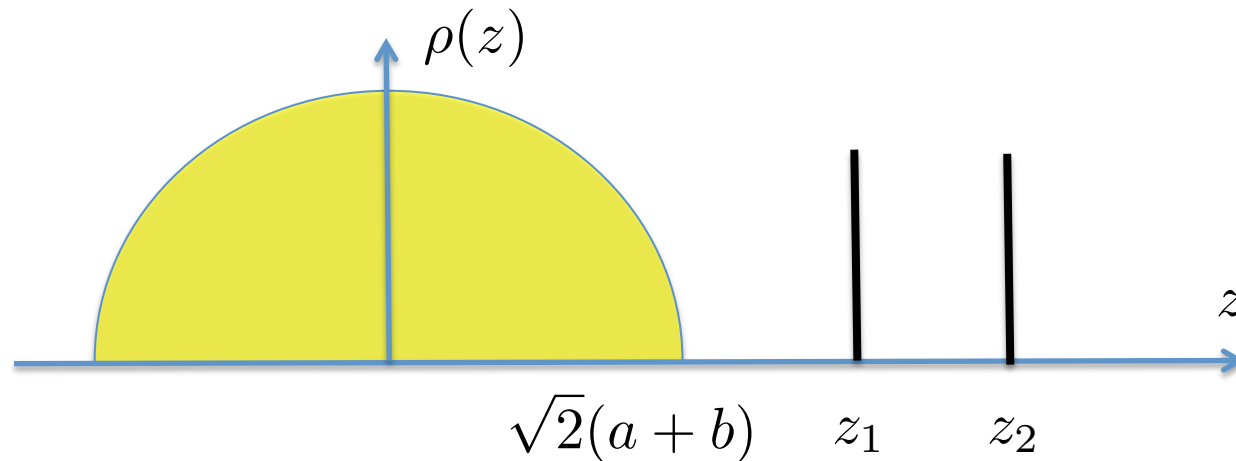
- Spectral density of the modularity matrix: $\frac{1}{2}(a - b)uu^T + X$



$$z_1 = \frac{1}{2}(a - b) + \frac{a + b}{a - b}$$

Non-rigorous spectral analysis

- Spectral density of the observed matrix:



- Communities are detectable if $z_1 > \sqrt{2}(a + b)$
- Method: find z_1 and the corresponding eigenvector u

Examples of algorithms

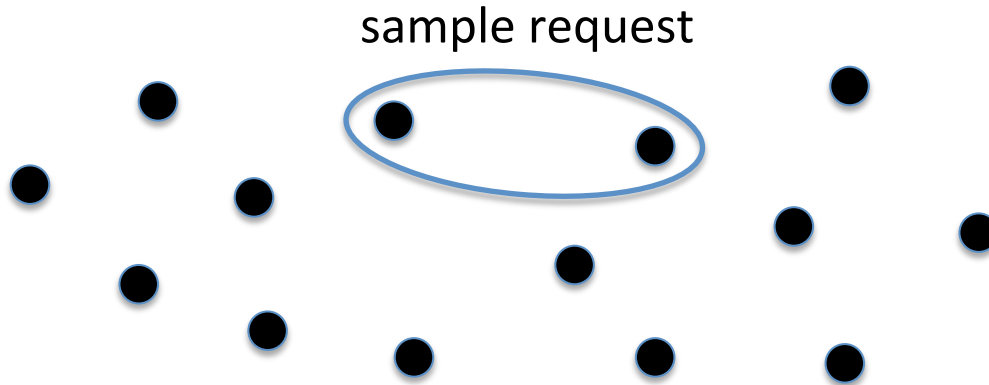
- Maximum Likelihood (NP hard problem)
 - Exact solution: Belief Propagation (no proof)
 - Convex relaxation (akin compressed sensing): poor performance guarantees
- Spectral method
 - Find a rank- K approximation of the adjacency matrix
(+ Trimming + Post-processing)

Open problems

- Very sparse graphs: condition for asymptotic detection with more than two communities
- General graphs: condition on p, q, n for asymptotically accurate detection? (**this talk**)
- General graphs: what is the optimal scaling of $\varepsilon^\pi(n)$?

2. General Sampling Framework

Sampling framework



- Large data set available: many *samples* for the interaction of each pair of nodes
- Sample for a given node pair: Bernoulli with mean p if nodes are in the same cluster, with mean q otherwise
- Sample budget: T

Sampling strategies

- **Non-adaptive random strategies**

- The pair of nodes sampled in round t does not depend on past observations, and is chosen uniformly at random
- S1: sampling with replacement
- S2: sampling without replacement

- **Adaptive strategies**

- The pair of nodes sampled in round t depends on past observations

- Classical SBM: random sampling without replacement, and

$$T = n(n - 1)/2$$

Objectives

- Performance metric: proportion of misclassified nodes $\varepsilon(n, T)$
Asymptotically accurate detection: $\lim_{n \rightarrow \infty} \mathbb{E}[\varepsilon(n, T)] = 0$
- Non-adaptive sampling:
 - Necessary conditions on n, T, p, q for the existence of asymptotically accurate algorithms
 - Asymptotically accurate clustering algorithms
- Adaptive sampling:
 - Necessary conditions on n, T, p, q for the existence of asymptotically accurate joint sampling and clustering algorithms
 - Asymptotically accurate sampling and clustering algorithms

