

SME0809 - Inferência Bayesiana - Prova 2

High-Dimensional Multivariate Bayesian Variable and Covariance Selection in Linear Regression (Zhao et al. 2021)

Grupo 13 - Francisco Miranda - 4402962 - Heitor Carvalho - 11833351

Dezembro 2021

Introdução

Com o desenvolvimento de técnicas de alto processamento na biologia molecular, a caracterização molecular em alta escala tornou-se um lugar comum, com o advento de técnicas como:

- genome-wide measurement of gene expression
- single nucleotide polymorphisms
- CpG methylation status
- pharmacological profiling for large-scale cancer drug screen.

A análise de associações conjuntas entre múltiplos fenótipos correlacionados e atributos moleculares de alta dimensionalidade é desafiadora.

Quando múltiplos fenótipos e informação genômica de alta dimensionalidade são analisados conjuntamente, a abordagem bayesiana permite especificar de maneira flexível as relações complexas entre os conjunto de dados altamente estruturados.

O pacote **BayesSUR** combina diversos modelos que foram propostos para a regressão multidimensional com resposta múltipla e introduz um novo modelo, que permite diferentes *prioris* na seleção de variáveis dos modelos de regressão e para diferentes pressupostos a respeito da estrutura de dependência entre as respostas.

Metodologia

- múltiplas opções de seleção de variáveis
- a matriz de covariância pode ser diagonal, densa ou esparsa.
- engloba três classes de modelos de regressão linear de múltipla resposta:
- HRR
- dSUR e SSUR
- MRF

O modelo de regressão é escrito como:

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{U} \quad (1)$$

$$\text{vec}(\mathbf{U}) \sim \mathcal{N}(\mathbf{0}, \mathbf{C} \otimes \mathbb{I}_n)$$

onde:

- \mathbf{Y} é uma matriz $s \times s$ das variáveis resposta com matriz de covariância \mathbf{C} ;
- \mathbf{X} é uma matriz $n \times p$ de preditores para todas as respostas;
- \mathbf{U} é a matriz dos resíduos;
- $\text{vec}(\cdot)$ denota a vetorização da matriz;
- $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denota uma distribuição normal multivariada com vetor de médias $\boldsymbol{\mu}$ e matriz de covariâncias $\boldsymbol{\Sigma}$;
- $\mathbf{0}$ denota um vetor coluna com todos os elementos nulos,
- \otimes é o produto de Kronecker e \mathbb{I}_n a matriz identidade de ordem n .

A seleção de variáveis é realizada através de uma matriz indicadora binária latente $\boldsymbol{\Gamma} = \{\gamma_{jk}\}$.

Uma *priori* “spike-and-slab” é utilizada para encontrar um subconjunto esparso relevante de preditores que expliquem a variabilidade de \mathbf{Y} : condicional em $\gamma_{jk} = 0$ ($j = 1, \dots, p$, e $k = 1, \dots, s$)

Definem-se $\beta_{jk} = 0$ condicionado em $\gamma_{jk} = 1$ seguem uma distribuição normal difusa:

$$\beta_\gamma \sim \mathcal{N}(\mathbf{0}, W_\gamma^{-1}) \quad (2)$$

Onde $\beta = \text{vec}(\mathbf{B})$, $\gamma = \text{vec}(\boldsymbol{\Gamma})$, β_γ consiste somente nos coeficientes selecionados (i.e. $\gamma_{jk} = 1$), assim W_γ é a sub matriz de \mathbf{W} formada pelos coeficientes selecionados correspondentes.

A matriz de precisão, \mathbf{W} , é geralmente decomposta em coeficientes de encolhimento e uma matriz que governa a estrutura de covariância dos coeficientes de regressão. É utilizado aqui $W = w^{-1}\mathbb{I}_{sp}$, o que significa que todos os coeficientes de regressão são independentes a priori, com uma *hiperpriori* no coeficiente de encolhimento w , i.e. $w \sim \mathcal{IGamma}(a_w, b_w)$.

