

Equações de estimação regularizadas

Vinicius Ricardo Riffel

Orientador: Prof. Dr. Wagner Hugo Bonat

Universidade Federal do Paraná

Departamento de Estatística

Laboratório de Estatística e Geoinformação

viniciusriffel@ufpr.br

Set. 14, 2022



- 1 Sobre a apresentação
- 2 Introdução
- 3 Metodologia
- 4 Resultados
- 5 Considerações e trabalho futuros
- 6 Referências

Sumário

- 1 Sobre a apresentação
- 2 Introdução
 - Do LM ao McGLM
- 3 Metodologia
 - Equações de estimação
 - Estimando um McGLM
 - Equações de estimação regularizadas
 - Estimando um McGLM regularizado
- 4 Resultados
 - Estudo de simulação
 - Aplicação em dados reais
- 5 Considerações e trabalho futuros
- 6 Referências

Sobre a apresentação

- TCC sob orientação do Prof. Dr. Wagner Hugo Bonat.
- Objetivo: propor um algoritmo de estimação baseado em equações de estimação regularizadas.
- Início em \approx dez/2021 e ainda não foi finalizado.
- Provisoriamente, a implementação está hospedada em <https://github.com/vriffel/mcglm>

Sumário

- 1 Sobre a apresentação
- 2 **Introdução**
 - Do LM ao McGLM
- 3 Metodologia
 - Equações de estimação
 - Estimando um McGLM
 - Equações de estimação regularizadas
 - Estimando um McGLM regularizado
- 4 Resultados
 - Estudo de simulação
 - Aplicação em dados reais
- 5 Considerações e trabalho futuros
- 6 Referências

Introdução

- Modelos de regressão são técnicas que permitem descrever problemas do mundo real de forma probabilística.
- Suposições são feitas na aplicação de tais técnicas.
- Há diferentes modelos disponíveis:
 - Modelo linear (LM)
 - Modelo linear generalizado (GLM)
 - Equações de estimação generalizadas (GEE)
 - Modelos multivariados de covariância linear generalizados (McGLM)

LM

- Especificação de um LM:

$$E(Y) = \mu = X\beta$$

$$Var(Y) = \tau I$$

- Suposições: normalidade, independência, variância constante.

GLM

- Especificação de um GLM:

$$E(Y) = g^{-1}(\mu) = g^{-1}(X\beta)$$

$$Var(Y) = \tau V(\mu)$$

- Suposições: família exponencial natural e independência.
- Ao utilizar quase-verossimilhança, não haverá a limitação da distribuição assumida.

GEE

- Especificação:

$$E(Y) = g^{-1}(\mu) = g^{-1}(X\beta)$$

$$Var(Y) = \tau V(\mu)^{\frac{1}{2}} \Omega(\rho) V(\mu)^{\frac{1}{2}}$$

- Limitações: dados multivariados e estrutura de covariância.
- Estimado via equações de estimação.

McGLM (Bonat e Jørgensen, 2016).

- Especificação:

$$E(Y) = (g_1^{-1}(X_1\beta_1), g_2^{-1}(X_2\beta_2), \dots, g_R^{-1}(X_R\beta_R))$$

$$Var(Y) = \Sigma_R \overset{G}{\otimes} \Sigma_b$$

$$\Sigma_R \overset{G}{\otimes} \Sigma_b = \text{Bdiag}(\tilde{\Sigma}_1, \tilde{\Sigma}_2, \dots, \tilde{\Sigma}_R) (\Sigma_b \otimes I) \text{Bdiag}(\tilde{\Sigma}_1^T, \tilde{\Sigma}_2^T, \dots, \tilde{\Sigma}_R^T)$$

$$\Sigma_r = V(\mu_r; p_r)^{\frac{1}{2}} (\Omega(\tau_r)) V(\mu_r; p_r)^{\frac{1}{2}}$$

$$h(\Omega(\tau_r)) = \tau_{r0}Z_{r0} + \tau_{r1}Z_{r1} + \dots + \tau_{rD}Z_{rD}$$

Sumário

- 1 Sobre a apresentação
- 2 Introdução
 - Do LM ao McGLM
- 3 Metodologia
 - Equações de estimação
 - Estimando um McGLM
 - Equações de estimação regularizadas
 - Estimando um McGLM regularizado
- 4 Resultados
 - Estudo de simulação
 - Aplicação em dados reais
- 5 Considerações e trabalho futuros
- 6 Referências

Equações de estimação

- As equações de estimação são métodos de estimação que se destacam por sua generalidade e propriedades.
- Diversos métodos de estimação são casos particulares:
 - mínimos quadrados;
 - método da máxima verossimilhança;
 - quase-verossimilhança.

