

Scaling MMSB to Large Datasets

Akshat Jindal(150075) Shivam Utreja(150682) Pratyush Garg(150521)

IIT Kanpur

April 25, 2019

- Relational data is ubiquitous nowadays and models to handle such data are thus important to study. In this project, we explore Mixed Membership Stochastic Blockmodels in detail starting from scratch and exploring ways to achieve scalability over baseline implementations.

- Relational data is ubiquitous nowadays and models to handle such data are thus important to study. In this project, we explore Mixed Membership Stochastic Blockmodels in detail starting from scratch and exploring ways to achieve scalability over baseline implementations.
- We explore and implement a Nested VI algorithm to achieve faster convergence and scalability over naive inference procedures for MMSB.

- Relational data is ubiquitous nowadays and models to handle such data are thus important to study. In this project, we explore Mixed Membership Stochastic Blockmodels in detail starting from scratch and exploring ways to achieve scalability over baseline implementations.
- We explore and implement a Nested VI algorithm to achieve faster convergence and scalability over naive inference procedures for MMSB.
- We also explore and implement a subsampling based approach to achieve scalability.

- Relational data is ubiquitous nowadays and models to handle such data are thus important to study. In this project, we explore Mixed Membership Stochastic Blockmodels in detail starting from scratch and exploring ways to achieve scalability over baseline implementations.
- We explore and implement a Nested VI algorithm to achieve faster convergence and scalability over naive inference procedures for MMSB.
- We also explore and implement a subsampling based approach to achieve scalability.
- To better report our results, we implement the baseline approach and a few scalable approaches for a qualitative comparison.

- **Stochastic Blockmodels** : The stochastic blockmodel is an extension or adaptation of mixture models. In that model, each object(data point x_n , belongs to a cluster and the relationships between objects are governed by the corresponding pair of clusters.
- **Limitation** : Here, each object can only belong to one cluster, or in other words, play a single latent role.
- When a protein or a social actor interacts with different partners, different functional or social contexts may apply.

The full generative model looks like

- For $n = 1, \dots, N$

$$\mathbf{z}_n \sim \text{multinoulli}(\boldsymbol{\pi})$$

- For $n = 1, \dots, N$

- For $m = 1, \dots, n - 1$

$$A_{nm} \sim \text{Bernoulli}(\eta_{\mathbf{z}_n, \mathbf{z}_m})$$

Mixed Membership Stochastic Block Model

Setting

- Graph $\mathcal{G} = (\mathcal{N}, \mathbf{Y})$ where \mathcal{N} is the number of vertices.
- $\mathbf{Y} : \mathcal{N} \times \mathcal{N}$ adjacency matrix
- K factions modelled by a $K \times 1$ mixture vector for every vertex and context dependent latent vectors for each vertex.
- Probabilities of interactions between different factions : \mathcal{B}

Generative Story

- 1 Draw a K dimensional mixed membership vector $\pi_p \sim \text{Dirichlet}(\alpha) \quad \forall p \in \mathcal{N}$
- 2 For each pair of vertices $(p, q) \in \mathcal{G}$:
 - Draw membership indicator for the initiator : $z_{p \rightarrow q} \sim \text{Multinomial}(\pi_p)$
 - Draw membership indicator for the initiator : $z_{q \rightarrow p} \sim \text{Multinomial}(\pi_q)$
 - Sample value of interaction : $Y(p, q) \sim \text{Bernoulli}(B_{z_{p \rightarrow q}, z_{q \rightarrow p}})$

Sparsity Parameter : ρ

- Adjacency matrices are often sparse.
- These non interactions or zeros in the matrix can either be due to rarity of interactions inherent in the data, or they may be an indication that the pair of relevant blocks rarely interacts.
- The 2^{nd} form of sparsity is captured by MMSB as it infers faction-faction interaction strength \mathbf{B} . 1^{st} form is due to our data.
- Sparsity parameter $\rho \in [0, 1]$ is introduced.
- we down-weight the probability of successful interaction by $(1 - \rho)$ where the weight ρ captures the portion of zeros that should not be explained by the blockmodel \mathbf{B} .