Fundamental limits

- Non-adaptive sampling:

$$\kappa_1(n, T) = T \frac{2(n-2)}{n(n-1)} \min\{KL(q, p), KL(p, q)\} \\ + 2 \sqrt{\frac{4T(n-2)}{n(n-1)} \left[\min\{q, 1-p\} \left(\log \frac{p(1-q)}{q(1-p)} \right)^2 + \left(\log(\min\{\frac{p}{q}, \frac{1-q}{1-p}\}) \right)^2 \right]}$$

Theorem 3 Under random sampling strategy S1 or S2, for any clustering algorithm π , we have:

$$\mathbb{E}[\varepsilon^\pi(n, T)] \geq \frac{1}{8} \exp(-\kappa_1(n, T)),$$

Fundamental limits

- Non-adaptive sampling -- necessary conditions for asymptotically accurate detection:

$$\frac{T}{n} = \omega(1), \quad \frac{T}{n} \min(KL(q, p), KL(p, q)) = \omega(1),$$

- Dense interaction: $p, q = \Theta(1)$

$$T(p - q)^2 / n = \omega(1)$$

- Sparse interaction: $p, q = o(1)$

$$T(p - q)^2 / (pn) = \omega(1)$$

Fundamental limits

- Adaptive sampling:

Theorem 4 For asymptotically accurate detection, we need:

$$\min\{p, 1 - q\} \frac{T}{n} = \Omega(1) \quad \text{and} \quad \frac{T}{n} \max(KL(q, p), KL(p, q)) = \omega(1).$$

- Example: $p = \frac{\log n}{n}$ $q = \frac{\sqrt{\log n}}{n}$
 - Non-adaptive sampling: $\frac{T}{n} = \omega\left(\frac{n}{\log(n)}\right)$
 - Adaptive sampling: $\frac{T}{n} = \Omega\left(\frac{n}{\log(n)}\right)$

Algorithms for non-adaptive sampling

- Spectral algorithms (an extension of Coja-Oghlan's algorithm)
 1. From random samples, build an observation matrix
 2. Trimming (remove nodes with too many interactions)
 3. Spectral decomposition (find the largest eigenvalues and corresponding eigenvectors)
 4. Greedy improvement (for each node compare the number of interactions with the various clusters)

Performance

Theorem 5 Assume that:

$$\frac{(p - q)^2}{p} \frac{\alpha T}{n} = \omega(1), \quad \frac{(p - q)^2}{p} \frac{\alpha T}{n} \geq \log\left(p \frac{T}{n}\right).$$

Then with high probability:

$$\varepsilon^{SP}(n, T) \leq 8 \exp\left(-\frac{(p - q)^2}{20p} \frac{\alpha T}{n}\right).$$

- The algorithm is asymptotically accurate under the necessary conditions for accurate detection in the case of random sampling
- The necessary conditions for accurate detection are tight!

Algorithms for adaptive sampling

- Spatial coupling idea: find reference kernels and build the clusters from these kernels
1. Kernels: select $n/\log(n)$ nodes and use $T/5$ samples to classify these nodes (using the previous spectral algorithm)
 2. Select one of remaining nodes. Sample $T/3n$ pairs between the selected nodes to each kernel. Classify the node.
 3. Repeat 2. until no remaining node or budget

Performance

Theorem 6 Assume that:

$$\frac{(p - q)^2}{p + q} \frac{T}{n} = \Omega(1), \quad \frac{T}{n} \max(KL(q, p), KL(p, q)) = \omega(1).$$

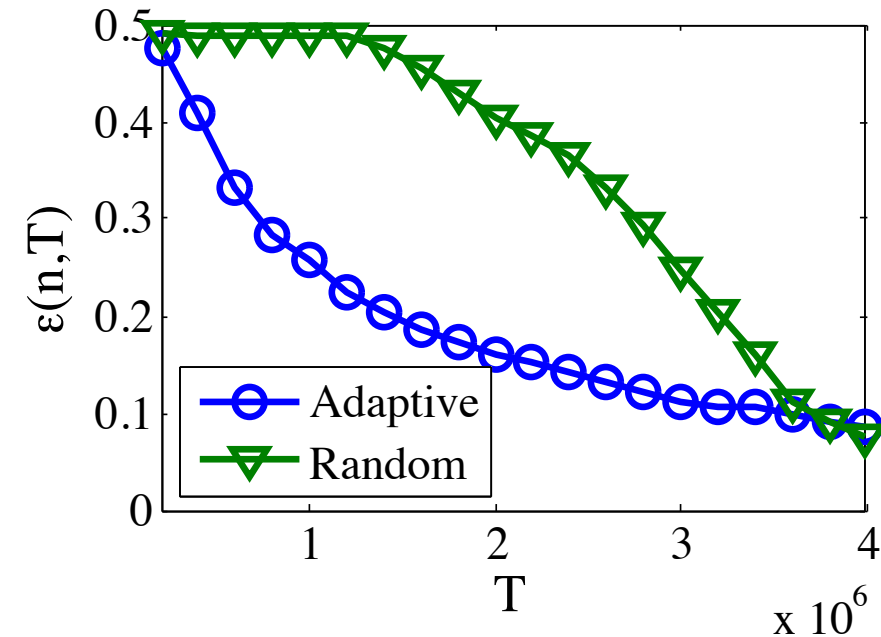
Then with high probability:

$$\varepsilon^{ASP}(n, T) \leq \exp \left(-\frac{T}{6n} (KL(q, p) + KL(p, q)) \right).$$