Definição equação de estimação

Uma equação de estimação de estimação é definida como:

$$g(Y; \beta) = 0$$

Estimando um McGLM

- Os McGLM são estimados utilizando o método *modified chaser* (Jørgensen e Knudsen, 2004).

Funções de estimação de um McGLM

$$\psi_{\beta}(\beta, \lambda) = \mathbf{D}^T \mathbf{C}^{-1} (\mathcal{Y} - \mathcal{M})$$

$$\psi_{\lambda_i}(\beta, \lambda) = \text{tr} \left(W_{\lambda_i} (\mathbf{r}^T \mathbf{r}) - \mathbf{C} \right)$$

Algoritmo de estimação de um McGLM

$$\beta^{(i+1)} = \beta^{(i)} - S_{\beta}^{-1} \psi_{\beta}(\beta^{(i)}, \lambda^{(i)})$$

$$\lambda^{(i+1)} = \lambda^{(i)} - S_{\lambda}^{-1} \psi_{\lambda_i}(\beta^{(i+1)}, \lambda^{(i)})$$

Equações de estimação regularizadas

- A regularização é uma restrição imposta no processo de estimação dos modelos.
- Torna possível lidar com *high dimensional data*, diminuir o erro preditivo e fazer seleção de covariáveis.
- Um exemplo de regularização em equações de estimação é (Fu, 2003):

$$g^*(Y; \beta) = g(Y; \beta) - \gamma p'(|\beta|) = g(Y; \beta) - \gamma p'(\beta) \frac{\beta}{|\beta|}$$

- γ pode ser escolhido via validação cruzada.

Estimando um McGLM regularizado

Equações de estimação do McGLM regularizado

$$\psi_{\beta}^*(\beta, \lambda) = \mathbf{D}^T \mathbf{C}^{-1} (\mathcal{Y} - \mathcal{M}) - \gamma \odot \Gamma(p'_1(|\beta|))$$

$$\psi_{\lambda_i}^*(\beta, \lambda) = \text{tr} \left(W_{\lambda_i} (r^T r) - \mathbf{C} \right) - \gamma_i \odot \Gamma(p'_2(|\lambda_i|))$$

Algoritmo de estimação

$$\beta^{(i+1)} = \beta^{(i)} - S_{\beta}^{*-1} \psi_{\beta}^*(\beta^{(i)}, \lambda^{(i)})$$

$$\lambda^{(i+1)} = \lambda^{(i)} - S_{\lambda}^{*-1} \psi_{\lambda_i}^*(\beta^{(i+1)}, \lambda^{(i)})$$

Por que regularizar?

- Torna possível estimar um McGLM em *high dimensional data*.
- Melhora capacidade preditiva.
- Permite diminuir o número de parâmetros e construir redes de associação.
- Pode haver viés nas estimativas penalizadas.

