

Homework 4: Gibbs Samplers for Stochastic Block Models of Relational Data

Brown University CS 242: Probabilistic Graphical Models

Homework due at 11:59pm on November 16, 2016

Markov chain Monte Carlo (MCMC) algorithms may be used to learn relational models of biological or social networks. Consider a set of N nodes, representing entities to be modeled. For node pairs $i \neq j$, we let $y_{ij} = 1$ if there is a relationship (e.g., friendship) from entity i to entity j , and $y_{ij} = 0$ otherwise. We focus on directed graphs, for which there is a distinct binary relationship y_{ji} from entity j to entity i . For some experiments, the y_{ij} variables will only be observed for a (known) subset of ordered entity pairs (i, j) . All of your derivations and code should support such partial observations.

We model such relational data as being generated from K unknown latent communities, indexed by integers $k = 1, \dots, K$. For each pair of communities k, ℓ , we define an interaction probability $W_{k\ell}$, and place a uniform prior distribution:

$$W_{k\ell} \sim \text{Beta}(1, 1), \quad 1 \leq k, \ell \leq K.$$

We will then use Gibbs samplers to learn this non-symmetric $K \times K$ matrix of community interaction probabilities, as well as community memberships for individual entities.

The *stochastic block models* we consider are closely related to probabilistic mixture models and the latent Dirichlet allocation (LDA) topic model, and you may find it helpful to review Gibbs samplers for those models. Also remember that the Dirichlet distribution is conjugate to categorical likelihoods, and the beta distribution to Bernoulli likelihoods.

Question 1: Stochastic Block Models

The basic stochastic block model assumes each entity i is a member of a single latent community z_i , sampled according to

$$z_i \sim \text{Cat}(\pi), \quad i = 1, \dots, N.$$

The K -dimensional distribution of community frequencies π has a symmetric Dirichlet prior:

$$\pi \sim \text{Dir}(\alpha, \dots, \alpha).$$

For all experiments involving the basic stochastic block model, we set $\alpha = 1$. For any pair of nodes $i \neq j$, their relationship link variables are generated according to

$$p(y_{ij} = 1 \mid z_i = k, z_j = \ell, W) = W_{k\ell}, \quad p(y_{ji} = 1 \mid z_i = k, z_j = \ell, W) = W_{\ell k}.$$

Remember that we may only have partial observations of these link variables.

- a) Given fixed parameters π, W , and observations y for some subset of ordered entity pairs, derive a formula for the posterior distribution $p(z_i \mid z_{\setminus i}, y, \pi, W)$. Here, $z_{\setminus i}$ denotes the community assignments for all entities except node i . **Hint:** Remember that z_i impacts the outgoing relationships y_{ij} from entity i to other entities j , as well as the incoming relationships y_{ji} from those other entities to node i .
- b) Given fixed entity assignments z , derive formulas for the posterior distributions of the model parameters, $p(\pi \mid y, z, W)$ and $p(W \mid y, z, \pi)$. These should be members of some standard exponential family of distributions.
- c) We have provided skeleton code for a block model Gibbs sampler in `sbGibbs.m`, which provides appropriate initializations and evaluates the joint log-probability at each iteration:

$$\begin{aligned} \log p(\pi, W, z, y) = & \log \text{Dir}(\pi \mid \alpha) + \sum_{k=1}^K \sum_{\ell=1}^K \log \text{Beta}(W_{k\ell} \mid 1, 1) \\ & + \sum_{i=1}^N \log \text{Cat}(z_i \mid \pi) + \sum_{i=1}^N \sum_{j \neq i} \log \text{Ber}(y_{ij} \mid W_{z_i z_j}) \end{aligned}$$

Using the formulas from parts (a,b), implement code which resamples each of the z, W, π variables once per iteration.

- d) Test your sampler using the code in `testToy.m`, which constructs a synthetic dataset with $N = 30$ entities, 10 in each of $K = 3$ communities. Generate links y_{ij} by assuming a within-community link probability of $W_{kk} = 0.95$, and a between-community link probability of $W_{k\ell} = 0.10, \ell \neq k$. Then given only the observed links y , explore how accurately the sampler recovers the underlying W, z variables. Run your sampler for 1000 iterations from each of 5 random initializations, and plot the resulting log-likelihood curves on a single set of axes. After each iteration, compute the Rand index (see provided script) between the true assignments and the sampled z , and plot these scores versus iteration.
- e) Using the code in `testSampson.m`, apply the stochastic block model to the Sampson monk data. Allow $K = 3$ communities, run the sampler for 1000 iterations from each of 5 random initializations, and again plot log-likelihoods and Rand indexes (from the true faction labels) versus iteration.