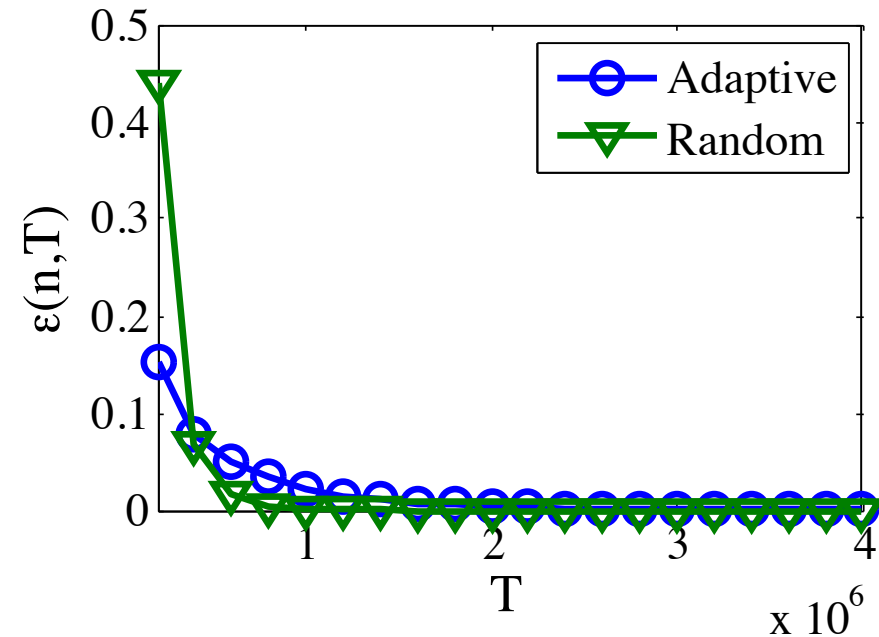
- The algorithm is asymptotically accurate under the necessary conditions for accurate detection in the case of adaptive sampling
- The necessary conditions for accurate detection are tight!

Random vs. adaptive sampling

- $n = 4,000$



$p = 0.01, q = 0.005$

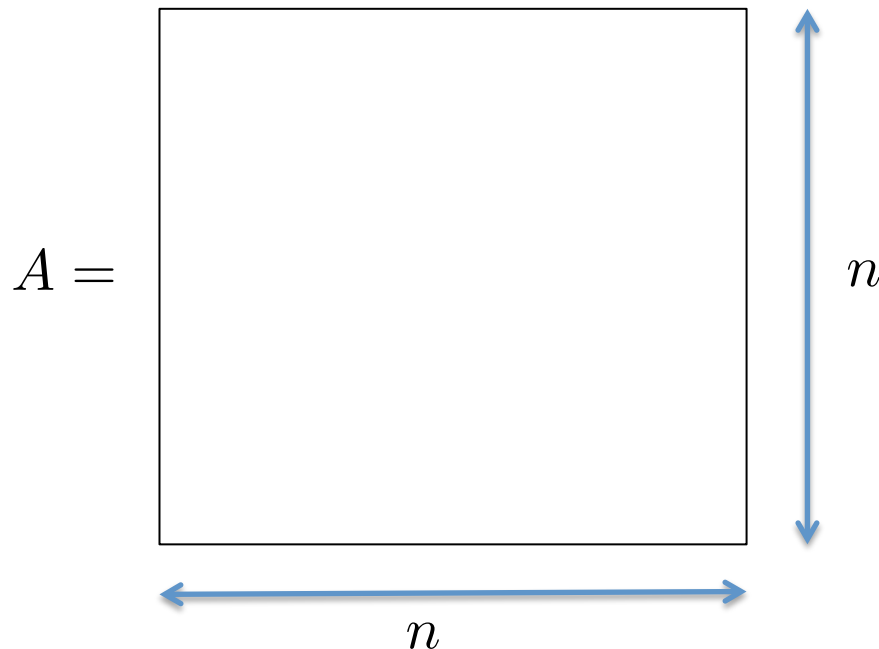


$p = 0.1, q = 0.05$

3. Streaming, Memory-limited Algorithms

Memory and streaming issues

- Storing and manipulating the adjacency matrix in RAM could be impossible
- Data about a node may arrive sequentially (one column at a time – e.g. recommendation systems)



How to deal with
 $n = 10^7$?

