

Variational Inference of Community Models using Belief Propagation

A Case Study: Model Selection for Stochastic Block Models

Xiaoran Yan

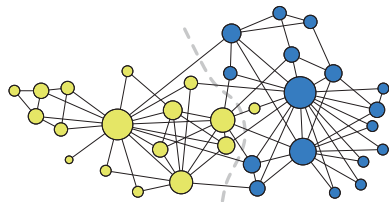
June 4, 2013

Joint work with Cris Moore, Cosma Shalizi, Aurelien Decelle, Lenka Zdeborová,
Florent Krzakala, Pan Zhang, Yaojia Zhu

Community structures in networks

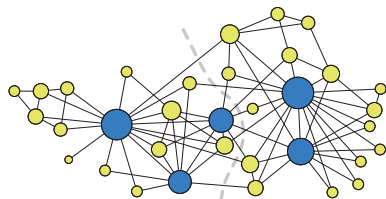
Assortative communities

- Dense connection within groups
- Sparse connection between groups
- Min cut, spectral partitioning, modularity, etc.



Functional communities

- Structure equivalence
- Disassortative communities
- Mixed structures, satellite structures
- Food webs, leaders and followers

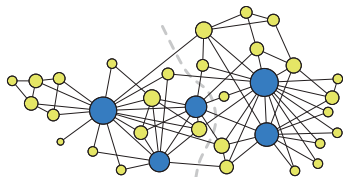


Assumptions:

- We represent our network as an undirected simple graph $G = (V, E)$ with the adjacency matrix A .
- Each node $u \in V$ has a hidden block label $g(u) \in \{1, \dots, k\}$.
- Each node's block label $g(u)$ is first generated according to $q_{g(u)}$. Let n_s be the number of nodes of type s , with $n = \sum_s n_s$.
- Between each pair of nodes $\{u, v\}$, an edge is generated independently with probability $p_{g(u)g(v)}$. Let m_{st} be the number of edges from block s to block t , with $\sum_{st} m_{st} = m$.

Given the parameters p_{st} and a block assignment, i.e., a function $g : V \rightarrow \{1, \dots, k\}$ assigning a label to each node, the likelihood of generating a given graph G in this model is:

$$\begin{aligned} & P(G, g | q, p) \\ &= \prod_u q_{g(u)} \prod_{u < v} (p_{g(u)g(v)})^{A_{uv}} (1 - p_{g(u)g(v)})^{1 - A_{uv}} \\ &= \prod_{s=1}^k q_s^{n_s} \prod_{s,t=1}^k (p_{st})^{m_{st}} (1 - p_{st})^{n_s n_t - m_{st}}. \end{aligned}$$



Summing over latent block assignment

- Overcoming discreteness
- A natural requirement in Bayesian approaches
- Avoid over-fitting by averaging over of latent states of different fit
- k^n number of terms in the sum

$$P(G | q, p) = \sum_g P(G, g | q, p),$$

In statistical physics terms

- Assuming temperature $T = 1/\beta = 1$
- State energy $E(g) = -\log P(G, g | q, p)$ —likelihood
- Ground state $\arg \min_g E(g)$ —Maximum likelihood state
- The Boltzmann distribution $P^*(g) = \frac{e^{-E(g)}}{\sum_g e^{-E(g)}}$
- The free energy $\sum_g e^{-E(g)} = -\log P(G | q, p)$ —the total likelihood
- Calorimetry Tricks for estimating the free energy
 - simulated annealing, population annealing, parallel tempering

The variational trick

$$\begin{aligned}\log P(G | q, p) &= \log \sum_g Q(g) \frac{P(G, g | q, p)}{Q(g)} \\ &\geq \sum_g \log \left[Q(g) \frac{P(G, g | q, p)}{Q(g)} \right] \\ &= \mathbb{E}_{Q(g)} \left[\log \frac{P(G, g | q, p)}{Q(g)} \right] \\ &= \mathbb{E}_{Q(g)} [\log P(G, g | q, p)] + \mathbb{S}[Q(g)] \\ &= -\langle \mathbb{E} \rangle + \mathbb{S}\end{aligned}$$