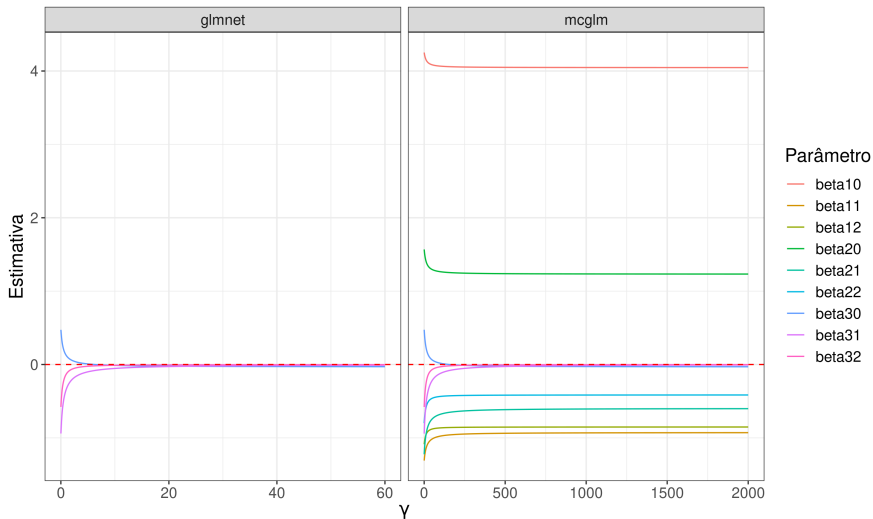
Sumário

- 1 Sobre a apresentação
- 2 Introdução
 - Do LM ao McGLM
- 3 Metodologia
 - Equações de estimação
 - Estimando um McGLM
 - Equações de estimação regularizadas
 - Estimando um McGLM regularizado
- 4 Resultados
 - Estudo de simulação
 - Aplicação em dados reais
- 5 Considerações e trabalho futuros
- 6 Referências

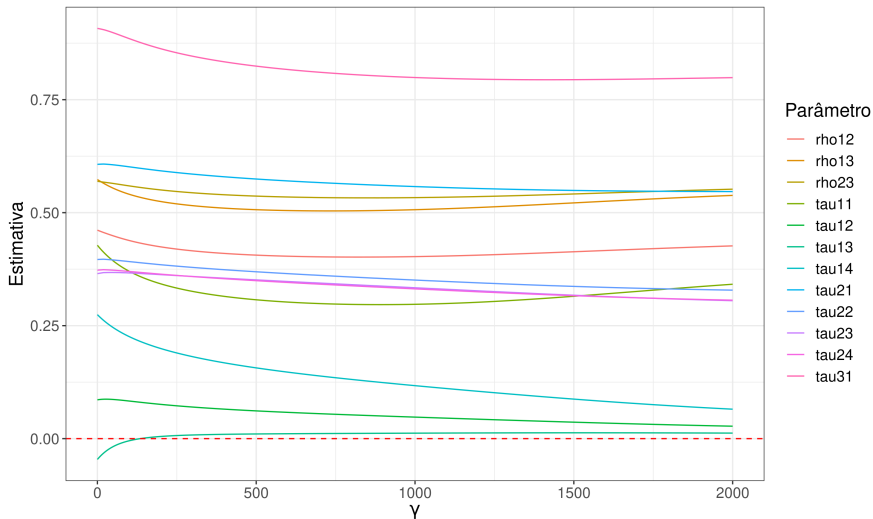
Estudo de simulação

- Conjunto de dados teca
 - Variáveis resposta: níveis de 3 elementos químicos no solo;
 - Covariável: Nível da medição no solo;
 - Disponível no pacote EACS (Zeviani, 2019).
- Objetivo: verificar se as estimativas penalizadas se comportam de maneira esperada.
- Penalização *Ridge*.
- Estrutura de covariância parecida com a não estruturada, mas com os interceptos.
- Estrutura de média linear.