Parameters to be estimated : $\{\alpha, z_{p \rightarrow q} s, \pi_p s, \mathbf{B}\}$. The authors use a Variational EM approach and estimate posterior distributions for $\pi_p s$ and $z_{p \rightarrow q} s$ using Mean-Field VB, and get point estimates for $\Theta : \mathbf{B}$ and α .

Posterior Inference : The E-Step

- We'll maximise $\mathcal{L}(q, \Theta^{old})$ in the E-Step wrt q , which gives q^{opt} as the joint CP of π_p s and $z_{p \rightarrow q}$ s.
- Since, this is intractable to calculate, we use Mean-Field VB to get q^{opt}
- Variational distribution taken :

$$q(\pi, \mathbf{Z} | \gamma, \phi) = \prod_p q_1(\pi_p | \gamma_p) \prod_{p,q} q_2(z_{p \rightarrow q} | \phi_{p \rightarrow q}) q_2(z_{q \rightarrow p} | \phi_{q \rightarrow p})$$

where q_1 is Dirichlet and q_2 is Multinomial. Thus the variational parameters to be inferred : $\{\gamma, \phi\}$

- Solving via the Mean-Field approach taught in class, the update equations for the parameters :

$$\phi_{p \rightarrow q, g}^{new} \propto e^{E_q[\log \pi_{p, g}]} \prod_h ((\mathbf{B}(g, h))^{Y(p, q)} (1 - \mathbf{B}(g, h))^{1 - Y(p, q)})^{\phi_{q \rightarrow p, h}^{old}}$$

$$\phi_{q \rightarrow p, h}^{new} \propto e^{E_q[\log \pi_{q, h}]} \prod_g ((\mathbf{B}(g, h))^{Y(p, q)} (1 - \mathbf{B}(g, h))^{1 - Y(p, q)})^{\phi_{p \rightarrow q, g}^{old}}$$

$$\gamma_{p, k}^{new} = \alpha_k + \sum_q \phi_{p \rightarrow q, k}^{new} + \sum_q \phi_{q \rightarrow p, k}^{new}$$

for all nodes $p = 1, 2, \dots, N$ and $k = 1, 2, \dots, K$. and $g, h = 1, \dots, K$.

Point Estimates: The M-Step

- Using the variational distribution, we maximise the $\mathcal{L}(q^{optimal}, \Theta)$ wrt $\Theta = \{\mathbf{B}\}$
- For \mathbf{B} , we have :

$$\hat{B}(g, h) = \frac{\sum_{p,q} Y(p, q) \cdot \phi_{p \rightarrow qg} \phi_{p \leftarrow qh}}{(1 - \rho) \cdot \sum_{p,q} \phi_{p \rightarrow qg} \phi_{p \leftarrow qh}},$$

- The value of K is set either by using BIC as discussed in class or cross-validation if the model size is not large enough for BIC to be effective.
- ρ is also estimated as $1 - \hat{d}$, where $\hat{d} = \frac{\sum_{p,q} Y(p, q)}{N^2}$

Scaling Up : Naive Inference

- ❶ In naive inference, γ and ϕ are initialized to non-informative values, and then we iterate the following two steps until convergence:
- ❷
 - Update $\phi_{p \rightarrow q}$ and $\phi_{q \rightarrow p}$ for all edges (p, q) .
 - Update γ_p for all nodes $p \in \mathcal{N}$
 - Now, update \mathbf{B} using the updated values of ϕ
- ❸ \mathbf{B} is updated only once at each iteration
- ❹ Every ϕ^t uses ϕ^{t-1} and γ^t
- ❺ $NK + 2N^2K$ scalars maintained across iterations
- ❻ Doesn't capture the dependence between γ s and \mathbf{B} .
- ❼ Slow convergence :(