The variational distribution

- If $Q(g) = P^*(g)$ the approximation is exact
- Use simpler forms of $Q(g)$ the approximation becomes a optimization

$$\log P(G | q, p) \approx \sup_Q \left[\mathbb{E}_{Q(g)} [\log P(G, g | q, p)] + \mathbb{S}[Q(g)] \right]$$

The choice of $Q(g)$

$$P(G, g | q, p) = \prod_u q_{g(u)} \prod_{u < v} (p_{g(u)g(v)})^{A_{uv}} (1 - p_{g(u)g(v)})^{1 - A_{uv}}.$$

$$\log P(G | q, p) \geq \sum_g \log \left[Q(g) \frac{P(G, g | q, p)}{Q(g)} \right]$$

Wish list

- Be able to factor $Q(g)$ into local terms
- Each individual local factor can be efficiently solved
- Can be optimized to achieve good approximation

The mean field free energy

$$Q(g) = \prod_u b_{g_u}^u$$

- Easy to solve
- Total independence
- Poor approximation for almost any graphs

Estimating the average energy

$$\begin{aligned} -\langle \mathbb{E} \rangle &= \mathbb{E}_{Q(g)}[\log P(G, g \mid q, p)] \\ &= \mathbb{E}_{Q(g)}\left[\sum_u \log q_{g(u)} + \sum_{(u,v) \in E} \log p_{g(u)g(v)} + \sum_{(u,v) \notin E} \log(1 - p_{g(u)g(v)})\right] \\ &= \sum_u \sum_s b_s^u \log q_s + \sum_{(u,v) \in E} \sum_{st} b_{st}^{uv} \log p_{st} + \sum_{(u,v) \notin E} \sum_{st} b_{st}^{uv} \log(1 - p_{st}) \end{aligned}$$

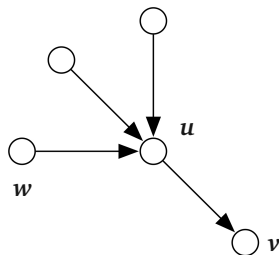
The cluster variational distribution

$$Q(g) = \frac{\prod_{u < v} b_{g_u g_v}^{uv}}{\prod_u (b_{g_u}^u)^{d_u - 1}}, \quad \text{with } b_{g_u}^u = \sum_{t, \forall v} b_{g_u t}^{uv}$$

- Exact for trees
- Empirical results show that it works pretty well for loopy graphs
- Conditional independence
- Corresponds to the 2nd order Kikuchi free energy
- Belief Propagation leads to the same fixed points (Yedidia, 2001)

The message passing algorithm

- Equivalent to the cavity method in statistical physics
- Each node u send a message to each neighbor v about u 's 1-point marginal
- The message $b_s^{u \rightarrow v}$ is based on all the other neighbors of u , as if v is absent
- Pass the messages around until convergence



Updating the messages

$$b_s^{u \rightarrow v} = \frac{1}{Z^{u \rightarrow v}} q_s \prod_{w \neq u, v} \sum_{t=1}^k b_t^{w \rightarrow u} (p_{st})^{A_{uw}} (1 - p_{st})^{1 - A_{uw}}$$

- Since we take non-edges into account, there are $O(n^2)$ messages to update
- In sparse networks, the messages along non-edges is of order $O(1/n)$, if we apply a mean field approximation for $\forall v, (u, v) \notin E$

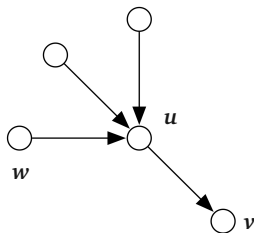
$$b_s^{u \rightarrow v} = b_s^u = \frac{1}{Z^u} q_s \prod_{w \neq u} \sum_{t=1}^k b_t^{w \rightarrow u} (p_{st})^{A_{uw}} (1 - p_{st})^{1 - A_{uw}}$$