Conjunto de dados teca



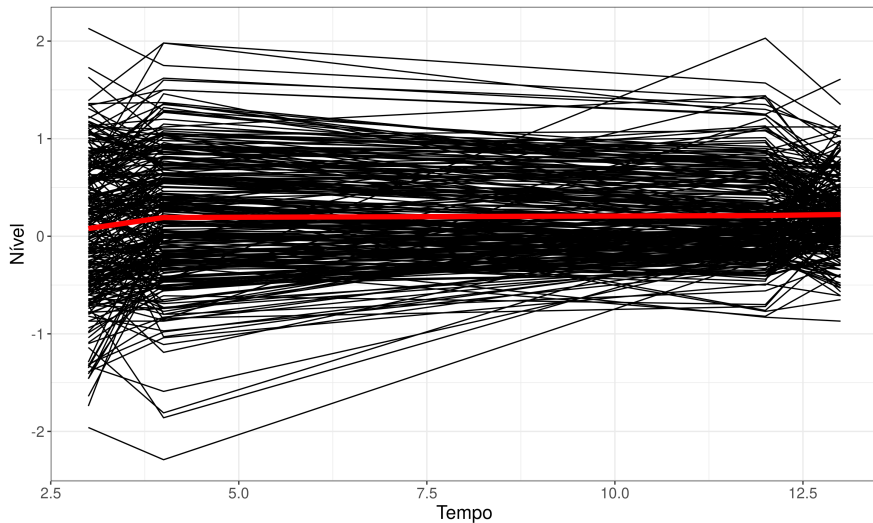
Conjunto de dados teca



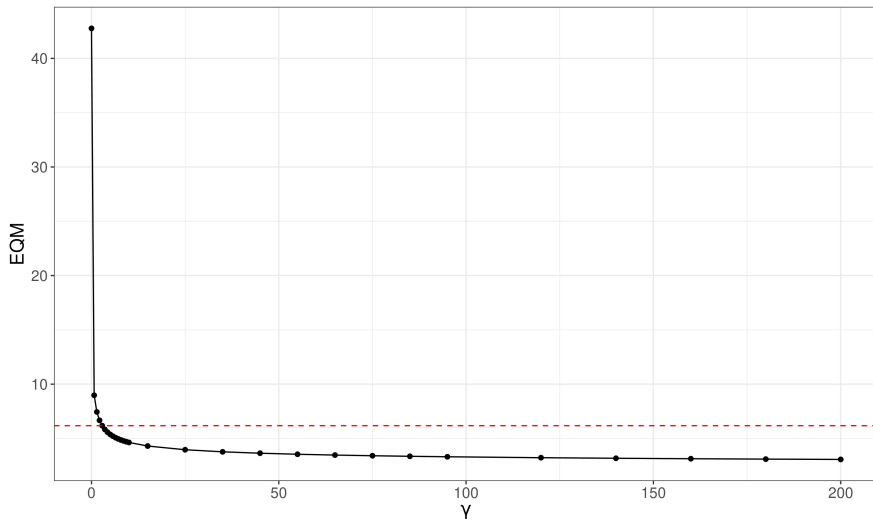
Conjunto de dados yeastG1

- Conjunto de dados yeastG1 (Spellman et al, 1998):
 - Variável resposta: nível de expressão genética em 4 tempo distintos;
 - 96 covariáveis;
 - 263 indivíduos;
 - Disponível no pacote PGEE (Inan et al, 2017).
- Objetivo: regularizar a estrutura de média para melhorar o erro preditivo.
- Penalização *Ridge*.
- Estrutura de covariância parecida com a não estruturada, mas com os interceptos.
- Estrutura de média linear.
- Seleção de γ via validação cruzada com 10 *folds*, utilizamos EQM.
- A estrutura de dependência foi preservada removendo os indivíduos para validação.