Nested VI

- Nested VI takes care of this by rescheduling the parameter updates.
- Always keep the block of free parameters, $\phi_{p \rightarrow q}$ and $\phi_{q \rightarrow p}$, optimized given the other variational parameters : Faster convergence
- B is updated at every pair p,q of vertices now. Dependence between B and γ thus retained : Faster convergence
- At each variational cycle we need to allocate $NK + 2K$ scalars only : Scalable

```

1. initialize  $\tilde{\gamma}_{pk}^0 = \frac{2N}{K}$  for all  $p, k$ 
2. repeat
3.   for  $p = 1$  to  $N$ 
4.     for  $q = 1$  to  $N$ 
5.       get variational  $\tilde{\phi}_{p \rightarrow q}^{t+1}$  and  $\tilde{\phi}_{p \leftarrow q}^{t+1} = f(Y(p, q), \tilde{\gamma}_p, \tilde{\gamma}_q, B^t)$ 
6.       partially update  $\tilde{\gamma}_p^{t+1}, \tilde{\gamma}_q^{t+1}$  and  $B^{t+1}$ 
7.   until convergence

```

```

5.1. initialize  $\phi_{p \rightarrow q, g}^0 = \phi_{p \leftarrow q, h}^0 = \frac{1}{K}$  for all  $g, h$ 
5.2. repeat
5.3.   for  $g = 1$  to  $K$ 
5.4.     update  $\phi_{p \rightarrow q}^{s+1} \propto f_1(\tilde{\phi}_{p \leftarrow q}^s, \tilde{\gamma}_p, B)$ 
5.5.     normalize  $\phi_{p \rightarrow q}^{s+1}$  to sum to 1
5.6.   for  $h = 1$  to  $K$ 
5.7.     update  $\phi_{p \leftarrow q}^{s+1} \propto f_2(\tilde{\phi}_{p \rightarrow q}^s, \tilde{\gamma}_q, B)$ 
5.8.     normalize  $\phi_{p \leftarrow q}^{s+1}$  to sum to 1
5.9.   until convergence

```

Figure: PseudoCode

Scaling MMSB : Method - II

- Algorithm suggested by Prem K. Gopalan and David M. Blei.
- Successfully scales MMSB to real-world social/citation/biological networks.
- Uses a combination of Mean Field VI and Stochastic Optimization to speed-up MMSB.

Model description.

The model used here is to model assortative undirected networks. The generative story is similar to the previous model, with only subtle changes.

- For each node, draw community memberships θ_i from $Dirichlet(\alpha)$
- For each pair of nodes i and j such that $i < j$;
 - Draw community indicator $z_{i \rightarrow j}$ from $Multinoulli(\theta_i)$.
 - Draw community indicator $z_{i < j}$ from $Multinoulli(\theta_j)$.
 - Draw connection between them using;

$$p(y_{ij} = 1 | z_{i \rightarrow j}, z_{i < j}) = \begin{cases} \beta_{z_{i \rightarrow j}} & \text{if } z_{i \rightarrow j} = z_{i < j} \\ \epsilon & \text{if } z_{i \rightarrow j} \neq z_{i < j} \end{cases} \quad (1)$$

Note that in this case, the nodes form a connection with high probability only if their latent factors are from the same community. It still supports overlapping communities, because the z s themselves are being drawn from a Multinoulli.

Posterior Approximation

The slight difference in the model is reflected in the slightly different, mean-field posterior approximation;

$$q(\beta, \theta, \mathbf{z}) = \prod_{k=1}^K q(\beta_k | \lambda_k) \prod_{n=1}^N q(\theta_n | \gamma_n) \prod_{i < j} q(\mathbf{z}_{i->j} | \phi_{i->j}) q(\mathbf{z}_{i<-j} | \phi_{i<-j}) \quad (2)$$

Each variable has its own distribution which makes the model very flexible, and allows us to learn each node's community memberships.