Belief Propagation (continued)

The messages with a mean field on non-edges

$$b_s^{u \rightarrow v} = \frac{1}{Z^{u \rightarrow v}} q_s \prod_{\substack{(w,u) \in E \\ w \neq u,v}} \sum_{t=1}^k b_t^{w \rightarrow u} p_{st} \prod_{\substack{(w,u) \notin E \\ w \neq u,v}} \sum_{t=1}^k b_s^u (1 - p_{st}) \text{ if } (u,v) \in E$$

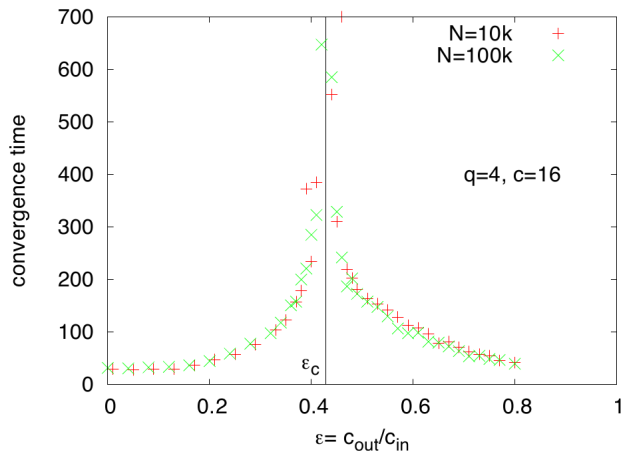
$$b_s^u = \frac{1}{Z^u} q_s \prod_{\substack{(w,u) \in E \\ w \neq u}} \sum_{t=1}^k b_t^{w \rightarrow u} p_{st} \prod_{\substack{(w,u) \notin E \\ w \neq u}} \sum_{t=1}^k b_s^u (1 - p_{st})$$



Running time analysis

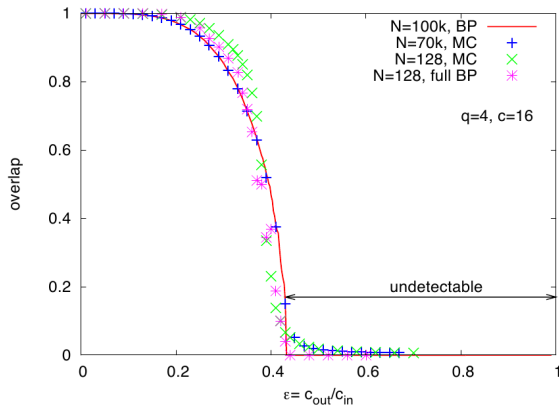
- $O(m + n)$ messages to update each sweep
- The messages converge within constant sweeps
- Linear time E-step
- The EM outer loop converges even faster
- Constant number of initial restarts
- $O(m + n)$ linear total time

Convergence time: finite correlation length



[Decelle, Krzakala, Moore, Zdeborová, PRL 2011]

Optimal detectability of Belief Propagation



Detectability

- BP fails when the p_{st} are too similar
- MCMC with extended running time also fails
- The same bound from spectral method in dense graphs
- Better bound than spectral method in sparse graphs
- Theoretical proof that BP is optimal in trees (Mossel, Neeman and Sly, 2001)
- No algorithm can do better!

The variational expectation maximization framework

Variational E-step

$$\log P(G | q, p) \approx \sup_Q [\mathbb{E}_{Q(g)}[\log P(G, g | q, p)] + \mathbb{S}[Q(g)]]$$

- Choose your favorite $Q(g)$
- Optimize $Q(g)$ while fixing the parameters q, p
- For SBMs: Bethe approximation with linear Belief Propagation

M-step

$$\hat{q}_s = \frac{\bar{n}_s}{n} = \frac{\sum_u b_s^u}{n}, \quad \hat{\omega}_{st} = \frac{\bar{m}_{st}}{n_s n_t} = \frac{\sum_{u \neq v} A_{uv} b_{st}^{uv}}{(\sum_u b_s^u)(\sum_u b_t^u)},$$

