

Optimization and learning techniques for clustering problems

Soledad Villar

Center for Data Science
Courant Institute of Mathematical Sciences



NEW YORK UNIVERSITY

Statistical Physics and Machine Learning back together
Cargèse, August 2018

Clustering

“Next AI revolution is unsupervised” Yann LeCun

Main task in unsupervised machine learning:

- ▶ Finding structure in unlabeled data.

k-means clustering

- ▶ Simple objective
- ▶ Useful for “generic data”

Clustering stochastic block model

- ▶ Specialized objective
- ▶ Algorithms get more precise but less robust to model changes

Data-driven clustering methods.

Clustering

“Next AI revolution is unsupervised” Yann LeCun

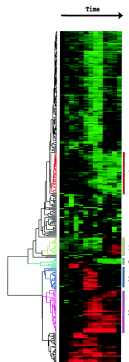
Main task in unsupervised machine learning:

- ▶ Finding structure in unlabeled data.

Example:

Cluster analysis and display of genome-wide expression patterns.

- ▶ Gene expression as a function of time for different conditions.
- ▶ Clustering gene expression data groups together genes of known similar functionality.



Summary

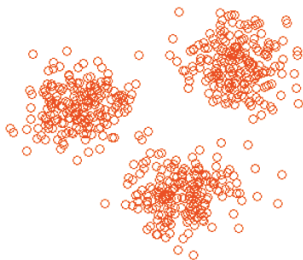
- ▶ k -means clustering
 - ▶ Landscape of manifold optimization
 - ▶ Lower bounds certificates from SDPs (data-driven)
- ▶ Clustering stochastic block model
 - ▶ AMP \longrightarrow Spectral methods \longrightarrow Graph neural networks
- ▶ Quadratic assignment

The k -means problem

Given a point cloud $\{x_i\}$,
partition the points in clusters
 C_1, \dots, C_k

k -means objective:

$$\min_{C_1, \dots, C_k} \sum_{t=1}^k \sum_{i \in C_t} \left\| x_i - \frac{1}{|C_t|} \sum_{j \in C_t} x_j \right\|^2$$



- NP-hard to minimize in general (even in the plane).

Lloyd's algorithm

1. Choose k centers at random.
2. Assign points to closest center.
3. Compute new centers are clusters centroids.

k -means++ is Lloyd's with smarter initialization.

Pros:

- ▶ Fast.
- ▶ Very easy to implement.
- ▶ Widely used and it works for most applications.

Cons:

- ▶ No guarantee of convergence (local minima).
- ▶ In extreme cases it may take exponentially many steps.
- ▶ Solutions depend heavily on initialization.
- ▶ Its output doesn't say how good of a solution it may be.

Optimization formulation

Taking $D_{ij} := \|x_i - x_j\|^2$, then

$$\sum_{t=1}^k \sum_{i \in C_t} \left\| x_i - \frac{1}{|C_t|} \sum_{j \in C_t} x_j \right\|^2 = \frac{1}{2} \text{Tr} \left(D \underbrace{\sum_{t=1}^k \frac{1}{|C_t|} \mathbf{1}_{C_t} \mathbf{1}_{C_t}^\top}_{\text{diagonal matrix}} \right)$$



=



$$\begin{aligned} &\text{minimize} && \text{Tr}(D \mathbf{Y} \mathbf{Y}^\top) \\ &\text{subject to} && \mathbf{Y} \in \mathbb{R}^{n \times k} \\ & && \mathbf{Y}^\top \mathbf{Y} = \mathbf{I}_k \\ & && \mathbf{Y} \mathbf{Y}^\top \mathbf{1} = \mathbf{1} \\ & && \mathbf{Y} \geq 0 \end{aligned}$$

Set of \mathbf{Y} 's is discrete.