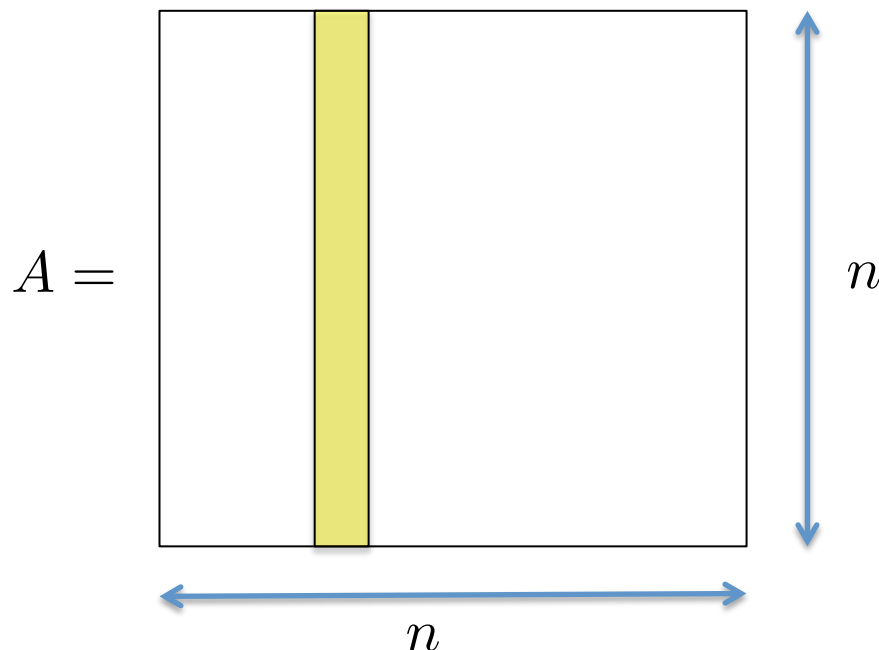
SBM with $f(n) = \omega(1)$

$$p = a \frac{f(n)}{n}$$

$$q = b \frac{f(n)}{n}$$

Memory and streaming issues

- Storing and manipulating the adjacency matrix in RAM could be impossible
- Data about a node may arrive sequentially (one column at a time – e.g. recommendation systems)

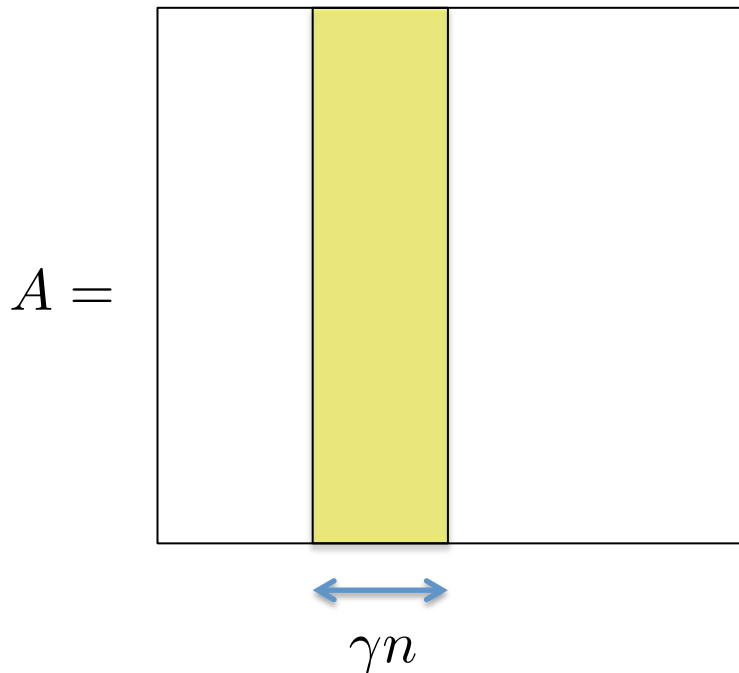


Offline algorithm: return the clusters after all the columns has been observed.

Online algorithm: classify the node immediately after its column is observed.

