

Gene Interaction Networks using the Gibbs Sampling Technique

A gene interaction network (GIN) in this method is the same as the one defined in Toyoshima et.al. [1]. Briefly, a GIN consists of a collection of P genes, denoted G_1, G_2, \dots, G_P , with observed values g_1, g_2, \dots, g_P . Define β_{ij} ($i, j=1, 2, \dots, P$) as parameters in the log-linear function form describing the linkage from gene i to gene j .

Mathematically, this is written as

$$E[\log(G_j)] = \sum_{i=1, i \neq j}^P I_{ij} \beta_{ij} \log(g_i) \quad (1)$$

where $E[\log(G_j)]$ represents the expectation for the natural logarithm of G_j and I_{ij} ($i, j = 1, 2, \dots, P$) is an indicator function that equals 1 if gene G_i has a link to gene G_j , otherwise it equals 0. If a gene has a regulatory effect on gene G_i , then that gene is referred to as a “Parent of gene G_i ” and we refer to it as belonging to the set $Pa(G_i)$.

Define \mathbf{T} to be the transition matrix whose (i, j) element consists of the indicator function I_{ij} and let \mathbf{B} represents the matrix whose (i, j) elements consists of β_{ij}

Prior distributions

In the log-linear model, the parameters are the variance for the distribution for each gene, denoted σ_j^2 , parameters in the log-linear form, β_{ij} ($i, j=1, 2, \dots, P, i \neq j$), and the

indicator functions defining the network, I_{ij} . Prior distributions can be assumed for each parameter using either informative or uninformative priors. Let $h(\sigma_j^2)$ denote the prior distribution for σ_j^2 . Define $h(\beta_{ij}|I_{ij}=1)$ to be the prior distribution for β_{ij} when, in a given network structure, $I_{ij}=1$ (i.e. there is a linkage from G_i to G_j). Similarly, define $h(\beta_{ij}|I_{ij}=0)$. Finally, define $h(I_{ij})$ to be the prior distribution for I_{ij} , $h(I_{ij})$ is a discrete distribution taking on only two values, 0 or 1.

In what follows, the prior distribution for the variance, $h(\sigma_j^2)$, is assumed to be inverse Gamma with parameters a_1 and a_2 (i.e. $h(\sigma_j^2)$ is $inv\Gamma(a_1/2, a_2/2)$). For β_{ij} , we assume $h(\beta_{ij}|I_{ij})$ is normal with mean 0 and variance $\sigma(\beta)^2(I_{ij}+(1-I_{ij})C)^2$ (i.e. $h(\beta_{ij}|I_{ij}=1)$ is $N(0, \sigma(\beta)^2(I_{ij}+(1-I_{ij})C)^2)$ where C is some constant to keep β_{ij} close to zero when $I_{ij}=0$. This is the key to initially guide the Gibbs sampler in choosing between networks with $I_{ij}=1$ and those with $I_{ij}=0$. In general, denoting β_j as the j^{th} column vector of \mathbf{B} with j^{th} element is absent, the prior distribution of β_j can be written by P-1 dimension multivariate normal distribution $N_{P-1}(0, \mathbf{C}\mathbf{C}|\mathbf{I}_j)$, where \mathbf{C} is a diagonal matrix with $\sigma(\beta)(I_{ij}+(1-I_{ij})C)$ is the (i, i) elements and \mathbf{I}_j is the j^{th} column vector of \mathbf{T} with the j^{th} element is absent. The prior distribution for I_{ij} was assumed to be a Bernoulli distribution with success ($I_{ij}=1$) probability p_{ij} . In the uninformative case, p_{ij} could be set to 0.5 and if there is some expectation that I_{ij} is not equal to zero, the prior

probability could be set higher or the correlation information between some linkages could be possible as well in more advance.