Manifold optimization implementation

$$\begin{aligned} &\text{minimize} && \text{Tr}(DYY^\top) + \lambda \|Y_-\|^2 \\ &\text{subject to} && Y \in M = \{Y \in \mathbb{R}^{n \times k} : Y^\top Y = I_k, \quad YY^\top \mathbf{1} = \mathbf{1}\} \end{aligned}$$

To implement the manifold optimization one needs:

$$\begin{aligned} \text{grad}_f &: M \rightarrow TM, \\ \text{retr}_Y &: T_Y M \rightarrow M, \end{aligned}$$

where

$$\text{retr}_Y(0) = 0, \quad \left. \frac{d}{dt} \right|_{t=0} \text{retr}_Y(tV) = V.$$

Gradient descent:

$$Y_{n+1} = \text{retr}_{Y_n}(-\alpha_n \text{grad}_f(Y_n)).$$

Manifold optimization, retraction implementation

Homogeneous structure of M .

Let $O(1_n)$ be the $n \times n$ orthogonal matrices that fix 1.

$$\begin{aligned} M \times O(1_n) \times O(k) &\rightarrow M \\ (Y, Q, R) &\mapsto QYR \quad (\text{transitive action}) \end{aligned}$$

Multiplication of Y on the right by R is equivalent to multiplication on the left by $R' = YRY^\top \in O(1_n)$.

$V \in T_Y M$ can be written as $V = BY + YA$ where $A \in \mathfrak{so}(k)$, $B \in \mathfrak{so}(1_n)$ where BY and YA are orthogonal. We compute A, B explicitly.

$$\begin{aligned} A &= Y^\top V \\ B &= VY^\top - YV^\top - 2YAY^\top \end{aligned}$$

$$\text{retr}_Y(V) = \exp(B) \exp(A') Y \in M$$

Manifold optimization algorithm

$$\lambda_0 \leftarrow 0$$

repeat

$$Y_{n+1} \leftarrow \text{GradientDescent}(f_\lambda) \text{ \{Initialized at } Y_n\}$$

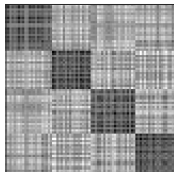
$$\lambda_{n+1} \leftarrow 2\lambda_n + 1$$

until $\|Y_{-}\|_F < \epsilon$

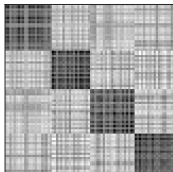
$$\text{minimize}_Y \quad f_\lambda(Y) = \text{Tr}(DYY^\top) + \lambda \|Y_{-}\|^2$$

$$\text{subject to} \quad Y \in M = \{Y \in \mathbb{R}^{n \times k} : Y^\top Y = I_k, \quad YY^\top \mathbf{1} = \mathbf{1}\}$$

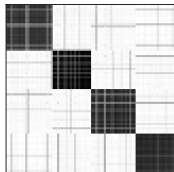
lambda=0



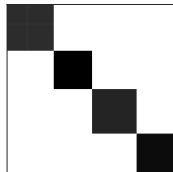
lambda=6300



lambda=2.047e+05



lambda=2.6214e+07

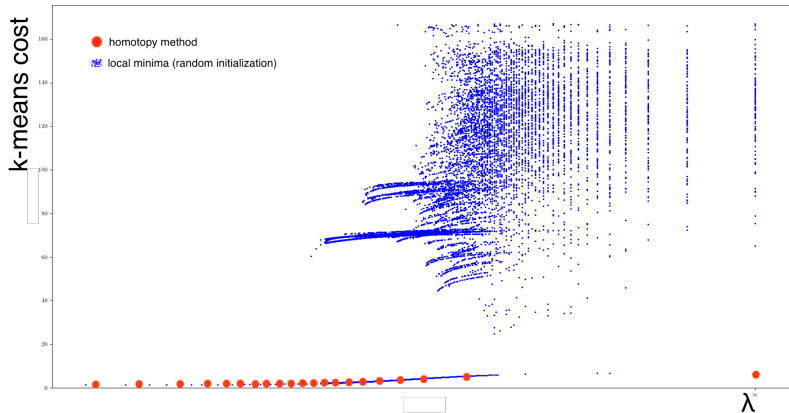


Problem structure allows to do each iteration in linear-time.