Posterior Updates

Given the above form of $q(\beta, \theta, z)$, we get the following form for the ELBO;

$$\begin{aligned}\mathcal{L} = & \sum_k \mathbb{E}_q [\log p(\beta_k | \eta)] - \sum_k \mathbb{E}_q [\log q(\beta_k | \lambda_k)] \\ & + \sum_n \mathbb{E}_q [\log p(\theta_n | \alpha)] - \sum_n \mathbb{E}_q [\log q(\theta_n | \gamma_n)] \\ & + \sum_{a,b} \mathbb{E}_q [\log p(z_{a \rightarrow b} | \theta_a)] + \mathbb{E}_q [\log p(z_{a \leftarrow b} | \theta_b)] \\ & - \sum_{a,b} \mathbb{E}_q [\log q(z_{a \rightarrow b} | \phi_{a \rightarrow b})] - \mathbb{E}_q [\log q(z_{a \leftarrow b} | \phi_{a \leftarrow b})] \\ & + \sum_{a,b} \mathbb{E}_q [\log p(y_{ab} | z_{a \rightarrow b}, z_{a \leftarrow b}, \beta)].\end{aligned}$$

[S2]

Figure: The ELBO

Posterior Updates

- In the above equation, the summations corresponding to the communities and the nodes are called the global terms, and those related to the edges are called the local terms.
- Their corresponding parameters are called the global and local parameters respectively.
- Hence, λ and γ become our global parameters, and Φ become our local parameters.
- We can now perform alternating updates between the global(Global Step) and local(Local Step) parameters, until we converge to a local optimal of the ELBO.
- The local parameters are set using the current estimates of the global parameters.(They, themselves need to be updated in an alternating fashion. Precise details covered in the report.)
- The global parameters are updated using the gradient computed from the current estimate of the local parameters.

Sampling techniques

- Note that in the above formulation, performing the global step would involve finding the gradient using N^2 -many local parameters. N can be of the order of millions.
- Hence, at each iteration, it is crucial that we sub-sample the graph's edges (i.e. local parameters) to make the algorithm tractable for real networks.
- Note that, any sampling technique will work, as long as the noisy gradient is unbiased, and its variance can be controlled.
- Also note, that this also acts as a natural way to be able to interleave data collection and model estimation.

Sampling techniques

The author suggests several sampling techniques in the paper.

- Random pair sampling.
- Random node sampling.
- Link Sampling

Link Sampling

- The other sampling schemes suggested above take into account, all pairs of nodes, irrespective of whether they are linked.
- This becomes increasingly redundant as the size of the graph increases. The number of links is far less than the number of non-links.
- The link sampling scheme overcomes this redundancy. However, we need to change the form of our approximation first.

$$q(\beta, \theta, z) = \prod_{k=1}^K q(\beta_k | \lambda_k) \prod_{n=1}^N q(\theta_n | \gamma_n) \prod_{(i,j) \in \text{links}} q(z_{i->j}, z_{i-<j} | \phi_{ij}) \prod_{(i,j) \in \text{nonlinks}} q(z_{i->j} | \phi_{i->j}) q(z_{i-<j} | \phi_{i-<j}) \quad (3)$$

Initialization and Complexity

- The authors provide an algorithm that initialises the global parameters, as well as sets the number of communities for the SVI algorithm.
- The local step for the SVI algorithm can be computed in $O(SK)$ operations.
- Note that it is not quadratic in K because we are working with the assortativity assumption.
- The time for the global step is $O(NK)$ per iteration. We can further reduce this to $O(nK)$, where n is the minibatch size; if we maintain distinct learning rates for each node.
- In the case of the link sampling algorithm, sampling all links together, gives a complexity of $O(MK + NK)$ per iteration. Its convergence is much faster, even without subsampling.