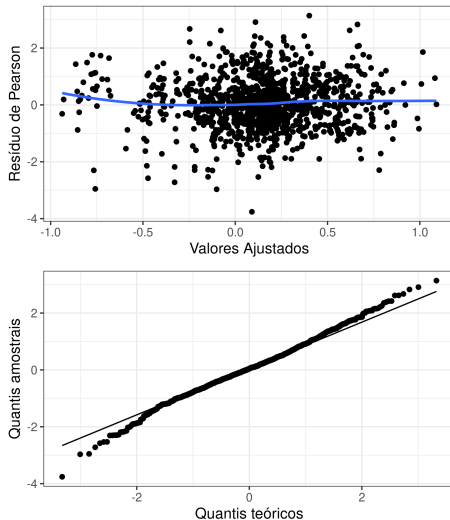
Conjunto de dados yeastG1



Conjunto de dados yeastG1



Conjunto de dados yeastG1



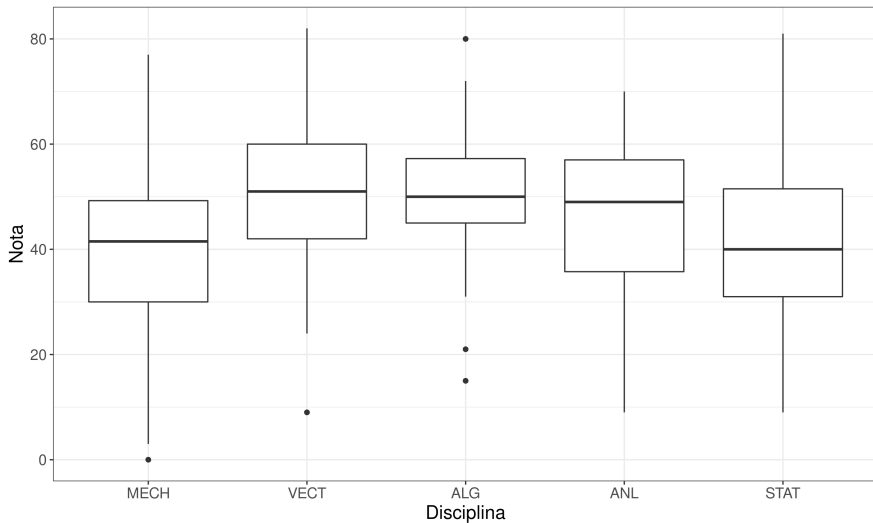
Conjunto de dados marks

- Conjunto de dados marks (Mardia et al, 1979):
 - Variável resposta: nota de 88 alunos em 6 disciplinas (Álgebra, Análise, Estatística, Mecânica, Vetores);
 - Disponível no pacote `bnlearn` (Scutari, 2010).
- Objetivo: construir uma rede de associações entre as disciplinas.
- Penalização *Ridge*.
- Estrutura de covariância parecida com a não estruturada, mas com os interceptos.
- Modelamos a precisão (função de ligação inversa).
- Estrutura de média linear.

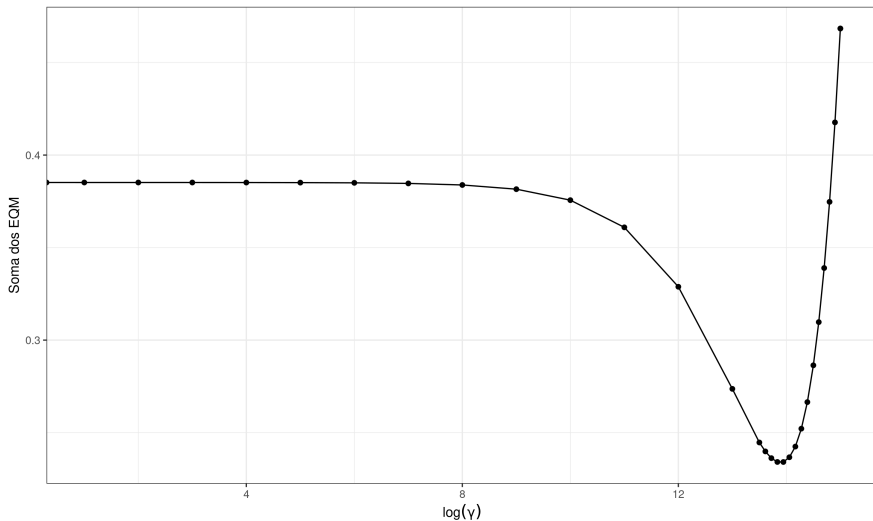
Conjunto de dados marks

- Seleção de γ via validação cruzada com uma fold.
- EQM do valor predito não foi possível.
- Utilizamos a soma dos EQMs entre a matriz de correlação da base de teste e as matrizes estimadas.

Conjunto de dados marks



Conjunto de dados marks

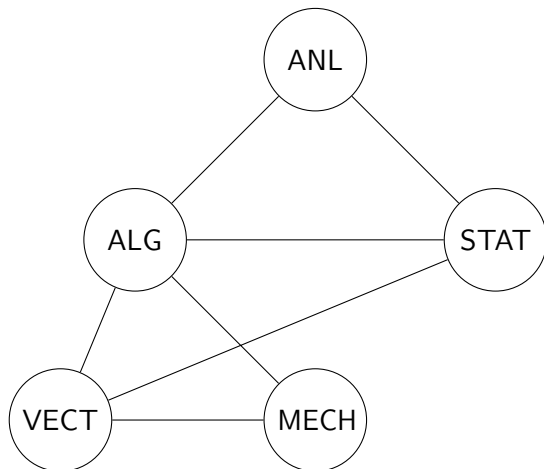


Conjunto de dados marks

- A conclusão pela existência de associação significativa é feita pelas estimativas dos coeficientes de dispersão.