With these prior distributions, and assuming that the natural log of G_j follows a normal distribution with mean $\sum_{i=1, \neq j}^P \beta_{ij} \cdot \log(g_i)$ and standard deviation σ_j , posterior distributions for each parameter can be estimated. It should be noted that the mean function $\sum_{i=1, \neq j}^P \beta_{ij} \cdot \log(g_i)$ is slightly different from the form shown in Equation (1). In the algorithm used here, β_{ij} is sampled even if $I_{ij}=0$. Hence, in the Gibbs sampling application, the summation range will be all genes except for gene G_j itself. The development presented here parallels the methods developed by George and McCulloch [2] using Gibbs sampling to select models for linear regression. Other prior setting had been studied and detail comparison had been made by the other researchers [3,4,5]. In our case, the model selection for linear regression will be extended to the model selection for Bayesian networks. Hence, some improvement to overcome the acyclic condition on Bayesian network will be needed to extend the method and that will be given in the below.

Posterior distributions

Let $\mathbf{g} = [\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_P]$ represent an $n \times P$ matrix containing the gene expression data after normalization and taking log-scale, where n is the sample size and \mathbf{g}_i is a vector of n observations from G_i , and define the complementation set $\mathbf{g}(\mathbf{i}) = [\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_{i-1}, \mathbf{g}_{i+1}, \dots, \mathbf{g}_P]$. The prior distribution assumed in the previous, the conditional posterior distribution for σ_j^2 is given by

$$\sigma_j^2 \sim f(\sigma_j^2 | \mathbf{g}(\mathbf{j}), \boldsymbol{\beta}_j, \mathbf{I}_j) = \text{Inv}\Gamma\left(\frac{n + a_1}{2}, \frac{|\mathbf{g}_j - \mathbf{g}(\mathbf{j}) \cdot \boldsymbol{\beta}_j| + a_2}{2}\right) \quad (2)$$

where $\boldsymbol{\beta}_j$ is the same as the above. In this equation, a_1 and a_2 could be chosen as zero for the non-informative prior. For $\boldsymbol{\beta}_j$, the conditional posterior is given by

$$\boldsymbol{\beta}_j \sim f(\boldsymbol{\beta}_j | \mathbf{g}, \sigma_j, \mathbf{I}_j) = N_P(\mathbf{A} \cdot \sigma_j^{-2} \cdot \mathbf{g}(\mathbf{j})^t \mathbf{g}_j, \mathbf{A}) \quad (3)$$

where the matrix \mathbf{A} is defined by its inverse as

$$\mathbf{A}^{-1} = \sigma_j^{-2} \cdot \mathbf{g}(\mathbf{j})^t \cdot \mathbf{g}(\mathbf{j}) + (\mathbf{C} \mathbf{R} \mathbf{C}^t)^{-1} \quad (4)$$

where \mathbf{R} could be $\mathbf{g}(\mathbf{j})^t \mathbf{g}(\mathbf{j})$, or $P-1$ dimension identical matrix, $\mathbf{E}(P-1)$ and \mathbf{C} is the same as above used in the covariance matrix for $\boldsymbol{\beta}_j$. Finally, I_{ij} is obtained by sampling consecutively from the conditional distribution,

$$I_{ij} \sim f(I_{ij} | \mathbf{g}, \boldsymbol{\beta}_j, \sigma_j, \mathbf{T}/I_{ij}) = f(I_{ij} | \boldsymbol{\beta}_j, \sigma_j, \mathbf{T}/I_{ij}), \quad (6)$$

where \mathbf{T}/I_{ij} represent the other all elements of \mathbf{T} except for I_{ij} . Here, it should be

noted that the posterior distribution of I_{ij} is conditioned by T/I_{ij} in Bayesian model selection. Because, if $I_{ij}=1$ makes cycle in a network, the success probability for I_{ij} will be zero. Hence, the following one step is put into the sampling process. Each distribution is Bernoulli with probability

$$P(I_{ij} = 1 | \beta_{\cdot j}, \sigma_j, T/I_{ij}) = \frac{a}{a + b}, \quad (7)$$

where a is equal to 0 and $b=1$ if $I_{ij}=1$ produces a cyclic in a network, else

$$a = f(\beta_{\cdot j} | T/I_{ij}, I_{ij} = 1) \times f(\sigma_j | T/I_{ij}, I_{ij} = 1) \times f(I_{ij} = 1) \quad (8)$$

$$b = f(\beta_{\cdot j} | T/I_{ij}, I_{ij} = 0) \times f(\sigma_j | T/I_{ij}, I_{ij} = 0) \times f(I_{ij} = 0). \quad (9)$$