A matriz indicadora binária latente $\boldsymbol{\Gamma}$ tem três opções de *priori*:

- Bernoulli independente hierárquica
- hotspot prior
- MRF prior

A matriz de covariância \mathbf{C} também possui três *prioris*:

- Gama inversa independente
- Wishart inversa
- hiper-inversa Wishart

São considerados no total nove possíveis modelos dentre as combinações de \mathbf{C} e $\boldsymbol{\Gamma}$

	$\gamma_{jk} \sim \text{Bernoulli}$	$\gamma_{jk} \sim \text{hotspot}$	$\gamma_{jk} \sim \text{MRF}$
$C \sim \text{indep}$	HRR-B	HRR-H	HRR-M
$C \sim IW$	dSUR-B	dSUR-H	dSUR-M
$C \sim HIW$	SSUR-B	SSUR-H	SSUR-M

Regressão Hierárquica Relacionada (HRR)

A Regressão Hierárquica Relacionada assume que \mathbf{C} é uma matriz diagonal, o que se traduz em independência condicional entre múltiplas variáveis resposta.

Uma *priori* gama inversa é especificada para a covariância dos resíduos, i.e

$$\sigma_k^2 \sim \mathcal{IGamma}(a_\sigma, b_\sigma)$$

Quando combinada com as *prioris* em (2), é conjugado com a verossimilhança do modelo (1). Podemos então amostrar a estrutura de seleção de variáveis $\boldsymbol{\Gamma}$ marginalmente com respeito a \mathbf{C} e \mathbf{B} .

HRR com uma *priori* Bernouli independente

Para uma *priori* simples de seleção do modelo de regressão, os indicadores binários latentes seguem uma *priori* de Bernoulli:

$$\gamma_{jk} | \omega_{jk} \sim \text{Ber}(\omega_{jk}) \quad (j = 1, \dots, p, \text{ e } k = 1, \dots, s) \quad (3)$$

Com uma *priori* hierárquica Beta em ω_j , i.e. $\omega_j \sim \text{Beta}(a_\omega, b_\omega)$, que quantifica a probabilidade de cada preditor ser associado com qualquer uma das variáveis resposta.

HRR com uma *priori* hotspot

É proposta a decomposição da probabilidade do parâmetro de associação ω_{jk} em (3), onde o_k é responsável pela esparsividade de cada modelo de resposta e π_j controla a propensão de cada preditor a ser associado a múltiplas respostas simultaneamente:

$$\gamma_{jk} | \omega_{jk} \sim \text{Ber}(\omega_{jk}) \quad (j = 1, \dots, p, \text{ e } k = 1, \dots, s) \quad (4)$$

$$\omega_{jk} = o_k \times \pi_j$$

$$o_k \sim \text{Beta}(a_0, b_0)$$

$$\pi_j \sim \text{Gamma}(a_\pi, b_\pi)$$

Regressão não relacionada aparentemente esparsa (SSUR)

Para modelar a matriz de covariância C é especificado uma *priori* hiper-Inversa Wishart, o que significa que as variáveis resposta têm por trás um grafo \mathcal{G} que codifica a dependência condicional entre as respostas.

Um grafo esparso corresponde à matriz esparsa de precisão C^{-1} . Do ponto de vista computacional, é impraticável especificar uma *priori* hiper-inversa Wishart diretamente em C^{-1} . É realizada uma transformação em C para fatorar a verossimilhança. A distribuição hiper inversa de Wishart i.e $C \sim \mathcal{HIW}_{\mathcal{G}}(\nu, \tau \mathbb{I}_s)$ transforma-se na variância escalar σ_{qt}^2 e no vetor de correlação associado $\boldsymbol{\rho}_{qt} = (\rho_{1,qt}, \rho_{2,qt}, \dots, \rho_{t-1,qt})^T$, com:

Amostragem MCMC e inferência *a posteriori*

Referências

Zhao, Zhi, Marco Banterle, Leonardo Bottolo, Sylvia Richardson, Alex Lewin, and Manuela Zucknick. 2021. “BayesSUR: An R Package for High-Dimensional Multivariate Bayesian Variable and Covariance Selection in Linear Regression.” *Journal of Statistical Software* 100 (11): 1–32. <https://doi.org/10.18637/jss.v100.i11>.

Apêndice: códigos

```
knitr::opts_chunk$set(echo = TRUE)
```