Claim: clustering is (maybe) not (that) hard

Optimization landscape

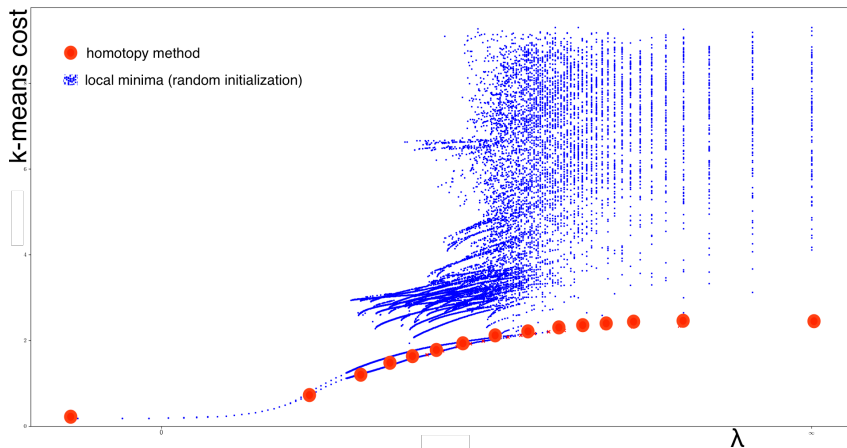
Points from three well-separated unit balls in small dimension.



Claim: clustering is (maybe) not (that) hard

No clustering structure

Points drawn from unit ball into three clusters in small dimension.



And it's optimal! How do I know that?

k-means optimization formulation

Recall $D_{ij} := \|x_i - x_j\|^2$, then

$$\sum_{t=1}^k \sum_{i \in C_t} \left\| x_i - \frac{1}{|C_t|} \sum_{j \in C_t} x_j \right\|^2 = \frac{1}{2} \text{Tr} \left(D \underbrace{\sum_{t=1}^k \frac{1}{|C_t|} \mathbf{1}_{C_t} \mathbf{1}_{C_t}^\top}_{\text{Diagram 1}} \right)$$



=



$$\begin{aligned} &\text{minimize} && \text{Tr}(D \mathbf{Y} \mathbf{Y}^\top) \\ &\text{subject to} && \mathbf{Y} \in \mathbb{R}^{n \times k} \\ & && \mathbf{Y}^\top \mathbf{Y} = \mathbf{I}_k \\ & && \mathbf{Y} \mathbf{Y}^\top \mathbf{1} = \mathbf{1} \\ & && \mathbf{Y} \geq 0 \end{aligned}$$

Set of \mathbf{Y} 's is discrete.

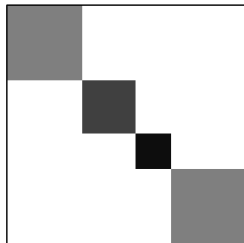
A semidefinite programming relaxation

Recall $D_{ij} := \|x_i - x_j\|^2$, then

$$\sum_{t=1}^k \sum_{i \in C_t} \left\| x_i - \frac{1}{|C_t|} \sum_{j \in C_t} x_j \right\|^2 = \frac{1}{2} \text{Tr} \left(D \underbrace{\sum_{t=1}^k \frac{1}{|C_t|} \mathbf{1}_{C_t} \mathbf{1}_{C_t}^\top}_{\text{block diagonal matrix}} \right)$$

Relax to SDP:

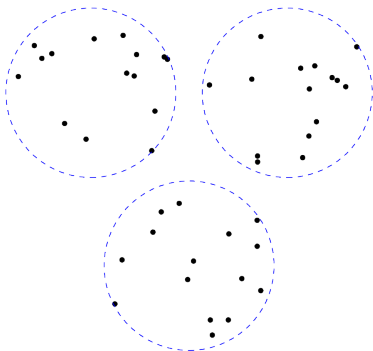
$$\begin{array}{ll} \text{minimize} & \text{Tr}(D\mathbf{X}) \\ \text{subject to} & \text{Tr}(\mathbf{X}) = k \\ & \mathbf{X}\mathbf{1} = \mathbf{1} \\ & \mathbf{X} \succeq 0 \\ & \mathbf{X} \preceq \mathbf{I} \end{array}$$



Stochastic ball model

(\mathcal{D}, γ, n) -stochastic ball model

- ▶ \mathcal{D} = rotation-invariant distribution over unit ball in \mathbb{R}^m
- ▶ $\gamma_1, \dots, \gamma_k$ = ball centers in \mathbb{R}^m
- ▶ Draw $r_{t,1}, \dots, r_{t,n}$ i.i.d. from \mathcal{D} for each $i \in \{1, \dots, k\}$
- ▶ $x_{t,i} = \gamma_t + r_{t,i}$ = i th point from cluster t



Stochastic ball model

(\mathcal{D}, γ, n) -stochastic ball model

- ▶ \mathcal{D} = rotation-invariant distribution over unit ball in \mathbb{R}^m
- ▶ $\gamma_1, \dots, \gamma_k$ = ball centers in \mathbb{R}^m
- ▶ Draw $r_{t,1}, \dots, r_{t,n}$ i.i.d. from \mathcal{D} for each $i \in \{1, \dots, k\}$
- ▶ $x_{t,i} = \gamma_t + r_{t,i}$ = i th point from cluster t



k-means SDP. *Stochastic ball model*

Exact recovery

SDP is tight for k unit balls in \mathbb{R}^m with centers $\gamma_1 \dots \gamma_k$

$$\min_{i \neq j} \|\gamma_i - \gamma_j\| \geq \min \left\{ 2\sqrt{2}\left(1 + \frac{1}{\sqrt{m}}\right), 2 + \frac{k^2}{m}, 2 + \sqrt{\frac{2k}{m+2}} \right\}$$

Conjecture (Li, Li, Ling, Strohmer)

Phase transition for exact recovery (under stochastic ball model) at

$$\min_{i \neq j} \|\gamma_i - \gamma_j\| \geq 2 + \frac{c}{m}$$

Awasthi, Bandeira, Charikar, Krishnaswamy, V., Ward, Proc. ITCS, 2015

Iguchi, Mixon, Peterson, V., Mathematical Programming, 2016

Li, Li, Ling, Strohmer, arXiv:1710.06008 2017

k-means SDP. *Subgaussian mixtures*

Relax and “round” (weak recovery).

Theorem

Centers: $\hat{\gamma}_1, \dots, \hat{\gamma}_k$. Estimated centers v_1, \dots, v_k

$$\frac{1}{k} \sum_{i=1}^k \|v_i - \hat{\gamma}_i\|^2 \lesssim k^2 \sigma^2 \text{ whp provided } \min_{i \neq j} \|\gamma_i - \gamma_j\| \gtrsim k\sigma.$$

Example: <http://solevillar.github.io/2016/07/05/Clustering-MNIST-SDP.html>

SDPs are slow!

Fast lower bound certificates

Want: Lower bound on k -means value,

Given $\{x_i\}_{i \in T} \subseteq \mathbb{R}^m$, let W = value of k -means++ initialization

$$\text{val}_{k\text{means}}(T) := \min \frac{1}{|T|} \sum_{t=1}^k \sum_{i \in C_t} \left\| x_i - \frac{1}{|C_t|} \sum_{j \in C_t} x_j \right\|^2 \geq \frac{1}{8(\log k + 2)} \mathbb{E} W$$

Running k -means++ on MNIST training set with $k = 10$ gives:

- ▶ Upper bound (by running the algorithm): 39.22
- ▶ Lower bound (estimated from above): 2.15

Can we get a better lower bound?