By equation (2) and the prior distribution assumed for I_{ij} , (8) and (9) can be given by,

$$a = f(\beta_{\cdot j} | T/I_{ij}, I_{ij} = 1) p_{ij}, \quad (10)$$

$$b = f(\beta_{\cdot j} | T/I_{ij}, I_{ij} = 0) (1 - p_{ij}). \quad (11)$$

Continuing the sampling method for all values of j results in a single observation from the posterior distribution for each parameter. The summary of the posterior distributions for I_{ij} s can be given as a matrix and provide an indication of the strongest or weakest possible linkages between genes.

Sampling Algorithm in Bayesian network

Gibbs sampling using the posterior distributions defined above could work in Bayesian network setting. The posterior distribution of the Bayesian network defined by \mathbf{T} , \mathbf{B} and $\mathbf{S}=[\sigma_1, \sigma_2, \dots, \sigma_P]$ given $\mathbf{g} = [\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_P]$ can be written in the following,

$$\begin{aligned} f(\mathbf{T}, \mathbf{B}, \mathbf{S} | \mathbf{g}) &= f(\mathbf{g} | \mathbf{T}, \mathbf{B}, \mathbf{S}) * h(\mathbf{T}) * h(\mathbf{B} | \mathbf{T}) * h(\mathbf{S}) \\ &= \prod_{j=1}^P f(\mathbf{g}_j | \mathbf{I}_j, \boldsymbol{\beta}_j, \sigma_j) h(\sigma_j) h(\boldsymbol{\beta}_j | \mathbf{I}_j) h(\mathbf{I}_j) \end{aligned} \quad (12)$$

It is well known that the whole posterior distribution over every parameter can be decomposed into sub-networks consisting of each G_j and its parents.

From the decomposition in equation (12), it is obvious that σ_j and $\boldsymbol{\beta}_j$ can be sampled from equation (2) and (3) because the posterior distributions for these parameters are not conditioned by other parameters except for \mathbf{I}_j , $\boldsymbol{\beta}_j$, σ_j . In the below, these are called j^{th} column parameters in \mathbf{T} , \mathbf{B} , and \mathbf{S} . In sampling \mathbf{I}_j , some caution should be needed, because the posterior distribution of \mathbf{I}_j is conditioned by \mathbf{T}/\mathbf{I}_j which is other column parameters in \mathbf{T} . This comes from the acyclic condition of Bayesian networks.

Hence, in sampling I_{ij} in \mathbf{I}_j , if $I_{ij}=1$ makes a cycle in the network, then the probability for $I_{ij}=1$ is set to be zero and if $I_{ij}=1$ does not make cycles, then equation (7) can be obtained by removing out the same terms in equation (12) which is the full posterior distribution and used for sampling I_{ij} . After sampling all \mathbf{I}_j elements (I_{ij} , $i=1,2,3,\dots,P$) then the other column parameters in \mathbf{T} , \mathbf{B} and \mathbf{S} will be sampled as the same manner as explained above. This is one step of the Gibbs sampling for Bayesian network.

References

1. Toyoshiba H, Yamanaka T, Sone H, Parham FM, Walker NJ, Martinez J, et al. 2004. Gene interaction network suggests dioxin induces a significant linkage between aryl hydrocarbon receptor and retinoic acid receptor beta. Environ Health Perspect 112(12):1217-1224.
2. George EI, McCulloch. RE. 1993. Variable Selection Via Gibbs Sampling. Journal of the American Statistics Association 88(423):881-889.
3. Dellaportas P, Forster JJ, Ntzoufras I. 2000. Bayesian Variable Selection Using the Gibbs Sampler Generalized Linear Models: A Bayesian Perspective (D. K. Dey, S. Ghosh, and B. Mallick, eds.). New York: Marcel Dekker.271-286.
4. Dellaportas P, Forster JJ, Ntzoufras I. 2002. On Bayesian Model and Variable Selection Using MCMC. Statistics and Computing 12:27-36.
5. Ntzoufras I. 2002. Gibbs Variable Selection Using BUGS. Journal of Statistical Software 7(7):1-19.