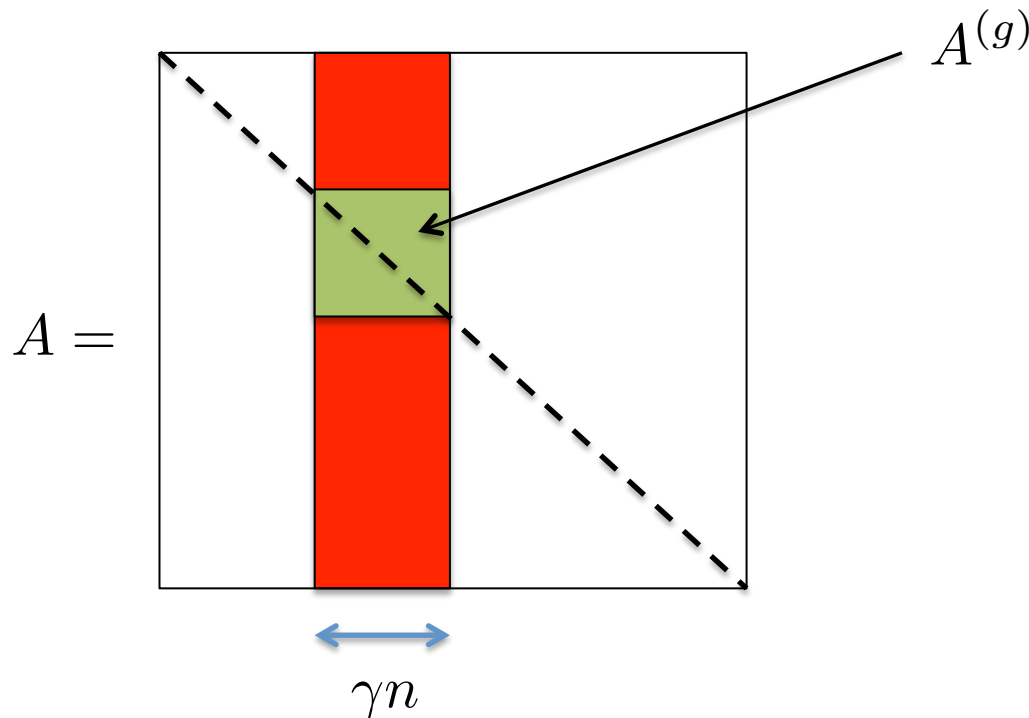
Classification with partial information

- Observe a fraction γ of columns
- Conditions on $\gamma, f(n)$ to classify the corresponding nodes or all nodes asymptotically accurately?



Classification with partial information

- Observe a fraction γ of columns
- Conditions on $\gamma, f(n)$ to classify the corresponding nodes or all nodes asymptotically accurately?



Fundamental limits

Theorem 7 Assume that $\sqrt{\gamma}f(n) = o(1)$. Then asymptotic detection is impossible.

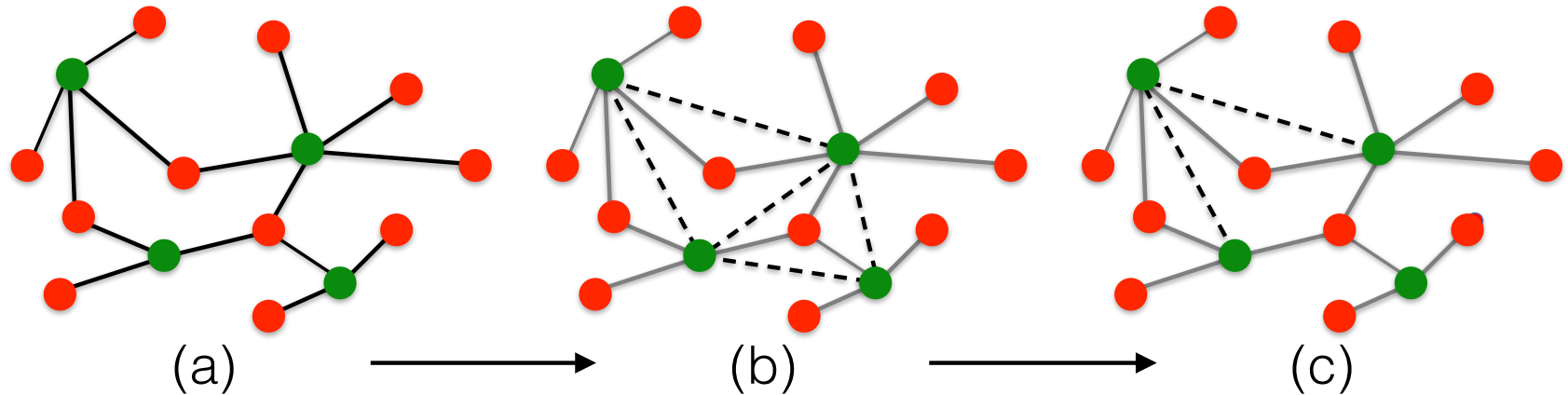
Theorem 8 (i) If there is an algorithm classifying green nodes asymptotically accurately, then $\sqrt{\gamma}f(n) = \omega(1)$.

(ii) If there is an algorithm classifying all nodes asymptotically accurately, then $\gamma f(n) = \omega(1)$.

Remark: if one uses information about green nodes only, then it is possible to classify these nodes only if $\gamma f(n) = \omega(1)$. We have to use side information provided by red nodes.