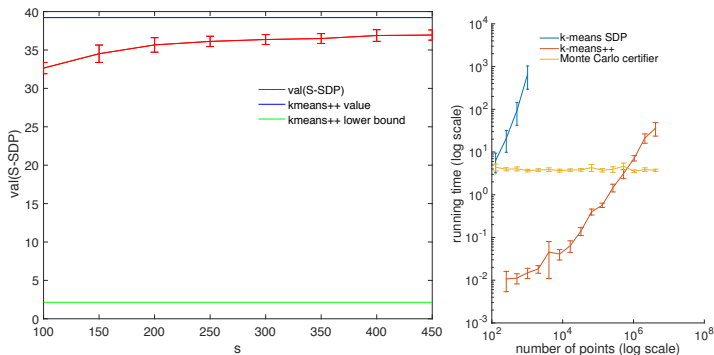
Idea

Given $\{x_i\}_{i \in T} \subseteq \mathbb{R}^m$, draw $S \sim \text{Unif} \in \binom{T}{s}$.

$$\begin{aligned}\mathbb{E} \text{val}_{\text{SDP}}(S) &\leq \mathbb{E} \text{val}_{\text{kmeans}}(S) \\ &\leq \mathbb{E} \left[\frac{1}{s} \sum_{t=1}^k \sum_{i \in C_t^* \cap S} \left\| x_i - \frac{1}{|C_t^* \cap S|} \sum_{j \in C_t^* \cap S} x_j \right\|^2 \right] \\ &\leq \mathbb{E} \left[\frac{1}{s} \sum_{t=1}^k \sum_{i \in C_t^* \cap S} \left\| x_i - \frac{1}{|C_t^*|} \sum_{j \in C_t^*} x_j \right\|^2 \right] = \text{val}_{\text{kmeans}}(T)\end{aligned}$$

Goal: Establish $\mathbb{E} \text{val}_{\text{SDP}}(S) > B$ with small p -value.

Better performance, and fast!



- ▶ Taking $s = 100 \ll 60,000$ MNIST points already works
- ▶ Constant runtime in N , faster than k -means++ for $N \gtrsim 10^6$.
- ▶ Theorem for mixtures of Gaussians (additive bound)

Generic algorithms vs. model-dependent algorithms

Clustering the stochastic block model

- ▶ Sparse graph clustering (non-geometric)
- ▶ Similar generic algorithmic (SDP, manifold optimization)
- ▶ Model-tailored algorithms (AMP, spectral methods on Non-backtracking matrices and Bethe Hessian)
- ▶ Can we learn the algorithm from data?
 - ▶ Use graph neural networks

Decelle, Krzakala, Moore, Zdeborová, 2013

Krzakala, Moore, Mossel, Neeman, Sly, Zdeborová, Zhang, 2013

Saade, Krzakala, Zdeborová, 2014

Abbe, 2018

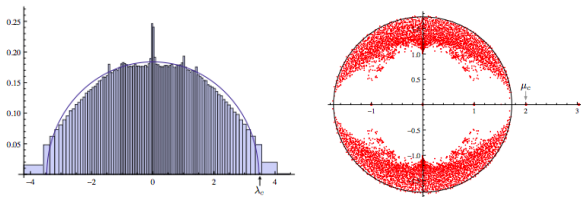
Spectral redemption

$A \sim SBM(a/n, b/n, n, 2)$ sparse.