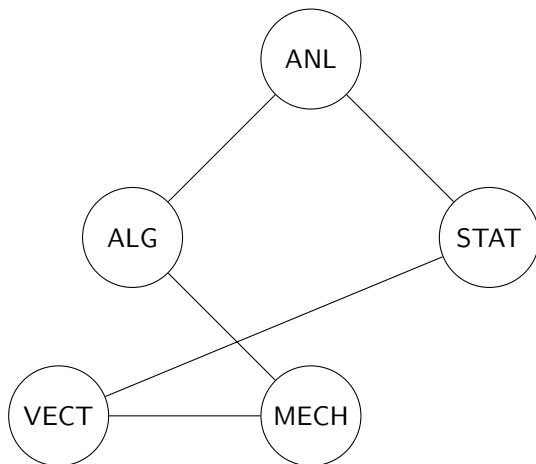
Coeficiente	Estimativa	Erro padrão	Estatística Z	Relação
$\hat{\tau}_{10}$	0.00644	0.00032	20.28	
$\hat{\tau}_{11}$	-0.00215	0.00040	-5.34	MECH-VECT
$\hat{\tau}_{12}$	-0.00124	0.00041	-3.03	MECH-ALG
$\hat{\tau}_{13}$	-0.00027	0.00041	-0.66	MECH-ANL
$\hat{\tau}_{14}$	-0.00055	0.00042	-1.33	MECH-STAT
$\hat{\tau}_{15}$	-0.00111	0.00041	-2.73	VECT-ALG
$\hat{\tau}_{16}$	-0.00068	0.00041	-1.67	VECT-ANL
$\hat{\tau}_{17}$	-0.00123	0.00041	-2.99	VECT-STAT
$\hat{\tau}_{18}$	-0.00201	0.00040	-5.02	ALG-ANL
$\hat{\tau}_{19}$	-0.00097	0.00041	-2.35	ALG-STAT
$\hat{\tau}_{110}$	-0.00164	0.00041	-4.01	STAT-ANL

Conjunto de dados marks



Conjunto de dados marks

- Aplicando a correção de Bonferroni a rede fica:



Sumário

- 1 Sobre a apresentação
- 2 Introdução
 - Do LM ao McGLM
- 3 Metodologia
 - Equações de estimação
 - Estimando um McGLM
 - Equações de estimação regularizadas
 - Estimando um McGLM regularizado
- 4 Resultados
 - Estudo de simulação
 - Aplicação em dados reais
- 5 Considerações e trabalho futuros
- 6 Referências

Considerações e trabalho futuros

- Foi apresentado um modelo baseado em equações de estimação regularizadas para lidar com os mais variados conjuntos de dados.
- Consistência dos estimadores.
- Outros tipos de penalização.
- Escolha dos valores de γ .

Sumário

- 1 Sobre a apresentação
- 2 Introdução
 - Do LM ao McGLM
- 3 Metodologia
 - Equações de estimação
 - Estimando um McGLM
 - Equações de estimação regularizadas
 - Estimando um McGLM regularizado
- 4 Resultados
 - Estudo de simulação
 - Aplicação em dados reais
- 5 Considerações e trabalho futuros
- 6 Referências

Referências

BONAT, W. H.; JØRGENSEN, B. Multivariate covariance generalized linear models. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, v. 65, n. 5, p. 649–675, 2016.

FU, W. J. Penalized estimating equations. *Biometrics*, v. 59, n. 1, p. 126–132, 2003.

JØRGENSEN, B.; KNUDSEN, S. J. Parameter orthogonality and bias adjustment for estimating functions. *Scandinavian Journal of Statistics*, v. 31, n. 1, p. 93–114, 2004.

INAN, G.; ZHOU, J.; WANG, L. PGEE: Penalized Generalized Estimating Equations in High-Dimension. [S.l.], 2017. R package version 1.5.

MARDIA, K.; KENT, J.; BIBBY, J. Multivariate analysis. London [u.a.]: Acad. Press, 1979. (Probability and mathematical statistics). ISBN 0124712509.

SCUTARI, M. Learning bayesian networks with the bnlearn R package. *Journal of Statistical Software*, v. 35, n. 3, p. 1–22, 2010.

SPELLMAN, P. T. et al. Comprehensive identification of cell cycle-regulated genes of the yeast *saccharomyces cerevisiae* by microarray hybridization. *Molecular Biology of the Cell*, v. 9, n. 12, p. 3273–3297, 1998. PMID: 9843569.

ZEVIANI, W. M. EACS: Estatística Aplicada à Ciência do Solo. [S.l.], 2019.