Algorithm for green nodes

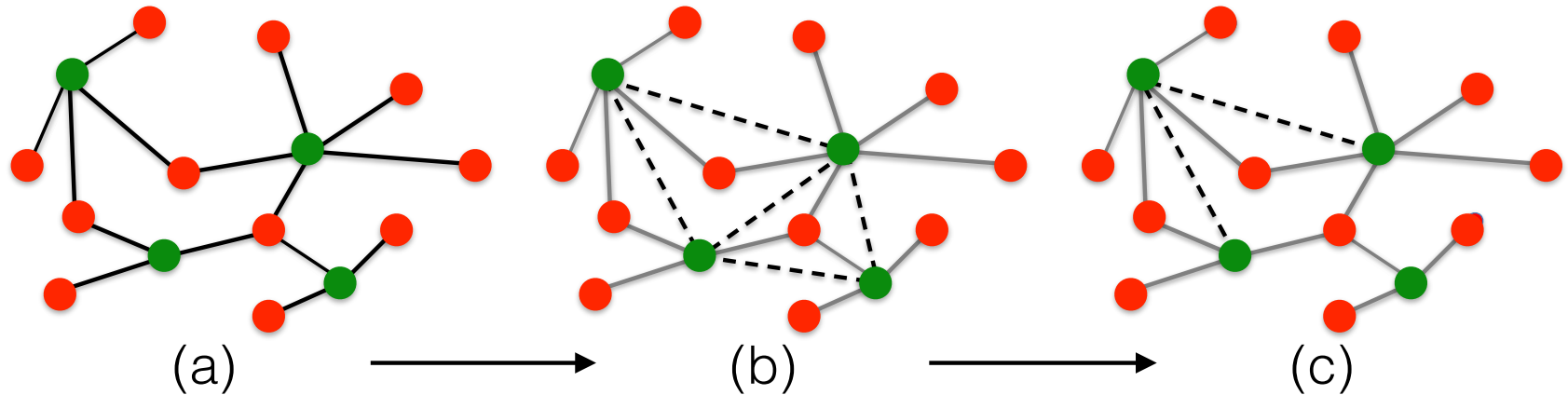
- Indirect edges through red nodes



Do it only for red nodes connected to exactly 2 green nodes
(avoid statistical dependence!)

Algorithm for green nodes

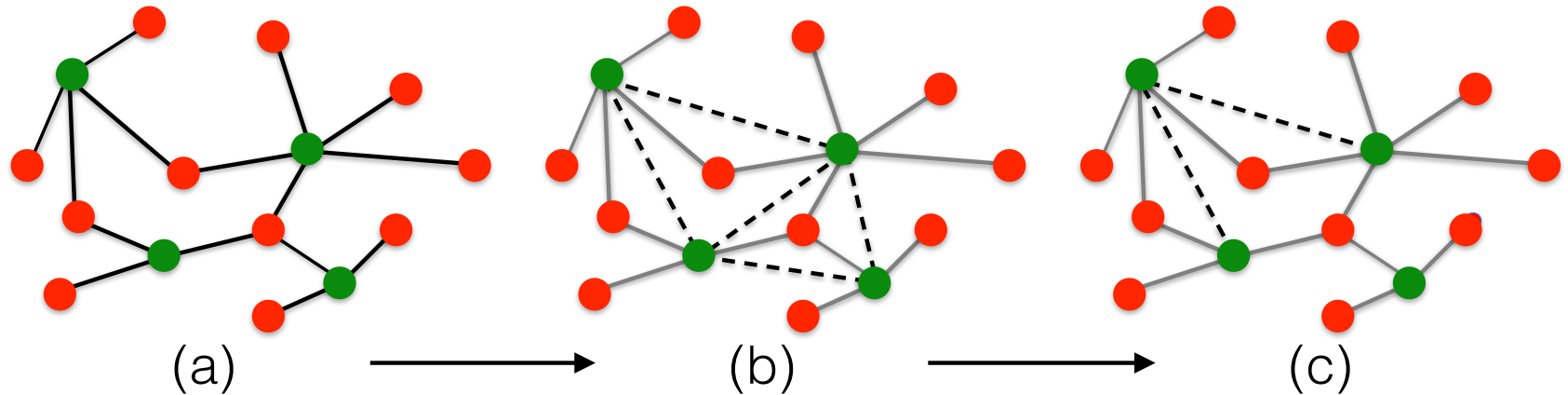
- Indirect edges through red nodes



- Result: a new adjacency matrix A'
- Algorithm: spectral method for $A^{(g)}$ and A' + keep the most informative matrix (with the highest normalized K -th eigenvalue)

Algorithm for green nodes

- Indirect edges through red nodes



Theorem 9 Assume that $\sqrt{\gamma}f(n) = \omega(1)$. The spectral algorithm exploiting indirect edges classifies the green nodes asymptotically accurately.