Results and Implementation

- We implemented the naive and nested VI algorithms for a qualitative comparison.
- We tested them on 2 graphs : A 5*5 graph (small) and a 55*55 graph(big)
- Reconstruction : $\mathbf{E}[\mathbf{Y}(\mathbf{p}, \mathbf{q})] = \phi_{\mathbf{p} \rightarrow \mathbf{q}}^{\mathbf{T}} \mathbf{B} \phi_{\mathbf{q} \rightarrow \mathbf{p}}$

Small Graph

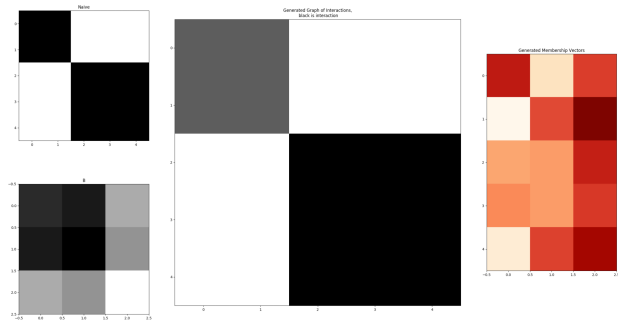


Figure: Results for Naive Inference on 5*5 graph, $K=2$

Small Graph

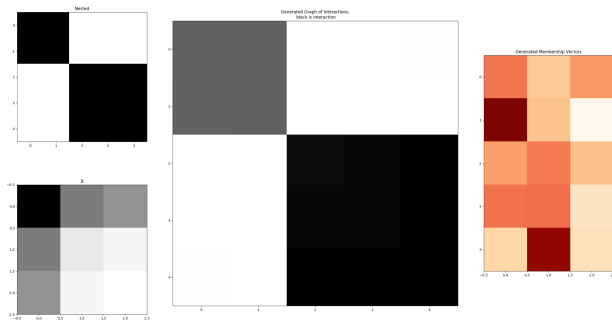


Figure: Results for Nested Inference on 5*5 graph, K=2

Large Graph

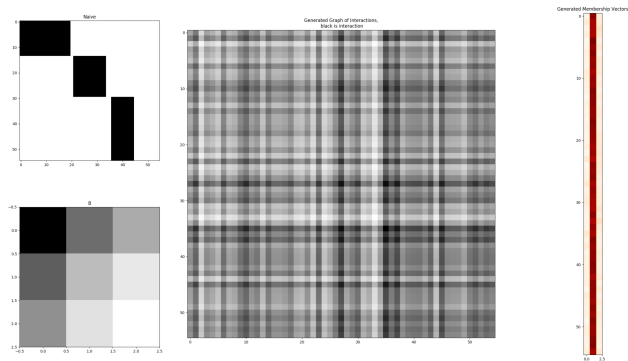


Figure: Results for Naive Inference on 55*55 graph, $K=3$. Notice poor result

Large Graph

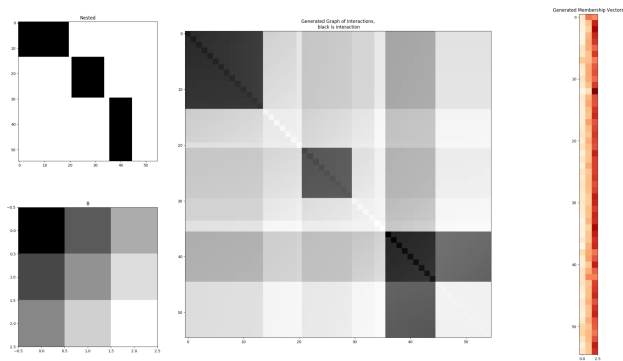


Figure: Results for Nested Inference on 55*55 graph, $K=3$. Notice that this has converged!

- MMSB is a complex model to perform inference on, especially when dealing with large, real world networks.
- The variational inference methods discussed above, prove to act as a powerful tool in this context; to infer the large number of variational parameters in this specific model.
- We also noted that scaling-up the method requires utilizing sampling methods and working with subgraphs.
- Link sampling showed that sampling methods that exploit the structure in the problem give better results.
- Finally, we feel that applying techniques such as BBVI to this problem (that reduce the variance of noisy gradients), should be looked into as part of the future work.

Questions ?