Question 2: Mixed Membership Stochastic Block Models

The mixed membership stochastic block model generalizes basic block models by allowing each entity to participate in multiple communities. We begin by sampling a K -dimensional community membership distribution for each entity:

$$\pi_i \sim \text{Dir}(\alpha, \dots, \alpha), \quad i = 1, \dots, N.$$

For all experiments involving mixed membership block models, assume $\alpha = 0.1$ to encourage membership in a sparse subset of communities. Observation y_{ij} is then determined by link-specific source and receiver community assignments, s_{ij} and r_{ij} , sampled as follows:

$$p(y_{ij} = 1 \mid s_{ij} = k, r_{ij} = \ell, W) = W_{k\ell}, \quad s_{ij} \sim \text{Cat}(\pi_i), \quad r_{ij} \sim \text{Cat}(\pi_j).$$

Intuitively, s_{ij}, r_{ij} are the communities which “explain” the interactions of node i with node j .

- a) Given fixed parameters π, W , and observations y for some subset of ordered entity pairs, derive a formula for the posterior distributions $p(r \mid s, y, \pi, W)$ and $p(s \mid r, y, \pi, W)$. Are the indicators r_{ij}, s_{ij} for different links conditionally independent?
- b) Given fixed link assignments r, s , derive formulas for the posterior distributions of the model parameters, $p(\pi \mid y, r, s, W)$ and $p(W \mid y, r, s, \pi)$. These should be members of some standard exponential family of distributions.
- c) We have provided skeleton code for a mixed membership block model Gibbs sampler in `mmsbGibbs.m`, which provides appropriate initializations and evaluates the joint log-probability at each iteration:

$$\begin{aligned} \log p(\pi, W, s, r, y) = & \sum_{i=1}^N \log \text{Dir}(\pi_i \mid \alpha) + \sum_{k=1}^K \sum_{\ell=1}^K \log \text{Beta}(W_{k\ell} \mid 1, 1) \\ & + \sum_{i=1}^N \sum_{j \neq i} \left[\log \text{Cat}(s_{ij} \mid \pi_i) + \log \text{Cat}(r_{ij} \mid \pi_j) + \log \text{Ber}(y_{ij} \mid W_{s_{ij}r_{ij}}) \right] \end{aligned}$$

Using the formulas from parts (a,b), implement code which resamples each of the r, s, W, π variables once per iteration.

- d) Apply your sampler to the synthetic data from part 1(d), assuming $K = 3$, running for 1000 iterations from each of 5 initializations, and plotting log-likelihood versus iteration. Looking at the result from the highest-likelihood Markov chain, does the mixed membership model provide a reasonable interpretation of this data?
- e) Apply your sampler to the Sampson monk data, assuming $K = 3$, running for 1000 iterations from each of 5 initializations, and plotting log-likelihood versus iteration. Looking at the result from the highest-likelihood Markov chain, how many monks have significant membership in more than one community?
- f) The basic sampler outlined above may mix slowly, due to correlations between the source and receiver link variables. To address this, derive a formula for the joint posterior distribution $p(r, s \mid y, \pi, W)$. Are the paired indicators (r_{ij}, s_{ij}) for different links conditionally independent? Hint: You should be able to sample from this conditional distribution by defining an appropriate K^2 -dimensional categorical distribution.
- g) Repeat the experiments from parts (d,e) using the blocked sampler, including creation of log-likelihood plots. Are there substantial performance differences between the samplers?

Question 3: Link Prediction

In this question, we use our mixed membership models to predict the presence of likely links in a partially sampled social network. We focus on a network describing advice relationships among $N = 71$ attorneys in a New England law firm. The `testLawyer.m` script provides code that randomly subsamples half of the y_{ij} link variables to use for training, and reserves the other half for testing. Note that the training data is half of the directed node pairs (potential links), *not* half of the actually present links.

- a) *Using a single train-test split of the attorney network, apply your Gibbs sampler from problem 1, as well as the blocked Gibbs sampler from problem 2. For both models assume $K = 6$, run the sampler for 1000 iterations from each of 3 random initializations, and reserve the final sample from the most probable Markov chain.*
- b) *For the standard block model, derive a formula for the posterior distribution of each of the test link variables y_{ij} , given the parameters z, π, W learned during training. Implement this prediction formula in `sbPredict.m`, and use the held-out test labels to create an ROC curve summarizing classification performance.*
- c) *For the mixed membership block model, compute the posterior distribution of each of the test link variables y_{ij} , given the parameters π, W learned during training. The source and receiver link variables s_{ij}, r_{ij} should be analytically marginalized. Implement this prediction formula in `mmsbPredict.m`, and use the held-out test labels to create an ROC curve summarizing classification performance. Compare your results to part (b).*
- d) *With more computational effort, is there a more sophisticated way you could predict test link variables based on the output of your Gibbs sampling algorithms?*