- Solving for the MLEs of the parameters q, p while fixing $Q(g)$
- Go back to E-step, rinse and repeat until a fixed point is reached
- Gradient ascent in the joint space, maximizing the total likelihood
- Linear approximation for an exponential size problem

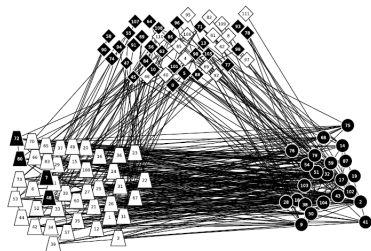
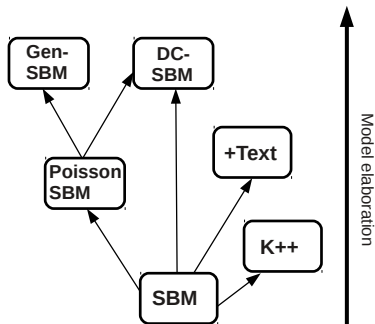
A case study: building the right stochastic block model

Variants and elaborations

- Degree corrected SBM
- Extensions for rich data
- More blocks!

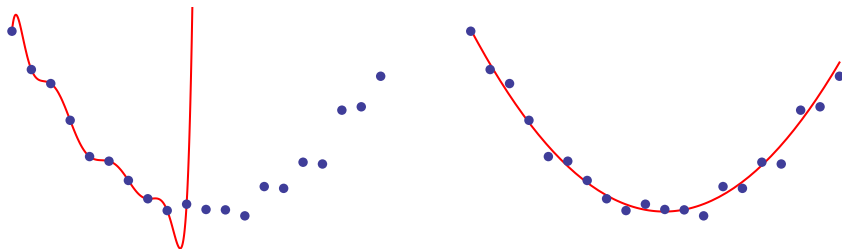
Model selection

- Which model to choose given the data?
- Number of blocks (order selection)?
- Over-fitting?



Occam's razor

- Complex models with more parameters have a natural advantage at fitting data.
- Simpler models have lower variability, thus less sensitive to noise in the data.
- Balance the trade-off between bias and variance.
- Excessive complexity not only increases the cost of the model, but also hurts the generalization performance.



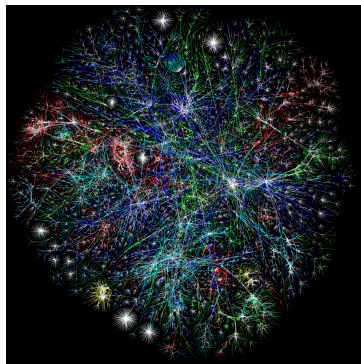
Model selection for stochastic block models

Common approaches

- Use the model you like.
- Make a choice based on domain expertise.
- Use off-the-shelf Information Criteria for independent data.
 - Akaike information criterion (AIC), Bayesian information criterion (BIC), etc.
- Non-parametric methods

Generalization test (cross validation)

- Node classification
- Link prediction
- The good
 - can compare any model
 - generalization performance focused
- The bad
 - require multiple data instances
 - or the ability to divide data into i.i.d. subsamples
 - multiple runs lead to inefficiency



Likelihood Ratio Test for block models

- Model selection between a pair of nested models as a hypothesis test
- Test results have proper confidence intervals
- The likelihood ratio test (LRT) is the uniformly most powerful test
- Basis for many off-the-shelf statistical tools (AIC)

Constructing a LRT

- Null model H_0 , nesting alternative H_1

$$\Lambda(G) = \log \frac{\sup_{H \in H_1} P(G | H)}{\sup_{H \in H_0} P(G | H)},$$

- Reject the null model when Λ exceeds some threshold, which is based on
 - our desired error rate
 - Null distribution of Λ
- To get the Null distribution of Λ , we can
 - analytic prediction
 - parametric bootstrapping