Algorithm for red nodes

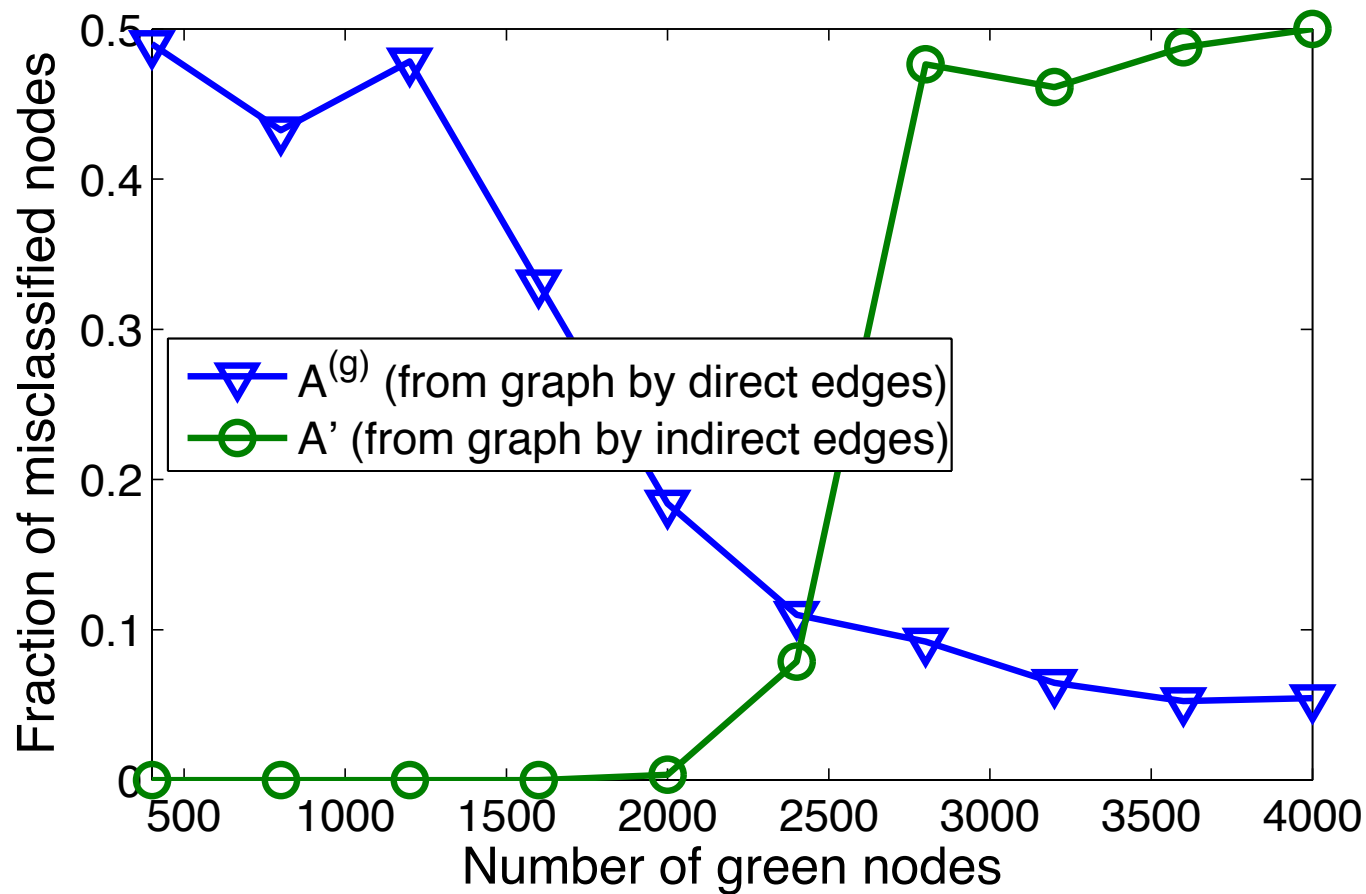
- Use the clusters of green nodes as kernels for the red nodes as in the adaptive sampling algorithm

Theorem 10 Assume that $\gamma f(n) = \omega(1)$. The above algorithm classifies the red nodes asymptotically accurately.

The algorithms are optimal (see fundamental limits)

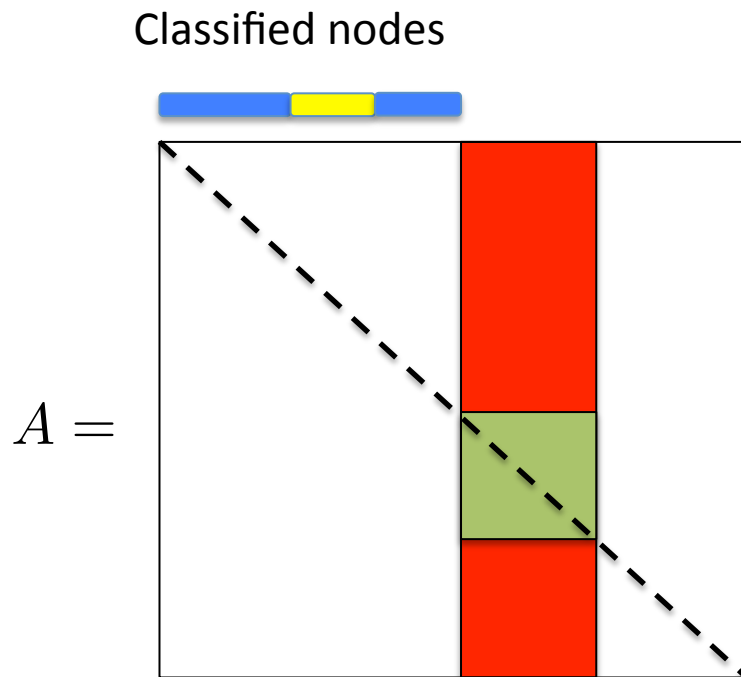
Example

$n = 1,000,000$ and $p = 0.005$, $q = 0.001$



Online memory-efficient algorithm

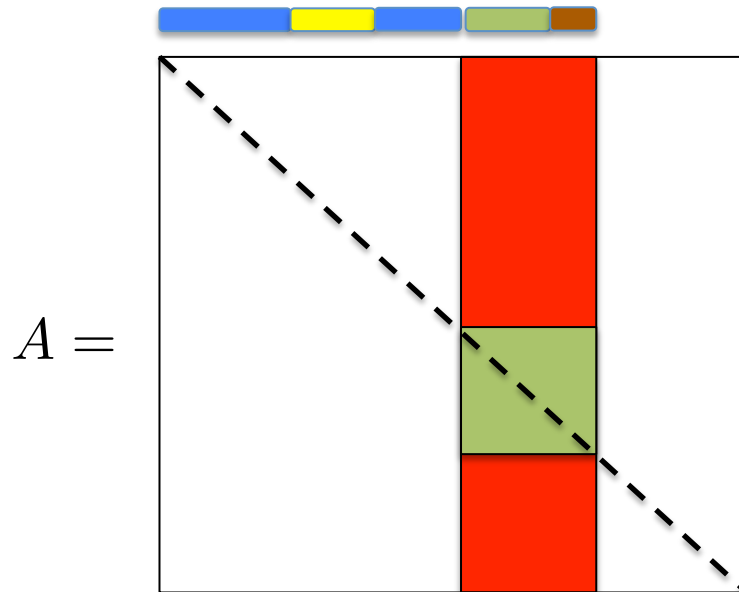
- Sequentially treat blocks of columns
 - Place the block of B columns in the memory