Spectrum doesn't concentrate (high degree vertices dominate it)
Laplacian is not useful for clustering

Consider the non-backtracking operator (from linearized BP)

$$B_{(i \rightarrow j)(i' \rightarrow j')} = \begin{cases} 1 & \text{if } j = i' \text{ and } j' \neq i \\ 0 & \text{otherwise} \end{cases}$$



Second eigenvector of B reveals clustering structure

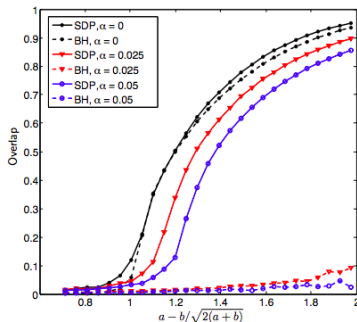
Bethe Hessian

$$BH(r) = (r^2 - 1)I - rA + D$$

Fixed points of BP \longleftrightarrow Stationary points of Bethe free energy

Second eigenvector reveals clustering structure

Pitfall: highly dependent in the model



Goal: Combine graph operators I, D, A, \dots to generate robust “data-driven spectral methods” for problems in graphs

Graph neural networks

Power method: $v^{t+1} = Mv^t \quad t = 1, \dots, T.$

Graph with adjacency matrix A . Set $\mathcal{M} = \{I_n, D, A, A^2, \dots, A^{2^J}\},$

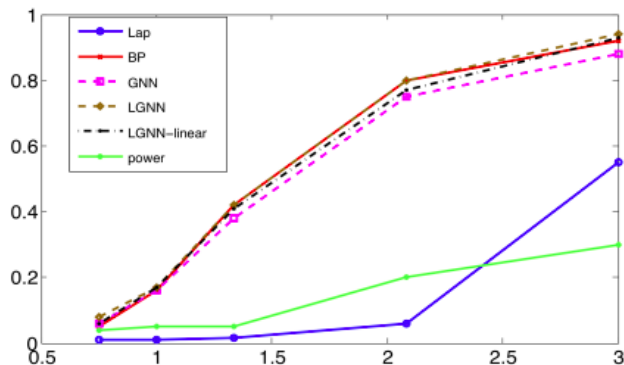
$$v_l^{t+1} = \rho \left(\sum_{M \in \mathcal{M}} M v^t \theta_{M,l}^t \right), \quad l = 1 \dots d_{t+1}$$

with $v^t \in \mathbb{R}^{n \times d_t},$

$\Theta = \{\theta_1^t, \dots, \theta_{|\mathcal{M}|}^t\}_t, \theta_M^t \in \mathbb{R}^{d_t \times d_{t+1}}$ trainable parameters.

- ▶ Extension to line graph (GNN with non-backtracking).
- ▶ Covariant wrt permutations $G \mapsto \phi(G)$ then $G_\Pi \mapsto \Pi \phi(G).$
- ▶ Preliminary theoretical analysis of energy landscape.
 - ▶ Strong assumptions \Rightarrow local minima have low energy.

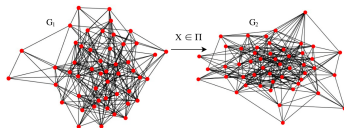
Numerical performance. SBM $k = 2$



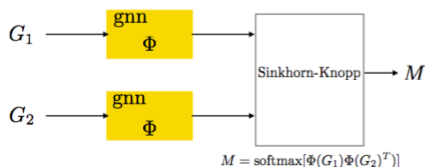
Overlap as function of SNR

Extension to quadratic assignment

Graph matching: $\min_{X \in \Pi} \|G_1 X - X G_2\|_F^2 = \|G_1 X\|^2 + \|X G_2\|^2 - 2\langle G_1 X, X G_2 \rangle$



Siamese neural network:

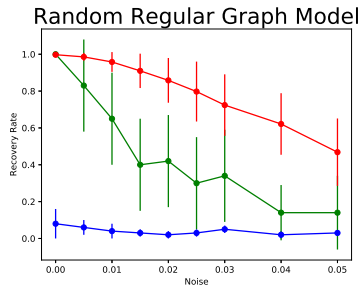
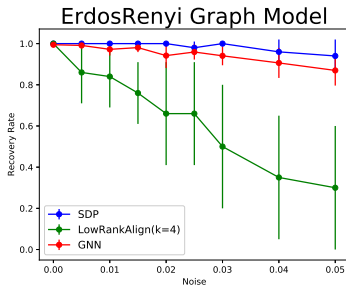


$$G_2 = \pi G_1 + N \quad N \sim \text{Erdos-Renyi}$$

$$G_1 \sim \text{Erdos-Renyi}$$

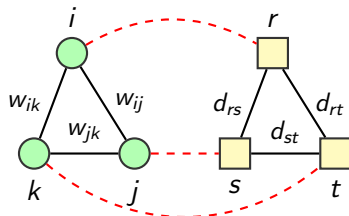
$$G_1 \sim \text{Random regular}$$

Numerical experiment



Our model runs in $O(n^2)$, LowRankAlign is $O(n^3)$, SDP in $O(n^4)$.

Quadratic Assignment Problem (QAP)

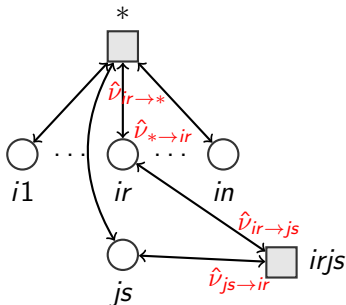


$$\begin{aligned} \min_{\pi} \quad & \sum_{i \in A} \sum_{j \in A} \sum_{r \in B} \sum_{s \in B} \underbrace{\sum_{E_{ijrs}} w_{ij} d_{rs}}_{E_{ijrs}} \pi_{ir} \pi_{js} \\ \text{subject to} \quad & \sum_r \pi_{ir} = 1 \quad \forall i \in A \\ & \sum_i \pi_{ir} = 1 \quad \forall r \in B \\ & \pi_{ir} \in \{0, 1\} \end{aligned}$$

Max-Product Belief Propagation for QAP

- Posterior probability:

$$p(\pi) = \frac{1}{Z} \prod_i \mathbb{1}\{\sum_t \pi_{it} = 1\} \prod_r \mathbb{1}\{\sum_k \pi_{kr} = 1\} \prod_{j \neq i} \prod_{s \neq r} e^{-\beta \pi_{ir} \pi_{js} E_{ijrs}}$$



Max-Product Message Updates

- ▶ Variable to check nodes:

$$\nu_{ir \rightarrow *}(1) \cong \prod_{kt} e^{-\beta \hat{\nu}_{kt \rightarrow ir}(1) E_{irkt}}$$

$$\nu_{ir \rightarrow *}(0) \cong 1$$

- ▶ Check to variable nodes:

$$\nu_{* \rightarrow ir}(1) \cong \nu_{ir \rightarrow *}(1) \prod_{\ell \neq i} \nu_{\ell r \rightarrow *}(0) \prod_{t \neq r} \nu_{it \rightarrow *}(0)$$

$$\nu_{* \rightarrow ir}(0) \cong \nu_{ir \rightarrow *}(0) \max_{t \neq r, \ell \neq i} \nu_{it \rightarrow *}(1) \nu_{\ell r \rightarrow *}(1) \prod_{u \neq t, r} \nu_{iu \rightarrow *}(0) \prod_{m \neq \ell, i} \nu_{mr \rightarrow *}(0)$$

- ▶ Variable back to variable nodes:

$$\hat{\nu}_{ir \rightarrow js}(1) \cong \nu_{* \rightarrow ir}(1) \prod_{ku \neq js} e^{-\beta \hat{\nu}_{ku \rightarrow ir}(1) E_{irku}}$$

$$\hat{\nu}_{ir \rightarrow js}(0) \cong \nu_{* \rightarrow ir}(0)$$

Simplifications (Min-Sum Algorithm)