LRT for block models

- Classic $\frac{1}{2}\chi_\ell^2$ result
- Key assumptions:
 - parameter estimates have Gaussian distributions
 - central limit theorems for IID data
 - large data limit
- Networks data is relational
- The latent block assignment variables are discrete
- Sparse networks far from large data limit

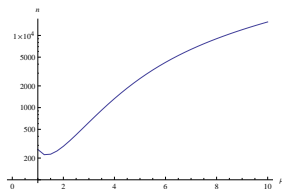
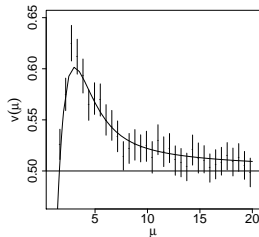
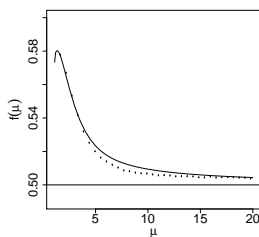
Bootstrapping results using BP

Likelihood Ratio Test for SBM vs DC-SBM

- H_0 : Graph is generated according to the Vanilla-SBM
- H_1 : Graph is generated according to the DC-SBM, where an edge is generated with probability $\theta_u \theta_v p_{g(u)g(v)}$.

$$\Lambda(G) = \log \frac{\sup_{H \in H_1} P(G | H)}{\sup_{H \in H_0} P(G | H)},$$

- According to classic χ^2 test, $\Lambda(G) \sim \frac{1}{2} \chi_\ell^2$ with $\ell = n - k$



The Bayesian approach

- Integrating over parameters of different fit
- The posterior is proportional to total likelihood
- BIC has close relation with Minimum Description Length

Posterior of the SBM

$$P(M_i | G) = \frac{P(M_i)}{P(G)} P(G | M_i)$$
$$\propto \sum_g \iiint_0^1 d\{p_{st}\} d\{q_s\} P(G, g | q, p)$$

- Uniform prior on $P(M_i)$
- Constant evidence $P(G)$
- Bayes factor

Bayesian model selection for SBMs

- The good
 - compare any model with proper posterior
 - combine domain prior with data
 - conjugate priors lead to tractability
- The bad
 - realistic priors often not conjugate
 - model selection is inherently NOT Bayesian
- Belief propagation for the sum over integrated likelihood

The variational expectation maximization framework

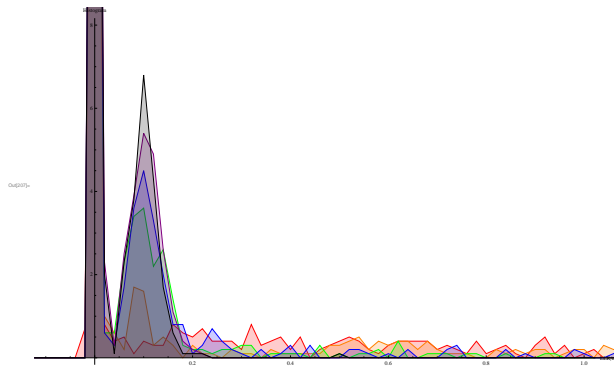
- A flexible framework for community based models
 - Mixed membership block model (Airoldi, Blei, Fienberg and Xing, 2008)
 - The Ball-Karrer-Newman model (2011)
- Choice of variational distribution balance between accuracy and scalability
- SBM: Linear Belief Propagation for the exponential sum
- The analogy between machine learning and statistical physics is powerful

Applications

- Works for both Frequentist and Bayesian model selection
- The right scalability and accuracy lead to better theories
- Project: Bayesian recommendation system based on the SBM (Roger Guimerà)
- Even faster algorithm with message sub-sampling

Looking for Postdoc opportunities

- Getting my Ph.D. this July
- Up for interesting projects in any discipline
- Preferably in the US, but open for other places
- everyxt@gmail.com



Order selection for SBMs

- Even a challenge for classic i.i.d mixture models
- Degeneracy at zero
- The non-zero peak scale with the size of network
- AIC and BIC are bad for order selection