Online memory-efficient algorithm

- Sequentially treat blocks of columns
 - Place the block of B columns in the memory
 - Classify the corresponding nodes using previous algorithm

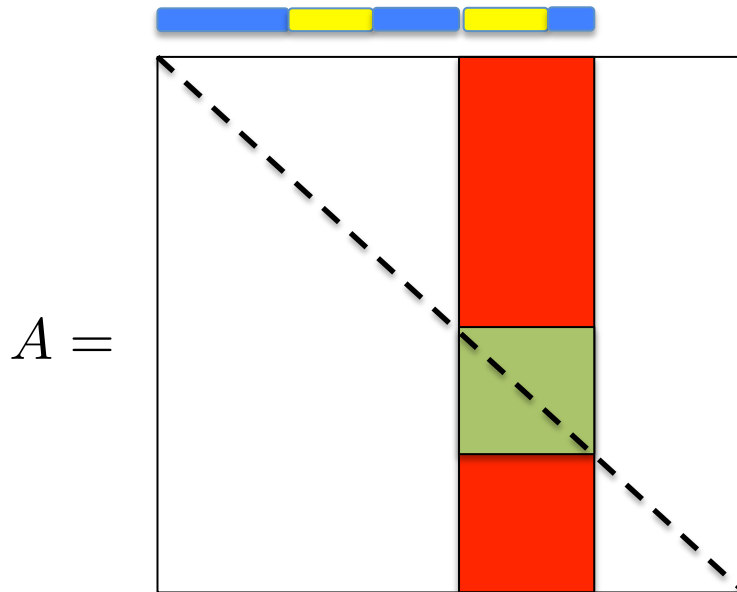
Classified nodes



Online memory-efficient algorithm

- Sequentially treat blocks of columns
 - Place the block of B columns in the memory
 - Classify the corresponding nodes using previous algorithm
 - Merge the obtained clusters with the clusters of the first block, and erase the block

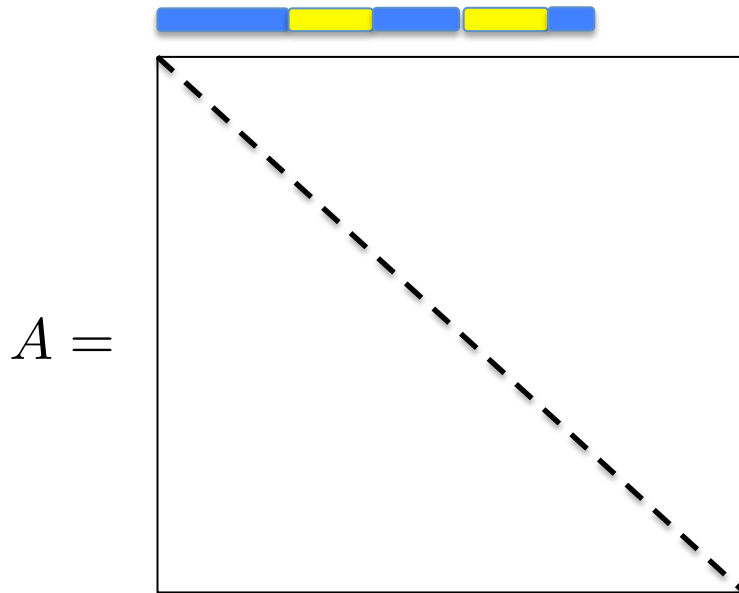
Classified nodes



Online memory-efficient algorithm

- Sequentially treat blocks of columns
 - Place the block of B columns in the memory
 - Classify the corresponding nodes using previous algorithm
 - Merge the obtained clusters with the clusters of the first block, and erase the block

Classified nodes



Online memory-efficient algorithm

- Block size: $B = \left(p\sqrt{np\log(n)}\right)^{-1}$
- Memory requirement: $O(Bnp\log(n))$ bits

Theorem 11 The previous algorithm detects communities asymptotically accurately in (less than) one pass, and using sublinear memory as soon as $p = \omega(\log(n)/n)$.

Conclusions

- A generic sampling framework extending the SBM
 - Necessary conditions for asymptotically accurate detection
 - Asymptotically optimal joint sampling and clustering algorithms
 - Our results hold in any regime! (Sparse or dense)
- Memory limited, streaming algorithm
 - Everything can be done with linear memory or even sublinear memory in (less than) one pass on the data
 - Memory-performance trade-off of our spectral approach
- Open questions
 1. Scaling of the proportion of misclassified nodes?
 2. Pareto boundary of the performance-memory trade-off?

Thanks!

Papers

Yun-Proutiere, COLT 2014

Yun-Lelarge-Proutiere, NIPS 2014