Introduce log-likelihood ratio:

$$\hat{\ell}_{ir \rightarrow js} = \frac{1}{\beta} \log \frac{\hat{\nu}_{ir \rightarrow js}(1)}{\hat{\nu}_{ir \rightarrow js}(0)}$$

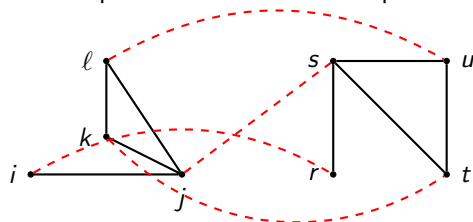
Then the algorithm simplifies to one update step as follows,

$$\begin{aligned} \hat{\ell}_{ir \rightarrow js} \leftarrow \min_{t \neq r, l \neq i} \sum_{ku} & \frac{e^{\beta \hat{\ell}_{ku \rightarrow it}}}{1 + e^{\beta \hat{\ell}_{ku \rightarrow it}}} E_{itku} + \frac{e^{\beta \hat{\ell}_{ku \rightarrow \ell r}}}{1 + e^{\beta \hat{\ell}_{ku \rightarrow \ell r}}} E_{\ell rku} \\ & - 2 \sum_{ku \neq js} \frac{e^{\beta \hat{\ell}_{ku \rightarrow ir}}}{1 + e^{\beta \hat{\ell}_{ku \rightarrow ir}}} E_{irku} + \frac{e^{\beta \hat{\ell}_{js \rightarrow ir}}}{1 + e^{\beta \hat{\ell}_{js \rightarrow ir}}} E_{irjs} \end{aligned}$$

A special case: Graph Matching (Isomorphism)

Graph-A

Graph-B



When $A \sim \text{ER}(n, 0.5)$ and

$$B = PAP^T \oplus Z$$

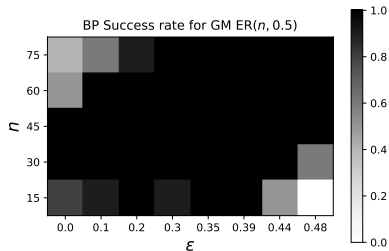
for some $P \in \Pi$ where $Z \sim \text{Ber}(\epsilon)$. Find P given unlabeled graphs.

$$\min_P \|A - PBP^T\|_F$$

subject to $P \in \Pi$

equivalent to QAP with

$$E_{ijrs} = A_{ij} \oplus B_{rs}.$$



Summary

- ▶ k -means clustering
 - ▶ Nice manifold optimization landscape
 - ▶ Lower bounds certificates from SDPs (data-driven)
- ▶ Clustering stochastic block model
 - ▶ AMP \longrightarrow Spectral methods \longrightarrow Graph neural networks
- ▶ Quadratic assignment

Thank you

Relax no need to round: Integrality of clustering formulations

P. Awasthi, A. S. Bandeira, M. Charikar, R. Krishnaswamy, S. Villar, R. Ward
arXiv:1408.4045, Innovations in Theoretical Computer Science (ITCS) 2015

Probably certifiably correct k -means clustering

T. Iguchi, D. G. Mixon, J. Peterson, S. Villar
arXiv:1509.07983, Mathematical Programming, 2016

Clustering subgaussian mixtures by semidefinite programming

D. G. Mixon, S. Villar, R. Ward
arXiv:1602.06612, Information and Inference: A Journal of the IMA, 2017

Monte Carlo approximation certificates for k -means clustering

D. G. Mixon, S. Villar
arXiv:1710.00956

Supervised Community Detection with Hierarchical Graph Neural Networks

Z. Chen, L. Li, J. Bruna
arXiv:1705.08415

A note on learning algorithms for quadratic assignment with Graph Neural Networks

A. Nowak, S. Villar, A. S. Bandeira, J. Bruna, ICML (PADL) 2017
arXiv:1706.07450